



ユーザーガイド

AWS Glue



AWS Glue: ユーザーガイド

Copyright © 2024 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Amazon の商標とトレードドレスは、Amazon 以外の製品またはサービスとの関連において、顧客に混乱を招いたり、Amazon の名誉または信用を毀損するような方法で使用することはできません。Amazon が所有していない他のすべての商標は、それぞれの所有者の所有物であり、Amazon と提携、接続、または後援されている場合とされていない場合があります。

Table of Contents

AWS Glue とは	1
AWS Glue の機能	2
AWS Glue のイノベーションについて学ぶ	3
AWS Glue の開始方法	3
AWS Glue へのアクセス	4
関連サービス	4
使用方法	5
独立で実行されるサーバーレス ETL ジョブ	6
概念	7
AWS Glue の用語	9
コンポーネント	12
AWS Glue コンソール	12
AWS Glue Data Catalog	13
AWS Glue クローラおよび分類子	14
AWS Glue ETL オペレーション	14
AWS Glue でのストリーミング ETL	14
AWS Glue ジョブシステム	15
ビジュアル ETL コンポーネント	15
AWS Glue for Spark と AWS Glue for Ray	21
AWS Glue for Ray とは	22
半構造化されたスキーマをリレーショナルスキーマに変換する	23
AWS Glue の種類	25
AWS Glue データカタログタイプ	25
AWS Glue with Spark スクリプトのタイプ	25
AWS Glue ークローラタイプ	26
使用開始方法	27
AWS Glue クローン作成の概要	27
IAM アクセス許可のセットアップ	29
次のステップ	34
ビジュアル ETL を使用するための IAM アクセス許可	34
AWS Glue Studio 中でのノートブックの使用開始	45
使用状況プロファイルの設定	48
使用状況プロファイルの管理	49
使用状況プロファイルとジョブ	62

AWS Glue Data Catalog の開始方法	63
概要	63
ステップ 1: データベースを作成する	63
ステップ 2. テーブルを作成する	65
次のステップ	66
データストアへのネットワークアクセスを設定する	70
AWS Glue の PyPI に接続する VPC のセットアップ	71
VPC での DNS のセットアップ	73
暗号化のセットアップ	73
開発のためのネットワーク設定	78
開発エンドポイント用にネットワークを設定する	78
ノートブックサーバー用の Amazon EC2 のセットアップ	80
データ検出とカタログ化	82
データカタログの入力	84
AWS Glue クローラー の使用	85
メタデータの手動定義	187
他の AWS のサービスとの統合	206
データカタログの設定	207
トランザクションテーブルへのデータ入力と管理	210
Iceberg テーブルの作成	210
Iceberg テーブルの最適化	214
データカタログの管理	227
スキーマを更新し、新規パーティションを追加する	228
列統計を使用したクエリのパフォーマンスの最適化	235
Data Catalog の暗号化	247
Lake Formation を使用したデータカタログの保護	248
データカタログにアクセスする	248
データカタログのベストプラクティス	249
AWS Glue スキーマレジストリ	250
スキーマ	251
レジストリ	253
スキーマのバージョニングと互換性	254
オープンソース Serde ライブラリ	260
Schema Registry のクォータ	260
仕組み	261
開始	263

AWS Glue Schema Registry との統合	286
AWS Glue スキーマレジストリへの移行	313
データへの接続	316
AWS Glue 接続プロパティ	317
必要な接続プロパティ	318
JDBC 接続プロパティ	319
MongoDB と MongoDB Atlas の接続プロパティ	324
Salesforce 接続プロパティ	324
Snowflake 接続	325
Vertica 接続	326
SAP HANA 接続	327
Azure SQL 接続	328
Teradata Vantage 接続	329
OpenSearch サービス接続	330
Azure Cosmos 接続	331
SSL 接続プロパティ	331
認証用の Kafka 接続プロパティ	334
Google BigQuery 接続	335
Vertica 接続	326
AWS Secrets Manager への接続認証情報の保存	335
AWS Glue 接続の追加	336
Redshift に接続する	337
Azure Cosmos DB に対する接続	341
Azure SQL に対する接続	345
への接続 BigQuery	348
MongoDB に対する接続	353
OpenSearch Service に対する接続	357
Salesforce への接続	360
SAP HANA に対する接続	372
Snowflake に接続する	376
Teradata に対する接続	380
Vertica に対する接続	383
コネクタと接続の使用	387
データソースへの接続	416
独自の JDBC ドライバーを使用した JDBC 接続の追加	425
AWS Glue 接続のテスト	430

すべての AWS コールを VPC を経由するように設定する	430
VPC の JDBC データストアに接続する	431
Elastic Network Interface を使用して VPC データにアクセスする	432
Elastic Network Interface プロパティ	433
MongoDB または MongoDB Atlas 接続を使用する	434
VPC エンドポイントを使用して Amazon S3 データストアをクローलする	434
前提条件	435
Amazon S3 への接続を作成する	436
Amazon S3 への接続をテストする	439
Amazon S3 データストアのクローラーを作成する	441
Amazon S3 backed データカタログテーブルのクローラーを作成する	443
クローラーの実行	444
トラブルシューティング	444
接続の問題のトラブルシューティング	444
チュートリアル: Elasticsearch 向けの AWS Glue コネクタを使用する	445
前提条件	446
ステップ 1: (オプション) OpenSearch クラスター情報の AWS シークレットを作成する	446
ステップ 2: コネクタをサブスクライブする	447
ステップ 3: AWS Glue Studio でコネクタをアクティブにして、接続を作成します。	448
ステップ 4: ETL ジョブの IAM ロールを設定する	449
ステップ 5: OpenSearch 接続を使用するジョブを作成する	450
ステップ 6: ジョブを実行する	451
インタラクティブセッションで AWS Glue ジョブを構築する	452
AWS Glue インタラクティブセッションの概要	452
制限事項	453
AWS Glue のインタラクティブセッションの開始方法	453
ローカルでインタラクティブセッションを設定するための前提条件	453
Jupyter AWS Glue とインタラクティブセッション Jupyter カーネルのインストール	453
Jupyter の実行	454
セッション認証情報とリージョンの設定	454
インタラクティブセッションプレビューからのアップグレード	456
SageMaker Studio でインタラクティブセッションを使用する	456
Microsoft Visual Studio Code によるインタラクティブセッションの使用	457
AWS Glue Jupyter および AWS Glue Studio ノートブックのインタラクティブセッションの設 定	460
Jupyter Magics の概要	460

Jupyter 向け AWS Glue インタラクティブセッションでサポートされているマジック	460
ネーミングセッション	477
インタラクティブセッションの IAM ロールの指定	478
名前付きプロファイルを使用したセッションの設定	478
AWS Glue for Ray インタラクティブセッション (プレビュー)	480
AWS Glue Studio コンソールでの Ray インタラクティブセッション	480
Jupyter Kernel を使用した Ray のインタラクティブセッション	481
Ray インタラクティブセッションのタイムアウトのデフォルト	481
AWS Glue の Ray インタラクティブセッションでサポートされているマジック	482
IAM を使用したインタラクティブセッション	483
インタラクティブセッションで使用する IAM プリンシパル	484
クライアントプリンシパルを設定する	484
ランタイムロールを設定する	485
以下の方法でセッションを非公開にします。 TagOnCreate	486
IAM ポリシーに関する考慮事項	491
スクリプトまたはノートブックの AWS Glue ジョブジョブへの変換	491
AWS Glue ストリーミングのためのインタラクティブセッション	492
ストリーミングセッションタイプの切り替え	492
インタラクティブな開発のためのサンプリング入力ストリーム	492
インタラクティブセッションでのストリーミングアプリケーションの実行	494
ローカルでの開発とテスト	495
を使用した開発AWS Glue Studio	496
インタラクティブセッションによる開発	496
Docker イメージによる開発	496
AWS Glue ETL ライブラリを使用した開発	508
開発者エンドポイント	516
開発エンドポイントからインタラクティブセッションへ移行する	518
開発エンドポイントを使用してスクリプトを開発する	520
ノートブックの管理	549
AWS Glue Studio でビジュアル ETL ジョブを作成する	551
コンソールにサインインする	551
AWS Glue Studio でジョブを作成するための次のステップ	552
ビジュアル ETL と AWS Glue Studio	552
AWS Glue Studio でのジョブの開始	552
ジョブエディタの機能	554
AWS Glue マネージドデータ変換ノードの編集	562

カスタムビジュアル変換	625
データレイクフレームワークを AWS Glue Studio で使用する	642
データターゲットノードの設定	654
ジョブスクリプトの編集またはアップロード	659
ジョブ図でノードの親ノードを変更する	664
ジョブ図からノードを削除する	664
AWS Glue データカタログノードにソースパラメータとターゲットパラメータを追加する	669
AWS Glue で Git のバージョン管理システムを使用する	671
AWS Glue Studio ノートブックによるコードの作成	680
ノートブックの使用の概要	680
AWS Glue Studio 中のノートブックを使用した ETL ジョブの作成	681
ノートブックエディタコンポーネント	682
ノートブックとジョブスクリプトの保存	683
ノートブックセッションの管理	684
CodeWhisperer との併用 AWS Glue Studio notebooks	686
ジョブ実行を表示する	686
ジョブモニタリングダッシュボードにアクセスする	686
ジョブモニタリングダッシュボードの概要	687
ジョブの実行ビュー	687
ジョブの実行ログの表示	691
ジョブの実行の詳細を表示する	691
Spark ジョブ実行の Amazon CloudWatch メトリクスを表示	694
Ray ジョブ実行の Amazon CloudWatch メトリクスを表示	695
機密データを検出して処理する	697
データのスキャン方法の選択	697
検出する PII エンティティの選択	699
検出感度のレベルの指定	703
特定された PII データによる対処方法の選択	704
詳細なアクションオーバーライドの追加	705
ジョブの管理	706
ジョブの実行の開始	706
ジョブの実行のスケジュール	707
ジョブスケジュールの管理	708
ジョブの実行の停止	709
ジョブの表示	709

最近のジョブの実行の情報を表示する	710
ジョブスクリプトの表示	711
ジョブのプロパティを変更する	712
ジョブの保存	715
ジョブのクローンを作成する	717
ジョブの削除	717
ジョブの使用	719
AWS Glue バージョン	719
AWS Glue バージョン	719
起動時間を短縮した Spark ETL ジョブの実行	734
Spark ジョブの AWS Glue の AWS Glue バージョン 3.0 への移行	739
Spark ジョブの AWS Glue の AWS Glue バージョン 4.0 への移行	748
AWS Glue for Ray (プレビュー) から AWS Glue for Ray への移行	763
AWS Glue バージョンサポートポリシー	764
Spark ジョブの使用	766
ジョブのパラメータ	766
スパークとジョブ PySpark	776
ストリーミング ETLジョブ	913
FindMatches によるレコードのマッチング	928
Spark プログラムの移行	963
Ray ジョブの使用	970
AWS Glue for Ray の使用を開始する	971
サポートされている Ray のランタイム環境	972
Ray ジョブのワーカーの考慮	973
Ray ジョブのパラメータ	973
Ray ジョブメトリクス	976
Python シェルジョブのプロパティの設定	978
制限事項	978
Python シェルジョブのジョブプロパティの定義	978
サポートされている Python シェルジョブのライブラリ	980
独自の Python ライブラリの提供	982
AWS Glue の Python シェルジョブで AWS CloudFormation を使用する	986
モニタリング	986
AWSタグ	988
CloudWatch Events を使用した自動化	993
AWS Glue リソースのモニタリング	995

CloudTrail によるログ記録	997
ジョブ実行ステータス	1001
AWS Glue ストリーミング	1004
ストリーミングのユースケース	1004
AWS Glue ストリーミングを使用する利点は何ですか？	1005
AWS Glue ストリーミングを使用するタイミング	1006
サポートされているデータソース	1007
サポートされるデータターゲット	1007
チュートリアル: AWS Glue Studio を使用して最初のストリーミングワークロードを構築する	1008
前提条件	1008
Amazon Kinesis からストリーミングデータを消費する	1008
チュートリアル: AWS Glue Studio ノートブックを使用して最初のストリーミングワークロードを構築する	1019
前提条件	1020
Amazon Kinesis からストリーミングデータを消費する	1020
Streaming の概念	1027
AWS Glue Streaming ジョブの構造	1028
Kafka 接続	1031
Kinesis 接続	1038
Streaming オプション	1045
AWS Glue Streaming 自動スケーリング	1046
AWS Glue Studio の Auto-Scaling を有効にする	1046
AWS CLI または SDK を使用した Auto Scaling の有効化	1047
仕組み	1048
メンテナンスウィンドウ	1050
メンテナンスウィンドウの設定	1050
メンテナンスウィンドウの動作	1051
ジョブのモニタリング	1052
データ損失処理	1054
高度な AWS Glue のストリーミングの概念	1055
ストリームを処理する際の時間に関する考慮事項	1055
ウィンドウ処理	1056
遅延データとウォーターマークの処理	1061
AWS Glue Streaming ジョブのモニタリング	1063
メトリクスの視覚化	1064

メトリクスの詳細	1065
最高のパフォーマンスを得る方法	1070
AWS Glue データ品質	1072
利点と主な特徴	1072
仕組み	1073
のデータ品質 AWS Glue Data Catalog	1073
AWS Glue ETL ジョブのデータ品質	1073
AWS Glue Data Quality エントリポイントの比較	1074
考慮事項	1076
用語	1076
制限	1077
AWS Glue Data Quality のリリースノート	1077
一般提供を開始: 新機能	1077
2023 年 11 月 27 日 (プレビュー)	1078
2024 年 3 月 12 日	1078
2024 年 6 月 26 日	1078
AWS Glue Data Quality での異常検出	1078
仕組み	1079
アナライザーを使用してデータを検査する	1080
DetectAnomaly ルールの使用	1081
異常検出の利点とユースケース	1081
AWS Glue Data Quality の IAM アクセス許可	1082
IAM アクセス許可	1082
評価実行のスケジューリングに必要な IAM の設定	1084
IAM ポリシーの例	1086
Data Catalog で AWS Glue Data Quality を使用する	1089
前提条件	1090
手順の例	1090
ルールの推奨事項の生成	1090
ルールの推奨のモニタリング	1092
推奨ルールセットの編集	1092
新しいルールセットの作成	1094
ルールセットを実行してデータ品質を評価する	1095
データ品質スコアと結果を表示する	1097
関連トピック	1097
AWS Glue Studio を使用したデータ品質の評価	1098

利点	1098
AWS Glue Studio で ETL ジョブのデータ品質を評価する	1098
Data Quality ルールビルダー	1104
異常検出の設定とインサイトの生成	1109
AWS Glue Studio ノートブックの ETL ジョブのデータ品質	1113
前提条件	1114
AWS Glue Studio での ETL ジョブの作成	1114
データ品質定義言語 (DQDL) リファレンス	1119
構文	1120
ルールタイプリファレンス	1134
API を使用したデータ品質の測定および管理	1180
前提条件	1181
AWS Glue Data Quality 推奨事項の操作	1181
AWS Glue Data Quality ルールセットの操作	1184
AWS Glue Data Quality 実行の操作	1186
AWS Glue Data Quality 評価結果の操作	1191
アラート、デプロイ、スケジュールの設定	1192
Amazon EventBridge 統合でのアラートと通知の設定	1192
CloudWatch 統合でアラートと通知を設定する	1200
データ品質評価の結果をクエリする	1202
データ品質ルールをデプロイする	1206
データ品質ルールのスケジューリング	1206
AWS Glue Data Quality エラーのトラブルシューティング	1206
エラー: モジュールが見つからない	1207
エラー: アクセス許可が不十分である	1207
エラー: ルールセットが一意ではない	1208
エラー: テーブルに特殊文字が含まれている	1208
エラー: 大規模なルールセットによるオーバーフロー	1208
エラー: ルールのステータスが失敗である	1208
AnalysisException: デフォルトデータベースの存在を確認できません	1208
指定されたキーマップは所定のデータフレームに適していません	1209
java.lang.RuntimeException : データの取得に失敗しました。	1209
LAUNCH ERROR: Error downloading from S3 for bucket	1210
InvalidInputException (ステータス: 400): DataQuality ルールを解析できません	1210
エラー: EventBridge が、設定したスケジュールに基づいて Glue DQ ジョブをトリガーしない	1211

CustomSQL エラー	1211
動的ルール	1212
ユーザークラス: org.apache.spark.sql.AnalysisException: org.apache.hadoop.hive.q1.metadata の例外。HiveException	1214
UNCLASSIFIED_ERROR; IllegalArgumentException: 解析エラー: ルールまたはアナライ ザーが指定されていません。、入力時に実行可能な代替手段がありません	1214
Amazon Q でのデータインテグレーション AWS Glue	1215
Amazon Q とは	1215
Amazon Q でのデータインテグレーション AWS Glue	1215
Amazon Q データ統合の操作	1216
ベストプラクティス	1218
サービス向上	1218
考慮事項	1219
Amazon Q データ統合の設定	1219
IAM 許可の設定	1219
サポート対象のコード生成	1221
インタラクションの例	1223
Amazon Q チャットインタラクション	1223
AWS Glue Studio ノートブックのインタラクション	1224
オーケストレーション	1228
トリガーを使用したジョブとクローラの開始	1228
AWS Glue トリガー	1228
トリガーの追加	1231
トリガーの有効化と無効化	1235
設計図とワークフローを使用した複雑な ETL アクティビティの実行	1236
ワークフローの概要	1236
ワークフローの手動による作成と構築	1240
EventBridge イベントを使用するワークフローの開始	1245
ワークフローを開始した EventBridge イベントの表示	1252
ワークフローの実行およびモニタリング	1253
ワークフロー実行の停止	1256
ワークフロー実行の修復と再開	1257
ワークフローの実行プロパティの取得と設定	1263
AWS Glue APIを使用したワークフローのクエリ	1264
設計図とワークフローの制限	1268
設計図エラーのトラブルシューティング	1269

設計図のペルソナとロールのアクセス許可	1274
設計図の開発	1279
設計図の概要	1279
設計図の開発	1283
設計図の登録	1308
設計図の表示	1311
設計図の更新	1312
設計図からのワークフローの作成	1315
設計図の実行の表示	1317
AWS Glue の場合は AWS CloudFormation	1319
サンプルデータベース	1321
サンプルのデータベース、テーブル、パーティション	1322
サンプルの Grok 分類子	1326
サンプルの JSON 分類子	1327
サンプルの XML 分類子	1328
Amazon S3 クローラーのサンプル	1329
サンプルの接続	1332
サンプルの JDBC クローラー	1333
Amazon S3 から Amazon S3 へのサンプルジョブ	1336
Amazon S3 に書き込む JDBC のサンプルジョブ	1337
サンプルのオンデマンドトリガー	1339
サンプルのスケジュールされたトリガー	1340
サンプルの条件付きトリガー	1341
機械学習変換のサンプル	1343
サンプルのデータ品質ルールセット	1344
EventBridge スケジューラを使用するサンプルのデータ品質ルールセット	1345
サンプルの開発エンドポイント	1347
AWS Glue プログラミングガイド	1350
独自のカスタムスクリプトの提供	1350
AWS Glue for Spark	1351
チュートリアル: Spark スクリプトの記述	1352
の ETL PySpark	1365
Scala での ETL	1593
機能と最適化	1678
AWS Glue for Ray	1928
チュートリアル: Ray スクリプトを記述する	1928

AWS Glue for Ray での Ray Core と Ray Data の使用	1934
ファイルと Python ライブラリ	1936
データへの接続	1941
AWS SDKs	1943
AWS Glue API	1945
セキュリティ	1967
– データ型 –	1968
DataCatalogEncryptionSettings	1968
EncryptionAtRest	1968
ConnectionPasswordEncryption	1969
EncryptionConfiguration	1970
S3Encryption	1970
CloudWatchEncryption	1970
JobBookmarksEncryption	1971
SecurityConfiguration	1971
GluePolicy	1972
— 操作 —	1972
GetDataCatalogEncryptionSettings (get_data_catalog_encryption_settings)	1973
PutDataCatalogEncryptionSettings (put_data_catalog_encryption_settings)	1973
PutResourcePolicy (put_resource_policy)	1974
GetResourcePolicy (get_resource_policy)	1975
DeleteResourcePolicy (delete_resource_policy)	1976
CreateSecurityConfiguration (create_security_configuration)	1977
DeleteSecurityConfiguration (delete_security_configuration)	1978
GetSecurityConfiguration (get_security_configuration)	1978
GetSecurityConfigurations (get_security_configurations)	1979
GetResourcePolicies (get_resource_policies)	1980
カタログ	1981
データベース	1981
テーブル	1990
パーティション	2027
接続	2051
ユーザー定義関数	2068
Athena カタログをインポートする	2075
テーブルオブティマイザー	2077
— データ型 —	2077

TableOptimizer	2077
TableOptimizerConfiguration	2077
TableOptimizerRun	2078
RunMetrics	2078
BatchGetTableOptimizerEntry	2079
BatchTableOptimizer	2079
BatchGetTableOptimizerError	2080
— 操作 —	2081
GetTableOptimizer (get_table_optimizer)	2081
BatchGetTableOptimizer (batch_get_table_optimizer)	2082
ListTableOptimizerRuns (list_table_optimizer_runs)	2083
CreateTableOptimizer (create_table_optimizer)	2084
DeleteTableOptimizer (delete_table_optimizer)	2085
UpdateTableOptimizer (update_table_optimizer)	2086
クローラーおよび分類子	2087
分類子	2087
クローラー	2101
列統計	2128
スケジューラー	2135
ETL スクリプトの自動生成	2137
– データ型 –	2138
CodeGenNode	2138
CodeGenNodeArg	2138
CodeGenEdge	2139
口ケーション	2139
CatalogEntry	2140
MappingEntry	2140
— 操作 —	2141
CreateScript (create_script)	2141
GetDataflowGraph (get_dataflow_graph)	2142
GetMapping (get_mapping)	2142
GetPlan (get_plan)	2143
ビジュアルジョブ API	2144
— データ型 —	2145
CodeGenConfigurationNode	2148
JDBCConectorOptions	2154

StreamingDataPreviewOptions	2156
AthenaConnectorSource	2156
JDBCConectorSource	2157
SparkConnectorSource	2158
CatalogSource	2159
MySQLCatalogSource	2159
PostgreSQLCatalogSource	2159
OracleSQLCatalogSource	2160
MicrosoftSQLServerCatalogSource	2160
CatalogKinesisSource	2161
DirectKinesisSource	2161
KinesisStreamingSourceOptions	2162
CatalogKafkaSource	2165
DirectKafkaSource	2165
KafkaStreamingSourceOptions	2166
RedshiftSource	2169
AmazonRedshiftSource	2169
AmazonRedshiftNodeData	2169
AmazonRedshiftAdvancedOption	2172
オプション	2172
S3CatalogSource	2173
S3SourceAdditionalOptions	2173
S3CsvSource	2173
DirectJDBCSource	2176
S3DirectSourceAdditionalOptions	2176
S3JsonSource	2177
S3ParquetSource	2179
S3DeltaSource	2180
S3CatalogDeltaSource	2181
CatalogDeltaSource	2181
S3HudiSource	2182
S3CatalogHudiSource	2183
CatalogHudiSource	2183
DynamoDBCatalogSource	2184
RelationalCatalogSource	2185
JDBCConectorTarget	2185

SparkConnectorTarget	2186
BasicCatalogTarget	2187
MySQLCatalogTarget	2187
PostgreSQLCatalogTarget	2188
OracleSQLCatalogTarget	2188
MicrosoftSQLServerCatalogTarget	2189
RedshiftTarget	2189
AmazonRedshiftTarget	2190
UpsertRedshiftTargetOptions	2190
S3CatalogTarget	2191
S3GlueParquetTarget	2192
CatalogSchemaChangePolicy	2192
S3DirectTarget	2193
S3HudiCatalogTarget	2193
S3HudiDirectTarget	2194
S3DeltaCatalogTarget	2195
S3DeltaDirectTarget	2196
DirectSchemaChangePolicy	2197
ApplyMapping	2197
Mapping	2198
SelectFields	2199
DropFields	2199
RenameField	2200
スピゴット	2200
Join	2201
JoinColumn	2202
SplitFields	2202
SelectFromCollection	2202
FillMissingValues	2203
フィルター	2203
FilterExpression	2204
FilterValue	2204
CustomCode	2205
SparkSQL	2205
SqlAlias	2206
DropNullFields	2206

NullCheckBoxList	2207
NullValueField	2207
Datatype	2208
マージ	2208
Union	2209
PIIDetection	2209
集計	2210
DropDuplicates	2211
GovernedCatalogTarget	2211
GovernedCatalogSource	2212
AggregateOperation	2212
GlueSchema	2213
GlueStudioSchemaColumn	2213
GlueStudioColumn	2213
DynamicTransform	2214
TransformConfigParameter	2215
EvaluateDataQuality	2216
DQResultsPublishingOptions	2216
DQStopJobOnFailureOptions	2217
EvaluateDataQualityMultiFrame	2217
レシピ	2218
RecipeReference	2219
SnowflakeNodeData	2219
SnowflakeSource	2221
SnowflakeTarget	2222
ConnectorDataSource	2222
ConnectorDataTarget	2223
ジョブ	2224
ジョブ	2224
ジョブ実行	2249
トリガー	2267
インタラクティブセッション	2281
— データ型 —	2281
セッション	2281
SessionCommand	2283
Statement	2284

StatementOutput	2284
StatementOutputData	2285
ConnectionsList	2285
— 操作 —	2285
CreateSession (create_session)	2286
StopSession (stop_session)	2289
DeleteSession (delete_session)	2290
GetSession (get_session)	2291
ListSessions (list_sessions)	2292
RunStatement (run_statement)	2293
CancelStatement (cancel_statement)	2294
GetStatement (get_statement)	2294
ListStatements (list_statements)	2295
DevEndpoints	2296
— データ型 —	2296
DevEndpoint	2296
DevEndpointCustomLibraries	2300
— 操作 —	2301
CreateDevEndpoint (create_dev_endpoint)	2301
UpdateDevEndpoint (update_dev_endpoint)	2307
DeleteDevEndpoint (delete_dev_endpoint)	2308
GetDevEndpoint (get_dev_endpoint)	2309
GetDevEndpoints (get_dev_endpoints)	2309
BatchGetDevEndpoints (batch_get_dev_endpoints)	2310
ListDevEndpoints (list_dev_endpoints)	2311
スキーマレジストリ	2312
— データ型 —	2313
RegistryId	2313
RegistryListItem	2313
MetadataInfo	2314
OtherMetadataValueListItem	2314
SchemaListItem	2315
SchemaVersionListItem	2316
MetadataKeyValuePair	2316
SchemaVersionErrorItem	2316
ErrorDetails	2317

SchemaVersionNumber	2317
Schemald	2317
— 操作 —	2318
CreateRegistry (create_registry)	2319
CreateSchema (create_schema)	2320
GetSchema (get_schema)	2324
ListSchemaVersions (list_schema_versions)	2326
GetSchemaVersion (get_schema_version)	2327
GetSchemaVersionsDiff (get_schema_versions_diff)	2328
ListRegistries (list_registries)	2329
ListSchemas (list_schemas)	2330
RegisterSchemaVersion (register_schema_version)	2331
UpdateSchema (update_schema)	2332
CheckSchemaVersionValidity (check_schema_version_validity)	2334
UpdateRegistry (update_registry)	2335
GetSchemaByDefinition (get_schema_by_definition)	2336
GetRegistry (get_registry)	2337
PutSchemaVersionMetadata (put_schema_version_metadata)	2338
QuerySchemaVersionMetadata (query_schema_version_metadata)	2339
RemoveSchemaVersionMetadata (remove_schema_version_metadata)	2341
DeleteRegistry (delete_registry)	2342
DeleteSchema (delete_schema)	2343
DeleteSchemaVersions (delete_schema_versions)	2344
ワークフロー	2345
– データ型 –	2345
JobNodeDetails	2346
CrawlerNodeDetails	2346
TriggerNodeDetails	2346
Crawl	2347
ノード	2347
Edge	2348
ワークフロー	2348
WorkflowGraph	2350
WorkflowRun	2350
WorkflowRunStatistics	2351
StartingEventBatchCondition	2352

Blueprint	2352
BlueprintDetails	2353
LastActiveDefinition	2354
BlueprintRun	2354
— 操作 —	2356
CreateWorkflow (create_workflow)	2356
UpdateWorkflow (update_workflow)	2358
DeleteWorkflow (delete_workflow)	2359
GetWorkflow (get_workflow)	2359
ListWorkflows (list_workflows)	2360
BatchGetWorkflows (batch_get_workflows)	2361
GetWorkflowRun (get_workflow_run)	2362
GetWorkflowRuns (get_workflow_runs)	2362
GetWorkflowRunProperties (get_workflow_run_properties)	2363
PutWorkflowRunProperties (put_workflow_run_properties)	2364
CreateBlueprint (create_blueprint)	2365
UpdateBlueprint (update_blueprint)	2366
DeleteBlueprint (delete_blueprint)	2367
ListBlueprints (list_blueprints)	2368
BatchGetBlueprints (batch_get_blueprints)	2368
StartBlueprintRun (start_blueprint_run)	2369
GetBlueprintRun (get_blueprint_run)	2370
GetBlueprintRuns (get_blueprint_runs)	2371
StartWorkflowRun (start_workflow_run)	2372
StopWorkflowRun (stop_workflow_run)	2373
ResumeWorkflowRun (resume_workflow_run)	2373
使用プロファイル	2374
— データ型 —	2374
ProfileConfiguration	2375
ConfigurationObject	2375
UsageProfileDefinition	2376
— 操作 —	2376
CreateUsageProfile (create_usage_profile)	2376
GetUsageProfile (get_usage_profile)	2377
UpdateUsageProfile (update_usage_profile)	2378
DeleteUsageProfile (delete_usage_profile)	2379

ListUsageProfiles (list_usage_profiles)	2380
機械学習	2380
– データ型 –	2381
TransformParameters	2381
EvaluationMetrics	2382
MLTransform	2382
FindMatchesParameters	2385
FindMatchesMetrics	2386
ConfusionMatrix	2388
GlueTable	2388
TaskRun	2389
TransformFilterCriteria	2390
TransformSortCriteria	2391
TaskRunFilterCriteria	2391
TaskRunSortCriteria	2392
TaskRunProperties	2392
FindMatchesTaskRunProperties	2393
ImportLabelsTaskRunProperties	2393
ExportLabelsTaskRunProperties	2394
LabelingSetGenerationTaskRunProperties	2394
SchemaColumn	2394
TransformEncryption	2395
MLUserDataEncryption	2395
ColumnImportance	2395
— 操作 —	2396
CreateMLTransform (create_ml_transform)	2396
UpdateMLTransform (update_ml_transform)	2400
DeleteMLTransform (delete_ml_transform)	2402
GetMLTransform (get_ml_transform)	2403
GetMLTransforms (get_ml_transforms)	2405
ListMLTransforms (list_ml_transforms)	2406
StartMLEvaluationTaskRun (start_ml_evaluation_task_run)	2408
StartMLLabelingSetGenerationTaskRun (start_ml_labeling_set_generation_task_run)	2409
GetMLTaskRun (get_ml_task_run)	2410
GetMLTaskRuns (get_ml_task_runs)	2411
CancelMLTaskRun (cancel_ml_task_run)	2412

StartExportLabelsTaskRun (start_export_labels_task_run)	2413
StartImportLabelsTaskRun (start_import_labels_task_run)	2414
Data Quality	2416
— データ型 —	2416
DataSource	2416
DataQualityRulesetListDetails	2417
DataQualityTargetTable	2417
DataQualityRulesetEvaluationRunDescription	2418
DataQualityRulesetEvaluationRunFilter	2418
DataQualityEvaluationRunAdditionalRunOptions	2419
DataQualityRuleRecommendationRunDescription	2419
DataQualityRuleRecommendationRunFilter	2420
DataQualityResult	2420
DataQualityAnalyzerResult	2421
DataQualityObservation	2422
MetricBasedObservation	2422
DataQualityMetricValues	2423
DataQualityRuleResult	2423
DataQualityResultDescription	2424
DataQualityResultFilterCriteria	2425
DataQualityRulesetFilterCriteria	2425
— 操作 —	2426
StartDataQualityRulesetEvaluationRun (start_data_quality_ruleset_evaluation_run)	2427
CancelDataQualityRulesetEvaluationRun (cancel_data_quality_ruleset_evaluation_run) ...	2428
GetDataQualityRulesetEvaluationRun (get_data_quality_ruleset_evaluation_run)	2429
ListDataQualityRulesetEvaluationRuns (list_data_quality_ruleset_evaluation_runs)	2431
StartDataQualityRuleRecommendationRun (start_data_quality_rule_recommendation_run)	2432
CancelDataQualityRuleRecommendationRun (cancel_data_quality_rule_recommendation_run)	2433
GetDataQualityRuleRecommendationRun (get_data_quality_rule_recommendation_run) .	2434
ListDataQualityRuleRecommendationRuns (list_data_quality_rule_recommendation_runs)	2435
GetDataQualityResult (get_data_quality_result)	2436
BatchGetDataQualityResult (batch_get_data_quality_result)	2438
ListDataQualityResults (list_data_quality_results)	2438
CreateDataQualityRuleset (create_data_quality_ruleset)	2439

DeleteDataQualityRuleset (delete_data_quality_ruleset)	2441
GetDataQualityRuleset (get_data_quality_ruleset)	2441
ListDataQualityRulesets (list_data_quality_rulesets)	2442
UpdateDataQualityRuleset (update_data_quality_ruleset)	2443
機密データ	2445
— データ型 —	2445
CustomEntityType	2445
— 操作 —	2445
CreateCustomEntityType (create_custom_entity_type)	2446
DeleteCustomEntityType (delete_custom_entity_type)	2447
GetCustomEntityType (get_custom_entity_type)	2447
BatchGetCustomEntityTypes (batch_get_custom_entity_types)	2448
ListCustomEntityTypes (list_custom_entity_types)	2449
API のタグ付け	2450
— データ型 —	2450
Tag	2450
— 操作 —	2451
TagResource (tag_resource)	2451
UntagResource (untag_resource)	2452
GetTags (get_tags)	2452
一般的なデータ型	2453
Tag	2453
DecimalNumber	2454
ErrorDetail	2454
PropertyPredicate	2454
ResourceUri	2455
ColumnStatistics	2455
ColumnStatisticsError	2456
ColumnError	2456
ColumnStatisticsData	2456
BooleanColumnStatisticsData	2457
DateColumnStatisticsData	2458
DecimalColumnStatisticsData	2458
DoubleColumnStatisticsData	2459
LongColumnStatisticsData	2459
StringColumnStatisticsData	2460

BinaryColumnStatisticsData	2460
文字列パターン	2461
例外	2463
AccessDeniedException	2463
AlreadyExistsException	2463
ConcurrentModificationException	2463
ConcurrentRunsExceededException	2463
CrawlerNotRunningException	2464
CrawlerRunningException	2464
CrawlerStoppingException	2464
EntityNotFoundException	2464
FederationSourceException	2465
FederationSourceRetryableException	2465
GlueEncryptionException	2465
IdempotentParameterMismatchException	2466
IllegalWorkflowStateException	2466
InternalServiceException	2466
InvalidExecutionEngineException	2466
InvalidInputException	2467
InvalidStateException	2467
InvalidTaskStatusTransitionException	2467
JobDefinitionErrorException	2468
JobRunInTerminalStateException	2468
JobRunInvalidStateTransitionException	2468
JobRunNotInTerminalStateException	2469
LateRunnerException	2469
NoScheduleException	2469
OperationTimeoutException	2469
ResourceNotReadyException	2470
ResourceNumberLimitExceededException	2470
SchedulerNotRunningException	2470
SchedulerRunningException	2470
SchedulerTransitioningException	2471
UnrecognizedRunnerException	2471
ValidationException	2471
VersionMismatchException	2471

AWS Glue API コード例	2472
アクション	2480
CreateCrawler	2481
CreateJob	2493
DeleteCrawler	2505
DeleteDatabase	2511
DeleteJob	2516
DeleteTable	2523
GetCrawler	2528
GetDatabase	2537
GetDatabases	2546
GetJob	2549
GetJobRun	2551
GetJobRuns	2559
GetTables	2567
ListJobs	2580
StartCrawler	2586
StartJobRun	2596
シナリオ	2606
クローラーとジョブを開始する	2606
セキュリティ	2719
データ保護	2719
保管中の暗号化	2720
転送中の暗号化	2738
FIPS 準拠	2739
キーの管理	2739
AWS の他のサービスへの AWS Glue の依存関係	2739
開発エンドポイント	2740
ID およびアクセス管理	2741
対象者	2742
アイデンティティを使用した認証	2742
ポリシーを使用したアクセスの管理	2746
AWS Glue と IAM の連携方法	2749
AWS Glue の IAM アクセス許可の設定	2757
AWS Glue のアクセスコントロールポリシーの例	2791
AWS 管理ポリシー	2817

リソース ARN	2825
クロスアカウントアクセス許可の付与	2832
トラブルシューティング	2839
ログ記録とモニタリング	2841
コンプライアンス検証	2842
耐障害性	2844
インフラストラクチャセキュリティ	2844
VPC エンドポイント (AWS PrivateLink)	2845
共有 Amazon VPC	2847
AWS Glue のトラブルシューティング	2848
AWS Glue トラブルシューティング情報の収集	2848
Spark に関連するエラーのトラブルシューティング	2849
エラー: リソースを利用できません。	2850
エラー: VPC の subnetId に S3 エンドポイントまたは NAT ゲートウェイが見つかりません でした	2850
エラー: 必要なセキュリティグループのインバウンドルール	2850
エラー: 必要なセキュリティグループのアウトバウンドルール	2851
エラー: 渡されたロールに AWS Glue サービスのロールを引き受けるアクセス許可を付与す る必要があるため、ジョブの実行に失敗しました	2851
エラー: DescribeVpcEndpoints アクションが不正です。VPC ID vpc-id を検証できませ ん	2851
エラー: DescribeRouteTables アクションが不正です。サブネット ID: VPC id: vpc-id のサ ブネット ID を検証できません	2852
エラー: ec2 の呼び出しに失敗しました: DescribeSubnets	2852
エラー: ec2 の呼び出しに失敗しました。DescribeSecurityGroups	2852
エラー: AZ のサブネットが見つかりませんでした	2852
エラー: JDBC ターゲットへの書き込み時のジョブ実行の例外	2852
Error: Amazon S3: The operation is not valid for the object's storage class	2853
エラー: Amazon S3 タイムアウト	2853
エラー: Amazon S3 へのアクセスが拒否されました	2854
エラー: Amazon S3 のアクセスキー ID が存在しません	2854
エラー: URI に s3a:// を使用しながら Amazon S3 にアクセスするとジョブ実行が失敗し ます	2854
エラー: Amazon S3 のサービストークンが有効期限切れです	2856
エラー: ネットワークインターフェイスのプライベート DNS が見つかりません	2856
エラー: 開発エンドポイントのプロビジョニングに失敗しました	2857

エラー: ノートブックサーバー CREATE_FAILED	2857
エラー: ローカルノートブックの起動に失敗する	2857
エラー: クローラの実行に失敗しました	2858
エラー: パーティションが更新されませんでした	2858
Error: Job bookmark update failed due to version mismatch (エラー: バージョンの不一致のため、ジョブのブックマークの更新に失敗しました)	2858
エラー: ジョブのブックマークが有効なときにジョブがデータを再処理しています	2859
エラー: の VPCs間のフェイルオーバー動作 AWS Glue	2860
クローラーが Lake Formation の認証情報を使用している場合のクローラーエラーのトラブルシューティング	2861
Ray エラーをトラブルシューティングする	2863
Ray ジョブのログを検査する	2864
Ray ジョブのエラーをトラブルシューティングする	2864
AWS Glue 機械学習の例外	2866
CancelMLTaskRunActivity	2866
CreateMLTaskRunActivity	2867
DeleteMLTransformActivity	2868
GetMLTaskRunActivity	2868
GetMLTaskRunsActivity	2868
GetMLTransformActivity	2869
GetMLTransformsActivity	2869
GetSaveLocationForTransformArtifactActivity	2869
GetTaskRunArtifactActivity	2870
PublishMLTransformModelActivity	2870
PullLatestMLTransformModelActivity	2871
PutJobMetadataForMLTransformActivity	2871
StartExportLabelsTaskRunActivity	2872
StartImportLabelsTaskRunActivity	2872
StartMLEvaluationTaskRunActivity	2873
StartMLLabelingSetGenerationTaskRunActivity	2874
UpdateMLTransformActivity	2875
AWS Glue のクォータ	2876
AWS Glue パフォーマンスの向上	2877
ジョブタイプのチューニング戦略	2877
Spark パフォーマンスの向上	2877
プッシュダウンによる読み取りの最適化	2878

Amazon S3 に保存されているファイルの述語プッシュダウン	2878
JDBC ソースを操作するときのプッシュダウン	2879
AWS Glue のプッシュダウンに関する注意事項と制限事項	2882
AWS Glue 向けの Auto Scaling の使用	2882
要件	2883
AWS Glue Studio の Auto-Scaling を有効にする	1046
AWS CLI または SDK を使用した Auto Scaling の有効化	1047
Amazon CloudWatch メトリクスを使用した Auto Scaling のモニタリング	2885
Spark UI による Auto Scaling のモニタリング。	2886
Auto Scaling ジョブ実行の DPU 使用量モニタリング	2886
制限事項	2887
実行に上限を設定してワークロードをパーティション化する	2887
ワークロードのパーティション化の有効化	2887
ジョブを自動的に実行する AWS Glue トリガーの設定	2889
既知の問題	2890
クロスジョブデータアクセスの防止	2890
ドキュメント履歴	2893
以前の更新	2948
AWS 用語集	2950
.....	mmcmli

AWS Glue とは

AWS Glue は、分析を行うユーザーが複数のソースからのデータを簡単に検出、準備、移動、統合できるようにするサーバーレスのデータ統合サービスです。分析、機械学習、アプリケーション開発に使用できます。また、ジョブの作成、実行、ビジネスワークフローの実装のための生産性向上に役立つツールやデータ運用ツールも追加されています。

AWS Glue を使用すれば、70 を超える多様なデータソースを検出して接続し、一元化されたデータカタログでデータを管理できます。抽出、変換、ロード (ETL) パイプラインを視覚的に作成、実行、モニタリングして、データをデータレイクにロードできます。また、Amazon Athena、Amazon EMR、Amazon Redshift Spectrum を使用して、カタログ化されたデータをすぐに検索し、クエリできます。

AWS Glue は、主要なデータ統合機能を単一のサービスに統合します。これは、データ検出、最新の ETL、クリーニング、変換、一元化されたカタログ作成が含まれます。また、サーバーレスなので、管理するインフラストラクチャがありません。AWS Glue は、ETL、ELT、ストリーミングなどのすべてのワークロードを 1 つのサービスで柔軟にサポートすることで、さまざまなワークロードやユーザータイプのユーザーをサポートします。

また、AWS Glue アーキテクチャ全体でデータを簡単に統合できます。AWS 分析サービスと Amazon S3 データレイクを統合します。AWS Glue では、さまざまな技術的スキルセットに合わせてカスタマイズされたソリューションを備えており、デベロッパーからビジネスユーザーまで、すべてのユーザーが使いやすい統合インターフェイスとジョブ作成ツールを用意しています。

オンデマンドで拡張できる AWS Glue を使用すれば、データの価値を最大化する価値の高いアクティビティに集中できます。あらゆるデータサイズに合わせてスケーリングでき、すべてのデータ型とスキーマの差異をサポートします。俊敏性を高め、コストを最適化するため、AWS Glue は、組み込みの高可用性と従量制料金を提供します。

料金については、[AWS Glue の料金](#)を参照してください。

AWS Glue Studio

AWS Glue Studio は、AWS Glueでのデータ統合ジョブの作成、実行、モニタリングを容易にするグラフィカルインターフェイスです。データ変換ワークフローを視覚的に作成し、AWS Glue の Apache Spark ベースのサーバーレス ETL エンジンでシームレスに実行することができます。

AWS Glue Studio を使用すると、データを収集、変換、クリーニングするジョブを作成および管理できます。AWS Glue Studio を使用して、ジョブスクリプトのトラブルシューティングや編集ができます。

トピック

- [AWS Glue の機能](#)
- [AWS Glue のイノベーションについて学ぶ](#)
- [AWS Glue の開始方法](#)
- [AWS Glue へのアクセス](#)
- [関連サービス](#)

AWS Glue の機能

AWS Glue の機能は、次の 3 つの主要なカテゴリに分類されます。

- データの検出と整理
- 分析用データの変換、準備、クリーニング
- データパイプラインの構築とモニタリング

データの検出と整理

- 複数のデータストアを統合して検索 – AWS ですべてのデータをカタログ化することで、複数のデータソースやシンクで保存、インデックス作成、検索を行うことができます。
- データを自動的に検出 – AWS Glue クローラーを使用して自動的にスキーマ情報を推測し、AWS Glue Data Catalog のスキーマ情報に統合します。
- スキーマとアクセス許可を管理 – データベースとテーブルへのアクセスを検証し、制御します。
- さまざまなデータソースに接続 – オンプレミスと AWS の両方で複数のデータソースを活用し、AWS Glue 接続を使用してデータレイクを構築します。

分析用データの変換、準備、クリーニング

- ジョブキャンバスインターフェイスでデータを視覚的に変換 – ビジュアルジョブエディターで ETL プロセスを定義し、データを抽出、変換、ロードするコードを自動的に生成します。
- シンプルなジョブスケジューリングで複雑な ETL パイプラインを構築 – スケジュール、オンデマンド、またはイベントに基づいて AWS Glue ジョブを呼び出します。
- 転送中のストリーミングデータのクリーニングと変換 – 継続的なデータ消費が可能になり、転送中のデータをクリーニングして変換します。これにより、ターゲットデータストアでの分析が数秒でできるようになります。

- 組み込みの機械学習によるデータの重複排除とクリーニング – FindMatches 機能を使用することで、機械学習の専門知識がなくても、分析用のデータをクリーニングして準備できます。この機能は、相互に不完全な一致であるレコードを重複排除して検索します。
- 組み込みのジョブノートブック – AWS Glue ジョブノートブックは、AWS Glue での最小限のセットアップでサーバーレスノートブックを提供するため、すぐに使用を開始できます。
- ETL コードの編集、デバッグ、テスト – AWS Glue インタラクティブセッションを使用することで、データをインタラクティブに探索して準備できます。任意の IDE またはノートブックを使用して、データをインタラクティブに探索、実験、処理できます。
- 機密データの定義、検出、修正 – AWS Glue の機密データ検出により、データパイプラインとデータレイク内の機密データを定義、識別、処理できます。

データパイプラインの構築とモニタリング

- ワークロードに基づいて自動的にスケーリング – ワークロードに基づいて、リソースを動的にスケールアップまたはスケールダウンできます。これにより、ワーカーは必要な場合にのみジョブに割り当てられます。
- イベントベースのトリガーでジョブを自動化 – イベントベースのトリガーでクローラーまたは AWS Glue ジョブを開始し、依存するジョブとクローラーのチェーンを設計します。
- ジョブの実行とモニタリング – 選択したエンジン (Spark または Ray) を使用して AWS Glue ジョブを実行します。また、自動モニタリングツール、AWS Glue ジョブ実行のインサイト、AWS CloudTrail を使用してモニタリングします。Apache Spark UI を使用して、Spark を利用したジョブのモニタリングを改善します。
- ETL と統合アクティビティのワークフローを定義 – ETL のワークフローと、複数のクローラー、ジョブ、トリガーの統合アクティビティを定義します。

AWS Glue のイノベーションについて学ぶ

AWS Glue の最新イノベーションについて学び、AWS Glue を使用して、セルフサービスによるデータの準備を組織全体で可能にする方法を確認できます。

AWS Glue を従来の設定にとらわれずにスケールする方法と、ジョブのモニタリングとパフォーマンスのために AWS Glue を設定する方法を確認できます。

AWS Glue の開始方法

以下のセクションから開始することが推奨されます。

- [AWS Glue クローン作成の概要](#)
- [AWS Glue の概念](#)
- [AWS Glue 用の IAM アクセス許可のセットアップ](#)
- [AWS Glue Data Catalog の開始方法](#)
- [AWS Glue でジョブを作成する](#)
- [AWS Glue のインタラクティブセッションの開始方法](#)
- [AWS Glue でのオーケストレーション](#)

AWS Glue へのアクセス

次のインターフェイスを使用して、AWS Glue ジョブの作成、表示、管理ができます。

- AWS Glue コンソール – AWS Glue ジョブを作成、表示、管理するためのウェブインターフェイスを提供します。コンソールにアクセスするには、[AWS Glue](#) を参照してください。
- AWS Glue Studio – AWS Glue ジョブを視覚的に作成および編集するためのグラフィカルインターフェイスを提供します。詳細については、「[AWS Glue Studio とは?](#)」を参照してください。
- AWS CLI リファレンスの AWS Glue セクション – AWS Glue で使用できる AWS CLI コマンドを提供します。詳細については、「[AWS Glue 向けの AWS CLI リファレンス](#)」を参照してください。
- AWS Glue API – デベロッパー向けの完全な API リファレンスを提供します。詳細については、「[AWS Glue API](#)」を参照してください。

関連サービス

AWS Glue のユーザーは、以下も使用します。

- [AWS Lake Formation](#) – AWS Glue Data Catalog 内のリソースへの詳細に設定されたアクセスコントロールを行う認証レイヤーであるサービス。
- [AWS Glue DataBrew](#) – コードを記述せずに、データのクリーニングおよび正規化に使用できるビジュアルデータ準備ツール。

AWS Glue: 仕組み

AWS Glue は他の AWS のサービスを使用して ETL (抽出、変換、ロード) ジョブを調整し、データウェアハウスとデータレイクを構築して、出カストリームを生成します。AWS Glue は API オペレーションを呼び出して、データの変換、ランタイムログの作成、ジョブロジックの保存を行い、ジョブ実行のモニタリングに役立つ通知を作成します。AWS Glue コンソールはこれらのサービスを管理アプリケーションに接続して、お客様が ETL ワークの作成とモニタリングに集中できるようにします。管理およびジョブ開発のオペレーションは、コンソールがお客様に代わって実行します。データソースへのアクセスとデータターゲットへの書き込みを行うために必要な、認証情報と他のプロパティは、お客様が AWS Glue に提供する必要があります。

AWS Glue は、ワークロードを実行するために必要なリソースのプロビジョニングおよび管理を行います。AWS Glue が代わって行うため、ETL ツールのインフラストラクチャを作成する必要はありません。リソースが必要な場合、起動時間を削減するために、AWS Glue はインスタンスのウォームプールからインスタンスを使用してワークロードを実行します。

AWS Glue では、データカタログにあるテーブル定義を使用してジョブを作成します。ジョブは、変換を実行するプログラミングロジックを含むスクリプトで構成されます。トリガーを使用し、スケジュールに基づいて、または指定されたイベントの結果としてジョブを開始します。ターゲットデータが存在する場所、およびターゲットに入力するソースデータを指定します。入力により、AWS Glue はデータをソースからターゲットに変換するのに必要なコードを生成します。AWS Glue コンソールまたは API でスクリプトを提供してデータを処理することもできます。

データソースと送信先

AWS Glue for Spark は、次のような複数のシステムやデータベースとの間でデータを読み書きできます。

- Amazon S3
- Amazon DynamoDB
- Amazon Redshift
- Amazon Relational Database Service (Amazon RDS)
- JDBC アクセスが可能なサードパーティのデータベース
- MongoDB と Amazon DocumentDB (MongoDB 互換)
- その他のマーケットプレイスコネクタと Apache Spark プラグイン

データストリーム

AWS Glue for Spark は、次のシステムからデータをストリーミングできます。

- Amazon Kinesis Data Streams
- Apache Kafka

AWS Glue は複数の AWS リージョンで利用可能です。詳細については、[AWS](#) の「Amazon Web Services 全般のリファレンス のリージョンとエンドポイント」を参照してください。

トピック

- [独立で実行されるサーバーレス ETL ジョブ](#)
- [AWS Glue の概念](#)
- [AWS Glue コンポーネント](#)
- [AWS Glue for Spark と AWS Glue for Ray](#)
- [AWS Glue を使用して半構造化されたスキーマをリレーショナルスキーマに変換する](#)
- [AWS Glue ータイプシステム](#)

独立で実行されるサーバーレス ETL ジョブ

AWS Glue は、Spark または Ray から選択したエンジンを使用して、サーバーレス環境で ETL ジョブを実行します。AWS Glue は、独自のサービスアカウントでプロビジョニングして管理する仮想リソースでこれらのジョブを実行します。

AWS Glue は、以下を実行するよう設計されています。

- お客様のデータを分離します。
- 伝送中と保管時のお客様のデータを保護します。
- 一時的な制限された認証情報を使用して、またはアカウント内の IAM ロールに対するお客様の同意を得て、お客様のリクエストに応え必要な時だけお客様のデータにアクセスします。

ETL ジョブのプロビジョニング時に、Virtual Private Cloud (VPC) にある入力データソースおよび出力データターゲットを提供します。また、データソースおよびターゲットにアクセスするために必要な、IAM ロール、VPC ID、サブネット ID、およびセキュリティグループを提供します。各タブ (顧客アカウント ID、IAM ロール、サブネット ID、およびセキュリティグループ) に、AWS Glue

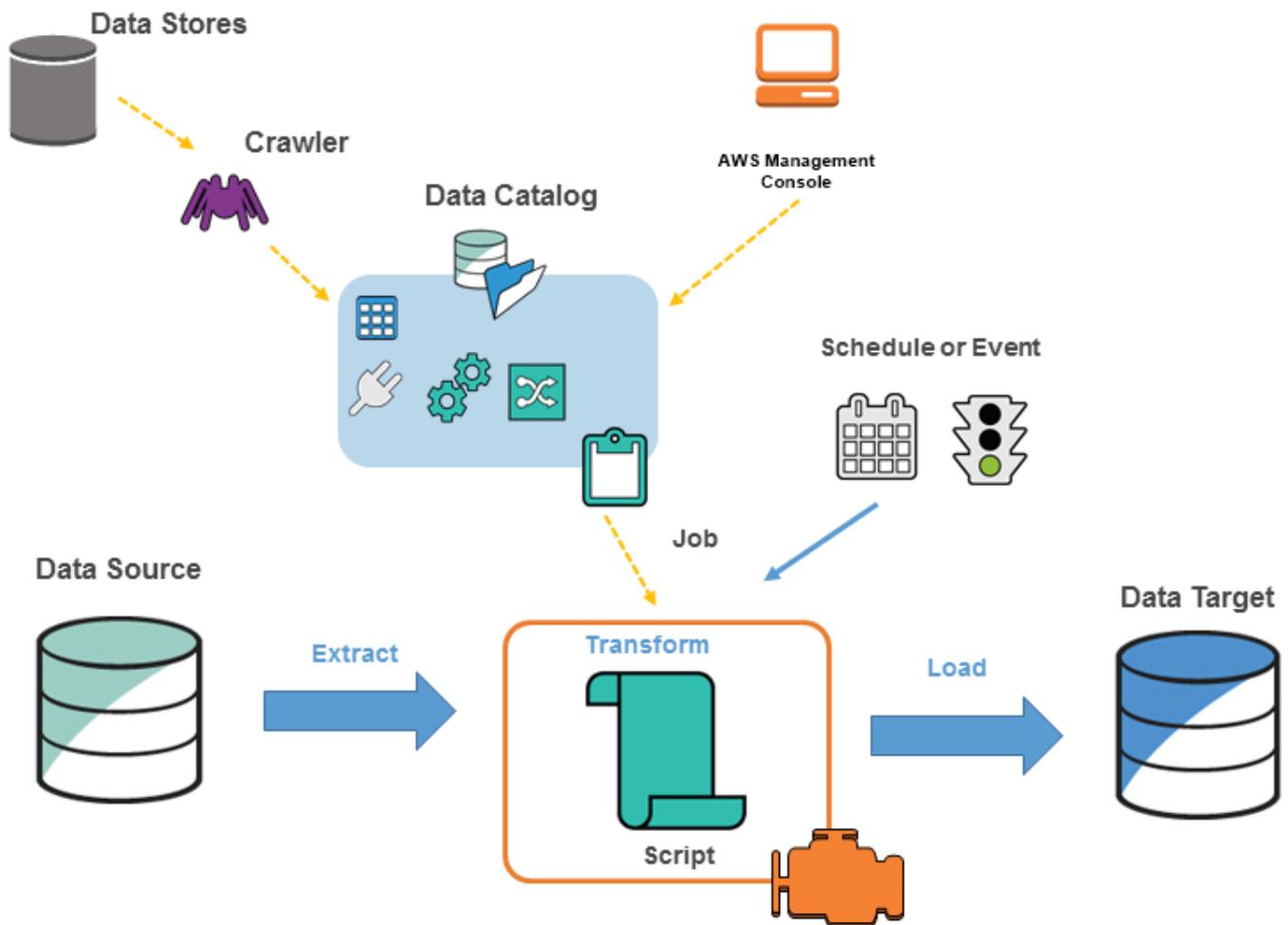
は、AWS Glue サービスアカウント内に存在する他のすべての環境からネットワークおよび管理レベルで分離された新しい環境を作成します。

AWS Glue はプライベート IP アドレスを使用して、サブネットで Elastic Network Interface を作成します。ジョブはこれらの Elastic Network Interface を使用して、データソースおよびデータターゲットにアクセスします。ジョブ実行環境内外へのトラフィックと、ジョブ実行環境内でのトラフィックは、VPC およびネットワークポリシーにより管理されます。ただし、1 つ例外があり、AWS Glue ライブラリに対する呼び出しは、AWS Glue VPC 経由で AWS Glue API オペレーションにトラフィックをプロキシできます。すべての AWS Glue API 呼び出しはログに記録されます。そのため、データの所有者は、監査ログをアカウントに配信する [AWS CloudTrail](#) を有効にすることで API アクセスを監査できます。

ETL ジョブを実行する AWS Glue 管理の環境は、他の AWS のサービスと同じセキュリティ実施方法で保護されています。プラクティスと共有されたセキュリティ責任の概要については、「[Introduction to AWS Security Processes](#)」のホワイトペーパーを参照してください。

AWS Glue の概念

次の図は、AWS Glue 環境のアーキテクチャを示します。



でジョブAWS Glueを指定し、データソースからデータターゲットへのデータの抽出、変換、ロード (ETL) に必要な作業を完了します。通常は、以下のアクションを実行します。

- データストアソースの場合は、クローラを指定し、AWS Glue Data Catalog にメタデータテーブルの定義を入力します。データストアでクローラを指定し、クローラは Data Catalog のテーブル定義を作成します。ストリーミングソースの場合は、Data Catalog テーブルを手動で定義し、データストリームのプロパティを指定します。

テーブル定義に加えて、AWS Glue Data Catalog には ETL ジョブを定義するために必要な他のメタデータが含まれています。このメタデータを使用して、データを変換するジョブを定義できます。

- AWS Glue はデータを変換するスクリプトを生成できます。または、AWS Glue コンソールまたは API でスクリプトを提供できます。

- ジョブをオンデマンドで実行する、または、指定したトリガーが発生すると開始するようにセットアップできます。トリガーは、時間ベースのスケジュールまたはイベントです。

ジョブが実行されると、スクリプトはデータソースからデータを抽出し、データを変換してデータターゲットにロードします。スクリプトは AWS Glue の Apache Spark 環境で実行されます。

Important

AWS Glue のテーブルとデータベースは、AWS Glue Data Catalog のオブジェクトです。それらにはメタデータが含まれ、データストアからのデータは含まれません。

AWS Glue で正常に処理するために、CSV などのテキストベースのデータは **UTF-8** でエンコードする必要があります。詳細については、Wikipedia の「[UTF-8](#)」を参照してください。

AWS Glue の用語

AWS Glue は、複数のコンポーネントの相互作用に依存して、抽出、変換、ロード (ETL) ワークフローを作成および管理します。

AWS Glue Data Catalog

AWS Glue の持続的なメタデータストア。これには、AWS Glue 環境を管理するためのテーブル定義、ジョブ定義、およびその他のコントロール情報が含まれています。各 AWS アカウントには、リージョンごとに 1 つの AWS Glue Data Catalog があります。

分類子

データのスキーマを決定します。AWS Glue は、一般的なファイルタイプの分類子を提供します (CSV、JSON、AVRO、XML など)。また、JDBC 接続を使用する一般的なリレーショナルデータベース管理システムの分類子を提供します。独自の分類子を記述するには、grok パターンを使用する、または、XML ドキュメント内の行タグを指定します。

Connection

特定のデータストアに接続するために必要なプロパティを含む Data Catalog オブジェクトです。

Crawler

データストア (ソースまたはターゲット) に接続し、分類子の優先順位リストを進行してデータのスキーマを判断し、AWS Glue Data Catalog にメタデータテーブルを作成するプログラムです。

データベース

論理グループに分類される、一連の関連付けられた Data Catalog テーブル定義です。

データストア、データソース、データターゲット

データストアは、データを永続的に保存するリポジトリです。例として、Amazon S3 バケット、リレーショナルデータベースなどがあります。データソースは、プロセスまたは変換への入力として使用されるデータストアです。データターゲットはプロセスまたは変換の書き込み先であるデータストアです。

開発エンドポイント

エンドポイントは、AWS Glue ETL スクリプトの開発およびテストに使用できる環境です。

動的フレーム

構造や配列などのネストされたデータをサポートする分散テーブルです。各レコードは自己記述型であり、半構造化データのスキーマの柔軟性を持つよう設計されています。各レコードには、データとそのデータを記述するスキーマの両方が含まれます。動的フレームと Apache Spark DataFrames の両方を ETL スクリプトで使用し、それらの間で変換できます。動的フレームは、データクリーニングと ETL 用の一連の高度な変換を提供します。

ジョブ

ETL 作業を実行するために必要なビジネスロジックです。変換スクリプト、データソース、およびデータターゲットで構成されます。ジョブ実行は、スケジュールされたトリガーにより、または、イベントにトリガーされることで開始されます。

ジョブパフォーマンスダッシュボード

AWS Glue は、ETL ジョブ用の包括的な実行ダッシュボードを提供します。ダッシュボードには、特定の時間枠からのジョブ実行に関する情報が表示されます。

ノートブックインターフェイス。

ジョブ作成とデータ探索を容易にするワンクリック設定により、ノートブック体験を強化します。ノートブックと接続は自動的に構成されます。Jupyter Notebook をベースとしたノートブックインターフェイスにより、AWS Glue のサーバーレス Apache Spark ETL インフラストラクチャを使用するスクリプトやワークフローを、対話的に開発、デバッグし、デプロイすることができます。また、ノートブック環境では、アドホッククエリ、データ分析、ビジュアライゼーション (表やグラフなど) を実行できます。

スクリプト

ソースからデータを抽出し、変換し、ターゲットにロードするコード。AWS Glue は PySpark または Scala スクリプトを生成します。

テーブル

データを表すメタデータ定義。データが、Amazon Simple Storage Service (Amazon S3) ファイル、Amazon Relational Database Service (Amazon RDS) テーブル、または別の一連のデータのどれにある場合でも、テーブルはデータのスキーマを定義します。AWS Glue Data Catalog のテーブルは、列名、データ型の定義、パーティション情報、および基本データセットに関するその他のメタデータで構成されています。データのスキーマは AWS Glue のテーブル定義で表されます。実際のデータは、ファイルまたはリレーショナルデータベーステーブルにあっても、元のデータストアに残ります。AWS Glue はファイルとリレーショナルデータベースのテーブルを AWS Glue Data Catalog に格納します。それらは、ETL ジョブを作成する際にソースおよびターゲットとして使用されます。

変換

データを操作して別の形式にするために使用するコードのロジック。

Trigger (トリガー)

ETL ジョブを開始します。トリガーはスケジュールされた時間またはイベントに基いて定義できます。

ビジュアルジョブエディタ

ビジュアルジョブエディターは、AWS Glue での抽出、変換、ロード (ETL) ジョブの作成、実行、およびモニタリングが簡単に行えるグラフィカルなインターフェイスです。データ変換ワークフローを

視覚的に構成し、AWS Glue の Apache Spark ベースのサーバーレス ETL エンジンでそれらをシームレスに実行して、ジョブの各ステップでスキーマとデータの結果を検査できます。

ワーカー

AWS Glue では、ETL ジョブの実行にかかる時間に対してのみお支払いが発生します。管理するリソースはなく、前払い料金もありません。また、起動時間やシャットダウン時間に対しての課金もありません。ETL ジョブの実行に使用されたデータ処理ユニット (DPUs) の数に基づいて時間あたりの料金が請求されます。単一のデータ処理ユニット (DPU) は、ワーカーとも呼びます。AWS Glue には、ジョブのレイテンシーとコスト要件を満たす構成を選択するのに役立つような 3 つのワーカータイプが用意されています。ワーカーには、Standard、G.1X、G.2X、および G.025X 構成があります。

AWS Glue コンポーネント

AWS Glue は、抽出、変換、ロード (ETL) ワークロードを設定し管理するためのコンソールと API オペレーションを備えています。いくつかの言語に固有な SDK と AWS Command Line Interface (AWS CLI) を介して API オペレーションを使用できます。AWS CLI の使用については、「[AWS CLI コマンドリファレンス](#)」を参照してください。

AWS Glue は AWS Glue Data Catalog を使用して、データソース、変換、およびターゲットについてのメタデータを保存します。Data Catalog は Apache Hive メタストアのドロップインリプレースメントです。AWS Glue Jobs system は、データの ETL オペレーションの定義、スケジューリング、および実行のためのマネージド型インフラストラクチャを備えています。AWS Glue API の詳細については、「[AWS Glue API](#)」を参照してください。

AWS Glue コンソール

AWS Glue コンソールを使用して、ETL ワークフローを定義しオーケストレーションします。コンソールは AWS Glue Data Catalog および AWS Glue Jobs system のいくつかの API オペレーションを呼び出して、次のタスクを実行します。

- ジョブ、テーブル、クローラ、接続などの AWS Glue オブジェクトを定義します。
- いくつかのクローラが実行するかをスケジュールします。
- ジョブトリガーのイベントやスケジュールを定義します。
- AWS Glue オブジェクトのリストを検索しフィルタリングします。
- 変換スクリプトを編集します。

AWS Glue Data Catalog

AWS Glue Data Catalog は AWS クラウド内にある永続的な技術メタデータストアです。

各 AWS アカウントには、AWS リージョンごとに 1 つの AWS Glue Data Catalog があります。各データカタログは、データベースに構成された高度にスケーラブルなテーブルのコレクションです。テーブルは、Amazon RDS、Apache Hadoop 分散ファイルシステム、Amazon OpenSearch Service などのソースに保存されている構造化データまたは半構造化データの集合をメタデータで表現したものです。AWS Glue Data Catalog は統一されたリポジトリを提供するため、異種システムはデータサイロのデータを追跡するためにメタデータを保存し検索することができます。その後は、メタデータを使用して、多種多様なアプリケーション全体で一貫したやり方でそのデータをクエリしたり変換したりすることができます。

データカタログを AWS Identity and Access Management ポリシーと Lake Formation と併用することで、テーブルやデータベースへのアクセスを制御します。この方法により、企業内の様々なグループが、機密情報を高度かつ詳細に保護しながら、より広範な組織に対してデータを安全に公開できるようになります。

また、データカタログを CloudTrail や Lake Formation を併用することで、スキーマ変更の追跡やデータアクセス制御を備えた包括的な監査およびガバナンス機能も提供します。これにより、データの不適切な変更や誤った共有を防止できます。

AWS Glue Data Catalog のセキュリティ保護と監査の詳細については、以下を参照してください。

- AWS Lake Formation – 詳細については、AWS Lake Formation デベロッパーガイドの「[AWS Lake Formation とは？](#)」を参照してください。
- CloudTrail – 詳細については、「AWS CloudTrail ユーザーガイド」の「[CloudTrail とは](#)」を参照してください。

次に示すのは、AWS Glue Data Catalog を使用する AWS の他のサービスとオープンソースプロジェクトです。

- Amazon Athena – 詳細については、Amazon Athena ユーザーガイドの「[テーブル、データベース、およびデータカタログの理解](#)」を参照してください。
- Amazon Redshift Spectrum – 詳細については、Amazon Redshift データベースデベロッパーガイドの「[Amazon Redshift Spectrum を使用した外部データのクエリ](#)」を参照してください。
- Amazon EMR – 詳細については、Amazon EMR マネジメントガイドの「[AWS Glue Data Catalog への Amazon EMR アクセスにリソースベースのポリシーを使用する](#)」を参照してください。

- Apache Hive メタストア用 AWS Glue Data Catalog クライアント – この GitHub プロジェクトの詳細については、「[AWS Glue Data Catalog Client for Apache Hive Metastore](#)」を参照してください。

AWS Glue クローラおよび分類子

AWS Glue では、あらゆる種類のリポジトリにあるデータのスキャン、分類、スキーマ情報の抽出、そのメタデータの AWS Glue Data Catalog への自動保存ができるクローラを設定することもできます。そこから AWS Glue Data Catalog は ETL オペレーションをガイドするのに使用できます。

クローラおよび分類子の設定方法については、「[クローラーを使用したデータカタログへの入力](#)」を参照してください。AWS Glue API を使用してクローラおよび分類子をプログラムする方法については、「[クローラーおよび分類子 API](#)」を参照してください。

AWS Glue ETL オペレーション

AWS Glue は、Data Catalog のメタデータを使用して、さまざまな ETL オペレーションを実行するために使用や変更ができる AWS Glue 拡張機能を備えた、Scala または PySpark (Apache Spark 用の Python API) スクリプトを自動生成できます。たとえば、未加工データを抽出、クリーンアップ、および変換してからその結果を別のリポジトリに保存して、クエリと分析を行うことができます。このようなスクリプトは、CSV ファイルをリレーショナル形式に変換し、Amazon Redshift に保存する場合があります。

AWS Glue ETL 機能の使用方法的詳細については、「[Spark スクリプトのプログラミング](#)」を参照してください。

AWS Glue でのストリーミング ETL

AWS Glue を使用すると、連続実行のジョブを使用して、ストリーミングデータに対して ETL 操作を実行できます。AWS Glue のストリーミング ETL は、Apache Spark 構造化ストリーミングエンジン上に構築され、Amazon Kinesis Data Streams、Apache Kafka、および Amazon Managed Streaming for Apache Kafka (Amazon MSK) からストリームを取り込むことができます。ストリーミング ETL では、ストリーミングデータのクリーニングと変換を行い、Simple Storage Service (Amazon S3) または JDBC データストアにロードできます。AWS Glue のストリーミング ETL を使用すると、IoT ストリーム、クリックストリーム、ネットワークログなどのイベントデータを処理できます。

ストリーミングデータソースのスキーマがわかっている場合は、Data Catalog テーブルで指定できます。そうでない場合は、ストリーミング ETL ジョブでスキーマ検出を有効にできます。ジョブは、受信データからスキーマを自動的に決定します。

ストリーミング ETL ジョブは、AWS Glue 組み込みの変換および Apache Spark 構造化ストリーミングネイティブの変換の両方を使用できます。詳細については、「[Operations on streaming DataFrames/Datasets](#)」を参照してください。

詳しくは、「[the section called “ストリーミング ETL ジョブ”](#)」を参照してください。

AWS Glue ジョブシステム

AWS Glue Jobs system は、ETL ワークフローをオーケストレーションするためのマネージド型インフラストラクチャを提供します。データを抽出したり変換したり異なる場所へ転送したりするのに使用するスクリプトを自動化するジョブを AWS Glue で作成できます。ジョブはスケジュールしたり連鎖させることができます。または新しいデータの到着などのイベントによってトリガーすることができます。

AWS Glue Jobs system の使用方法の詳細については、「[AWS Glue のモニタリング](#)」を参照してください。AWS Glue Jobs system API を使用したプログラミングについては、「[ジョブ API](#)」を参照してください。

ビジュアル ETL コンポーネント

AWS Glue では、操作可能なビジュアルキャンバスから ETL ジョブを作成できます。

Visual | Script | Job details | Runs | Data quality **New** | Schedules | Version Control

Try new UI | Actions | Save | Run

Data source - S3 bucket
S3 bucket

Data source - S3 bucket
Amazon S3

Transform - ApplyMapping
ApplyMapping

Data target - S3 bucket
S3 bucket

Unsaved job found
We found an unsaved graph, do you wish to restore it? **Restore**

ETL ジョブメニュー

キャンバスの上部にあるメニューオプションから、ジョブに関するさまざまなビューや設定の詳細にアクセスできます。

- [Visual] - ビジュアルジョブエディタのキャンバスです。ここでは、ノードを追加してジョブを作成できます。
- [スクリプト] - ETL ジョブのスクリプト表現です。AWS Glue では、ジョブのビジュアル表現に基づいてスクリプトを生成します。スクリプトを編集したり、ダウンロードしたりもできます。

i Note

スクリプトを編集することを選択した場合、ジョブ作成体験は完全にスクリプト専用モードに変更されます。それ以後、ビジュアルエディタを使用してジョブを編集することは

できません。そのため、スクリプトを編集することを選択する前に、すべてのジョブソース、変換、ターゲットを追加し、ビジュアルエディタで必要な変更をすべて行ってください。

- [ジョブの詳細] - [ジョブの詳細] タブでは、ジョブのプロパティを設定することでジョブを設定できます。基本的なプロパティとして、ジョブの名前と説明、IAM ロール、ジョブタイプ、AWS Glue バージョン、言語、ワーカータイプ、ワーカー数、ジョブブックマーク、Flex 実行、廃止数、ジョブタイムアウトなどがあります。高度なプロパティとして、接続、ライブラリ、ジョブパラメータ、タグなどがあります。
- [Runs] - ジョブの実行後に、このタブにアクセスして過去のジョブを表示できます。
- [Data quality] - [Data quality] は、データセットの品質を評価およびモニタリングします。このタブでは、データ品質の活用方法を詳しく確認したり、データ品質変換をジョブに追加したりできます。
- [Schedules] - このタブには、スケジュールされたジョブが表示されます。このジョブにスケジュールがアタッチされていない場合、このタブにはアクセスできません。
- [Version control] - ジョブを Git リポジトリに設定することで、ジョブで Git を使用できます。

ビジュアル ETL パネル

キャンバスで作業するとき、ノードの設定、データのプレビュー、出力スキーマの表示などに役立つパネルがいくつか用意されています。

- [プロパティ] - キャンバス上のノードを選択すると、[プロパティ] パネルが表示されます。
- [Data preview] - [Data preview] パネルには、データ出力のプレビューが表示されるため、ジョブを実行して出力を確認する前に判断を行えます。
- [Output schema] - [Output schema] タブでは、変換ノードのスキーマを表示および編集できます。

パネルサイズを変更する

画面の右側にある [プロパティ] パネルと、[Data preview] タブや [Output schema] タブを含む下部パネルのサイズを変更するには、パネルの端をクリックして左右または上下にドラッグします。

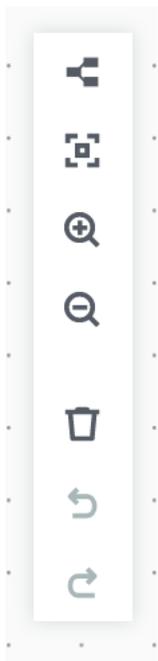
- [プロパティ] パネル - [プロパティ] パネルのサイズを変更するには、画面の右側にあるキャンバスの端をクリックしてドラッグし、左にドラッグして幅を拡大します。デフォルトでは、パネルは折りたたまれており、ノードを選択すると、[プロパティ] パネルがデフォルトサイズで表示されます。

- [Data preview] と [Output schema] パネル - 下部パネルのサイズを変更するには、画面の下部にあるキャンバスの下端をクリックしてドラッグし、上にドラッグして高さを拡大します。デフォルトでは、パネルは折りたたまれており、ノードを選択すると、下部パネルがデフォルトサイズで表示されます。

ジョブキャンバス

ビジュアル ETL キャンバスでは、ノードの追加、削除、移動/順序変更を直接行えます。これは、データソースで始まりデータターゲットで終わる、完全に機能する ETL ジョブを作成するためのワークスペースと考えてください。

キャンバス上のノードを操作するとき、ツールバーを使用できます。ツールバーでは、ズームインとズームアウト、ノードの削除、ノード間の接続の作成または編集、ジョブフローの向きの変更、アクションを元に戻す/やり直すなどが行えます。

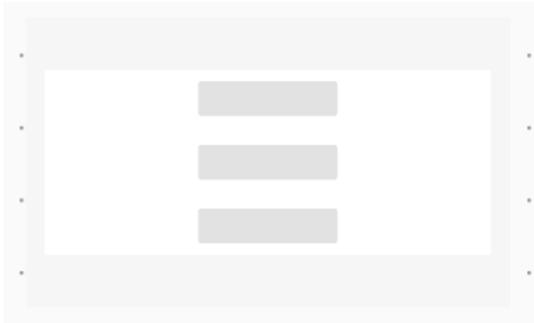


フローティングツールバーはキャンバス画面の右上に固定されており、アクションを実行する複数のアイコン画像が表示されています。

- レイアウトアイコン - レイアウトアイコンは、ツールバーの一番上のアイコンです。デフォルトでは、ビジュアルジョブの方向は上から下です。ノードを左から右に水平に配置することで、ビジュアルジョブの配置方向を変更します。レイアウトアイコンを再度クリックすると、方向は上から下に戻ります。

- リセンタリングアイコン - リセンタリングアイコンは、キャンバスの表示位置を中央に変更します。大規模なジョブで、中央位置に戻るときに活用できます。
- ズームインアイコン - ズームインアイコンは、キャンバス上のノードのサイズを拡大します。
- ズームアウトアイコン - ズームアウトアイコンは、キャンバス上のノードのサイズを縮小します。
- ゴミ箱アイコン - ゴミ箱アイコンは、ビジュアルジョブからノードを削除します。先にノードを選択しておく必要があります。
- 元に戻すアイコン - 元に戻すアイコンは、ビジュアルジョブで最後に実行したアクションを取り消します。
- やり直すアイコン - やり直すアイコンは、ビジュアルジョブで最後に実行されたアクションをやり直します。

ミニマップを使用する



リソースパネル

リソースパネルには、使用可能なデータソース、変換アクション、接続がすべて含まれています。キャンバスでリソースパネルを表示するには、[+] アイコンをクリックします。これで、リソースパネルが開きます。

リソースパネルを閉じるには、リソースパネルの右上隅にある [X] をクリックします。これにより、再び開く準備ができるまでパネルが非表示になります。

+ Add nodes
✕

▼ Popular transforms & data

 Amazon S3 (source)	 SQL Query
 Amazon Redshift (source)	 Aggregate
 Change Schema	 Custom Transform
 Join	 Filter

Transforms

Data

▼ Sources

- 

AWS Glue Data Catalog

AWS Glue Data Catalog table as the data source.
- 

Amazon S3

JSON, CSV, or Parquet files stored in S3.
- 

Amazon Kinesis

Read from an Amazon Kinesis Data Stream.
- 

Apache Kafka

Read from an Apache Kafka or Amazon MSK topic.
- 

Relational DB

AWS Glue Data Catalog table with a relational database as the data source.
- 

Amazon Redshift

Read your data from Amazon Redshift.
- 

MySQL

AWS Glue Data Catalog table with MySQL as the data source.
- 

PostgreSQL

AWS Glue Data Catalog table with PostgreSQL as the data source.
- 

Oracle SQL

AWS Glue Data Catalog table with Oracle SQL as the data source.
- 

Microsoft SQL Server

AWS Glue Data Catalog table with SQL Server as the data source.
- 

Amazon DynamoDB

AWS Glue Data Catalog table with DynamoDB as the data source.
- 

Snowflake

Read your data from Snowflake.

[Popular transforms & data]

パネルの上部に、[Popular transforms & data] というコレクションがあります。これらのノードは、AWS Glue で一般的に使用されるものです。いずれかを選択してキャンバスに追加します。また、[Popular transforms & data] の見出しの横にある三角アイコンをクリックすると、[Popular transforms & data] を非表示にできます。

[Popular transforms & data] セクションの下で、変換およびデータソースノードを検索できます。入力すると結果が表示されます。検索クエリに追加する文字が多いほど、表示される結果は少なくなります。検索結果には、ノードの名前または説明が表示されます。結果に表示されたノードを選択してキャンバスに追加します。

[Transforms] と [Data]

ノードを [Transforms] と [Data] に整理する 2 つのタブがあります。

[Transforms] – [Transforms] タブを選択すると、使用可能なすべての変換を選択できます。変換を選択してキャンバスに追加します。また、[Transforms] リストの下部にある [変換の追加] をクリックすると、[カスタムビジュアル変換](#)を作成するためのドキュメントへの新しいページが開きます。手順に従うことで、独自の変換を作成できます。この変換はその後、使用可能な変換のリストに表示されません。

[Data] – [Data] タブには、[ソース] と [ターゲット] のノードがすべて含まれています。[ソース] と [ターゲット] を非表示にするには、[ソース] または [ターゲット] の見出しの横にある三角アイコンをクリックします。三角アイコンを再度クリックすると、[ソース] と [ターゲット] を再表示できます。ソースノードまたはターゲットノードを選択し、キャンバスに追加します。また、[Manage Connections] をクリックすると、新しい接続を追加できます。このとき、コンソールのコネクタページが開きます。

AWS Glue for Spark と AWS Glue for Ray

AWS Glue on Apache Spark (AWS Glue ETL) では、PySpark を使用して、大規模なデータを処理する Python コードを書くことができます。Spark はこの問題に対する一般的な解決策ですが、Python 専門の経歴を持つデータエンジニアは、直感的な移行ではないと感じるかもしれません。Spark DataFrame モデルはシームレスに「Python らしい」というわけではなく、Scala 言語と、それが構築されている Java ランタイムを反映しています。

AWS Glue では、Python シェルジョブを使用してネイティブな Python データ統合を実行できます。これらのジョブは、単一の Amazon EC2 インスタンス上で実行され、そのインスタンスの容量に

よって制限されます。これにより、処理できるデータのスループットが制限され、ビッグデータを処理する場合は維持費が高くなります。

AWS Glue for Ray では、Spark の学習に多額の投資をすることなく Python ワークロードをスケールアップできます。Ray のパフォーマンスが向上する特定のシナリオを使用できます。選択肢が提供されることで、Spark と Ray の両方の強みを活用できます。

AWS Glue ETL と AWS Glue for Ray は根本的に異なるため、サポートしている機能も異なります。サポートされている機能については、ドキュメントを確認してください。

AWS Glue for Ray とは

Ray は Python を中心に、ワークロードのスケールアップに使用できるオープンソースの分散計算フレームワークです。Ray の詳細については、[Ray のウェブサイト](#)を参照してください。AWS GlueRay ジョブおよびインタラクティブセッションにより、AWS Glue 内で Ray を使用できます。

AWS Glue for Ray を使用して、複数のマシンで並列に実行される計算用の Python スクリプトを作成できます。Ray ジョブおよびインタラクティブセッションでは、pandas などの使い慣れた Python ライブラリを使用して、ワークフローを簡単に記述して実行できます。Ray データセットの詳細については、Ray ドキュメントの「[Ray データセット](#)」を参照してください。pandas の詳細については、[Pandas のウェブサイト](#)を参照してください。

AWS Glue for Ray を使用すると、わずか数行のコードだけで、企業規模のビッグデータに対する pandas のワークフローを実行できます。Ray ジョブは、AWS Glue コンソールまたは AWS SDK から作成できます。AWS Glue インタラクティブセッションを開いて、サーバーレスの Ray 環境でコードを実行することもできます。AWS Glue Studio のビジュアルジョブはまだサポートされていません。

AWS Glue for Ray ジョブを使用すると、スケジュールに従って、または Amazon EventBridge からのイベントに応じてスクリプトを実行できます。ジョブはログ情報とモニタリング統計を CloudWatch に保存するため、スクリプトの正常性と信頼性を確認できます。AWS Glue ジョブのシステムについての詳細は、「[the section called “Ray ジョブの使用”](#)」を参照してください。

AWS Glue for Ray のインタラクティブセッション (プレビュー) では、プロビジョニングされた同じリソースに対してコードのスニペットを次々と実行できます。これを使用して、スクリプトのプロトタイプ作成や開発を効率的に行ったり、独自のインタラクティブアプリケーションを構築できます。AWS Glue のインタラクティブセッションは、AWS Management Console で AWS Glue Studio ノートブックから使用できます。詳細については、「[AWS Glue Studio と AWS Glue でのノートブックの使用](#)」を参照してください。また、Jupyter カーネルを介して使用することもできます。こ

れにより、VSCode などの Jupyter ノートブックをサポートする既存のコード編集ツールからインタラクティブセッションを実行できます。詳細については、「[the section called “AWS Glue for Ray インタラクティブセッション \(プレビュー\)”](#)」を参照してください。

Ray は、負荷に基づいてリアルタイムで再設定を行うマシンのクラスターに処理を分散することで、Python コードのスケールアップ作業を自動化します。これにより、特定のワークロードで 1 ドルあたりのパフォーマンスが向上します。Ray ジョブでは、AWS Glue のジョブモデルに自動スケールアップがネイティブに組み込まれているため、この機能を最大限に活用できます。Ray ジョブは AWS Graviton 上で実行されるため、全体的な価格パフォーマンスが向上します。

コストの削減に加えて、ネイティブの自動スケールアップを使用すると、クラスターのメンテナンス、調整、管理に時間を費やすことなく Ray のワークロードを実行できます。pandas や AWS SDK for Pandas など、使い慣れたオープンソースライブラリをそのまま使用できます。これにより、AWS Glue for Ray で開発する際のイテレーション速度が向上します。AWS Glue for Ray を使用すると、費用対効果の高いデータ統合ワークロードをすばやく開発して実行できるようになります。

AWS Glue を使用して半構造化されたスキーマをリレーショナルスキーマに変換する

半構造化データをリレーショナルテーブルに変換することが一般的です。概念的には、階層的なスキーマをリレーショナルスキーマに平坦化します。AWS Glue は、この変換を臨機応変に実行できます。

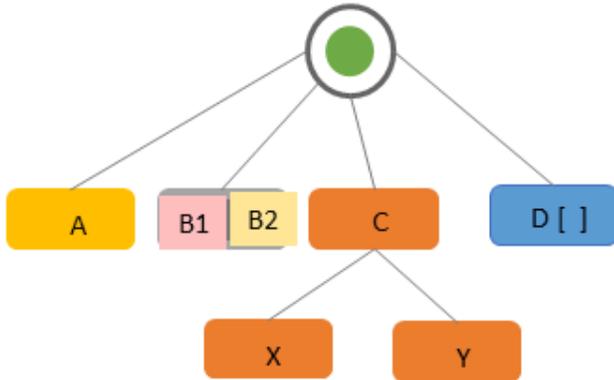
通常、半構造化されたデータには、データ内のエンティティを識別するためのマークアップが含まれています。固定されたスキーマのない、ネスト化されたデータ構造を持つことができます。半構造化データの詳細については、Wikipedia の [半構造化データ](#) を参照してください。

リレーショナルデータは、行と列で構成されるテーブルで表されます。テーブル間の関係は、プライマリキー (PK) と外部キー (FK) の関係によって表すことができます。詳細については、Wikipedia の [リレーショナルデータベース](#) を参照してください。

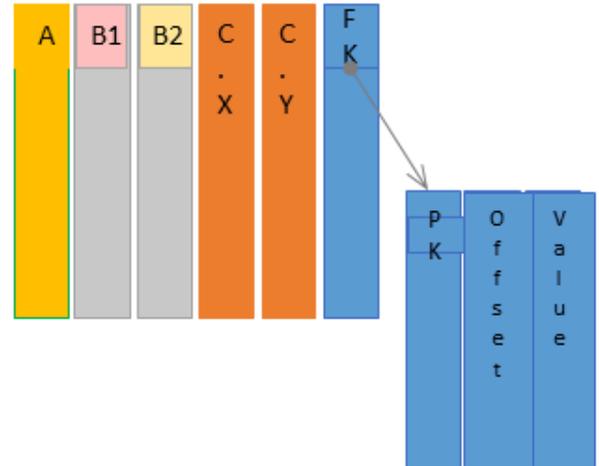
AWS Glue は、クローラを使用して半構造化データのスキーマを推測します。次に、ETL (抽出、変換、ロード) ジョブを使用してデータをリレーショナルスキーマに変換します。例えば、JSON データを解析して、Amazon Simple Storage Service (Amazon S3) ソースファイルから Amazon Relational Database Service (Amazon RDS) テーブルに変換したい場合があります。AWS Glue がスキーマの違いを処理する方法を理解すると、変換プロセスを理解するうえで役立ちます。

この図は、AWS Glue が半構造化スキーマをリレーショナルスキーマに変換する方法を示しています。

Semi-structured schema



Relational schema



この図表は、以下を示すものです：

- 単一の値 A は、直接リレーショナル列に変換されます。
- 値のペアである B1 および B2 は、2 つのリレーショナル列に変換されます。
- 子の X と Y を持つ C 構造は、2 つのリレーショナル列に変換されます。
- 配列 D[] は、別のリレーショナルテーブルを指す外部キー (FK) 列のリレーショナル列に変換されます。2 番目のリレーショナルテーブルには、プライマリキー (PK) に加えて、オフセットと配列の項目の値を含む列があります。

AWS Glue ータイプシステム

AWS Glue は、さまざまな方法でデータを保存するデータシステムに対して、複数のタイプシステムを使用して汎用性の高いインターフェースを提供しています。このドキュメントでは、AWS Glue eタイプのシステムとデータ標準を明確にしています。

AWS Glue データカタログタイプ

データカタログは、さまざまなデータシステムに格納されているテーブルとフィールドのレジストリであり、メタストアです。AWS Glue クローラーや AWS Glue with Spark ジョブなどの Glue コンポーネントがデータカタログに書き込む場合、フィールドの型を追跡するための内部型システムを使用して書き込みを行います。これらの値は、AWS Glue コンソールのテーブルスキーマのデータタイプ列に表示されます。このタイプのシステムは、Apache Hive の型のシステムに基づいています。Apache Hive の型のシステムの詳細については、Apache Hive wiki の「[型](#)」を参照してください。特定のタイプとサポートの詳細については、スキーマビルダーの一部として AWS Glue Console に例が用意されています。

検証、互換性、その他の用途

データカタログは、型フィールドに書き込まれた型を検証しません。AWS Glue コンポーネントがデータカタログの読み取りと書き込みを行うと、相互に互換性が保たれます。AWS Glue コンポーネントは、ハイブタイプとの高い互換性を維持することも目的としています。ただし、AWS Glue コンポーネントはすべての Hive タイプとの互換性を保証するものではありません。これにより、データカタログ内のテーブルを操作するときに Athena DDL などのツールとの相互運用が可能になります。

データカタログは型を検証しないため、他のサービスは、Hive の型のシステムに厳密に準拠するシステムやその他のシステムを使用することで、型を追跡するためにデータカタログを利用できます。

AWS Glue with Spark スクリプトのタイプ

AWS Glue with Spark スクリプトがデータセットを解釈または変換する場合 DynamicFrame、スクリプトで使用されているデータセットのメモリ内表現を提供します。DynamicFrame の目的は、Spark DataFrame と似ています。つまり、Spark がデータに対して変換をスケジュールして実行できるようにデータセットをモデル化します。toDF および fromDF メソッドを提供することにより、DynamicFrame の型表現が DataFrame と相互互換であることを保証します。

型情報を DataFrame に推測または提供できる場合は、特に文書化されていない限り、DynamicFrame に推測または提供できます。特定のデータ形式向けに最適化されたリーダーま

たはライターを提供するとき、Spark がお客様のデータを読み書きできる場合は、文書に記載されている制限に従って、提供されているリーダーおよびライターも読み取りまたは書き込みを行うことができます。リーダーおよびライターの詳細については、「[the section called “データ形式に関するオプション”](#)」を参照してください。

Choice 型

DynamicFrames は、ディスク上の行ごとの値が一貫性のない型を持つ可能性があるデータセットのフィールドをモデル化するメカニズムを提供します。例えば、フィールドで特定の行に文字列として格納された数値が、他の行では整数として格納されている場合があります。このメカニズムは Choice と呼ばれるインメモリ型です。Choice 列を具体的な型に変換するための、ResolveChoiceメソッドなどの変換機能を提供しています。AWS Glue ETLは、通常の実操作ではChoiceタイプをData Catalogに書き込みません。選択肢タイプは、DynamicFrame データセットのメモリモデルのコンテキストにのみ存在します。Choice 型の使用例については、「[the section called “データ準備サンプル”](#)」を参照してください。

AWS Glue ークローラータイプ

クローラーは、データセット用の一貫性のある使いやすいスキーマを作成し、それをデータカタログに保存して、他の AWS Glue コンポーネントや Athena で使用できるようにすることを目的としています。クローラーは、データカタログの前のセクション「[the section called “AWS Glue データカタログタイプ”](#)」で説明したように、型を扱います。列に 2 つ以上の型の値が含まれる「Choice」型のシナリオで使用可能な型を作成するために、クローラーは潜在的な型をモデル化した struct 型を作成します。

AWS Glue の開始方法

以下のセクションでは、AWS Glue の設定について説明します。AWS Glue の使用を開始するにあたり、すべての設定事項が必要になるわけではありません。IAM アクセス許可、暗号化、VPC を使用する場合は DNS、データストアにアクセスするための環境、インタラクティブセッションを使用する場合など、必要に応じて指示に従って設定することができます。

トピック

- [AWS Glue クローン作成の概要](#)
- [AWS Glue 用の IAM アクセス許可のセットアップ](#)
- [AWS Glue 使用状況プロファイルの設定](#)
- [AWS Glue Data Catalog の開始方法](#)
- [データストアへのネットワークアクセスを設定する](#)
- [AWS Glue での暗号化のセットアップ](#)
- [AWS Glue のための開発用ネットワークの設定](#)

AWS Glue クローン作成の概要

AWS Glue では、AWS Glue Data Catalog にメタデータを格納します。このメタデータを使用して、データソースを変換してデータウェアハウスやデータレイクをロードする ETL ジョブを調整します。以下の手順では、一般的なワークフローと、AWS Glue を使用して作業する際に行う選択肢のいくつかについて説明します。

Note

以下の手順に従うことも、手順 1 ~ 3 を自動的に実行するワークフローを作成することもできます。詳しくは、「[the section called “設計図とワークフローを使用した複雑な ETL アクティビティの実行”](#)」を参照してください。

1. AWS Glue Data Catalog にテーブル定義を入力します。

永続データストアの場合は、コンソールでクローラを追加して AWS Glue Data Catalog にデータを追加できます。[Add crawler (クローラを追加)] ウィザードは、テーブルのリストまたはクローラのリストから開始できます。クローラがアクセスするための 1 つ以上のデータストアを選択

します。スケジュールを作成して、クローラの実行頻度を決定することもできます。データストリームの場合は、テーブル定義を手動で作成して、ストリームプロパティを定義できます。

必要に応じて、データのスキーマを推測するカスタム分類子を提供できます。grok パターンを使用してカスタム分類子を作成できます。ただし、AWS Glue には、カスタム分類子がデータを認識しない場合にクローラによって自動的に使用される組み込み分類子が用意されています。クローラを定義する時に、分類子を選択する必要はありません。AWS Glue の分類子の詳細については、「[AWS Glue でのクローラーへの分類子の追加](#)」を参照してください。

一部のタイプのデータストアをクローラするには、認証とロケーション情報を提供する接続が必要です。必要に応じて、AWS Glue コンソールでこの必要な情報を提供する接続を作成できます。

クローラはデータストアを読み取り、データ定義と名前付きテーブルを AWS Glue Data Catalog に作成します。これらのテーブルは、選択したデータベースに整理されます。また、手動で作成したテーブルを Data Catalog に入力することもできます。この方法では、スキーマおよびその他のメタデータを提供して、Data Catalog にテーブル定義を作成します。この方法は少し面倒でエラーが発生しやすいいため、より良い方法として、クローラにテーブル定義を作成させることができます。

AWS Glue Data Catalog にテーブル定義を入力する方法の詳細については、「[テーブルの作成](#)」を参照してください。

2. ソースからターゲットへのデータの変換を記述するジョブを定義します。

一般に、ジョブを作成するには、次の選択をする必要があります。

- ジョブのソースとなるテーブルを AWS Glue Data Catalog から選択します。ジョブでは、このテーブル定義を使ってデータソースにアクセスし、データの型式を解釈します。
- ジョブのターゲットとなるテーブルまたは場所を AWS Glue Data Catalog から選択します。ジョブはこの情報を使用して、データストアにアクセスします。
- ソースをターゲットに変換するスクリプトを生成するように AWS Glue に指示します。AWS Glue は、ソーススキーマからターゲットスキーマ形式にデータを変換する組み込み変換を呼び出すコードを生成します。これらの変換は、データのコピー、列の名前の変更、データのフィルタリングなどの操作を実行し、必要に応じてデータを変換します。このスクリプトは、AWS Glue コンソールで変更できます。

AWS Glue でジョブを定義する方法の詳細については、「[AWS Glue Studio でビジュアル ETL ジョブを作成する](#)」を参照してください。

3. ジョブを実行してデータを変換します。

オンデマンドでジョブを実行するか、次のいずれかのトリガータイプに基づいてジョブを開始することができます。

- cron スケジュールに基づいたトリガー。
- イベントベースのトリガー。たとえば、別のジョブが正常に完了すると、AWS Glue ジョブを開始できます。
- オンデマンドでジョブを開始するトリガー。

AWS Glue のトリガーの詳細については、「[トリガーを使用したジョブとクローラの開始](#)」を参照してください。

4. スケジュールされたクローラとトリガーされたジョブをモニタリングします。

AWS Glue コンソールを使用して以下を表示します。

- ジョブの実行の詳細とエラー。
- クローラは詳細とエラーを実行します。
- AWS Glue アクティビティに関する通知

AWS Glue でクローラとジョブをモニタリングする方法の詳細については、「[AWS Glue のモニタリング](#)」を参照してください。

AWS Glue 用の IAM アクセス許可のセットアップ

このトピックの手順は、AWS Glue の AWS Identity and Access Management (IAM) アクセス許可をすばやくセットアップするのに役立ちます。以下のタスクを実行します。

- IAM ID に AWS Glue リソースへのアクセスを許可します。
- ジョブの実行、データへのアクセス、AWS Glue Data Quality タスクの実行のためのサービスロールを作成します。

AWS Glue の IAM アクセス許可をカスタマイズするために使用できる詳細な手順については、「[AWS Glue の IAM アクセス許可の設定](#)」を参照してください。

AWS Management Console で AWS Glue の IAM アクセス許可をセットアップするには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。

2. [開始方法] を選択します。
3. [AWS Glue のアカウントの準備] で、[IAM アクセス許可のセットアップ] を選択します。
4. AWS Glue アクセス許可を付与する IAM ID (ロールまたはユーザー) を選択します。AWS Glue は [AWSGlueConsoleFullAccess](#) マネージドポリシーをこれらの ID にアタッチします。これらのアクセス許可を手動で設定する場合、またはデフォルトのサービスロールを設定するだけでよい場合は、このステップをスキップできます。
5. [Next] を選択します。
6. ロールとユーザーが必要とする Amazon S3 アクセス権のレベルを選択します。このステップで選択したオプションは、選択したすべての ID に適用されます。
 - a. [S3 ロケーションの選択] で、アクセスを許可する Amazon S3 ロケーションを選択します。
 - b. 次に、上記で選択したロケーションに対して各 ID に付与するアクセス権を [読み取り専用] (推奨) にするか [読み取りおよび書き込み] にするかを選択します。AWS Glue は、選択した読み取りまたは書き込みアクセス許可とロケーションとの組み合わせに基づいて、アクセス許可ポリシーを ID に追加します。

次の表は、AWS Glue が Amazon S3 アクセスに付与するアクセス許可を示しています。

選択内容	AWS Glue がアタッチする内容
変更なし	アクセス許可は付与されません。AWS Glue は ID のアクセス許可を何も変更しません。
特定の Amazon S3 ロケーションへのアクセスを許可する (読み取り専用)	<p>選択した IAM ID に埋め込まれたインラインポリシー。詳細については、「IAM ユーザーガイド」の「インラインポリシー」を参照してください。</p> <p>AWS Glue は、次の規則を使用してポリシーに名前を付けます: AWSGlueConsole <i><Role/Use r></i> InlinePolicy-read-specific-access- <i><UUID></i>。例: AWSGlueConsoleRoleInlinePol</p>

選択内容	AWS Glue がアタッチする内容
	<p data-bbox="912 214 1432 298">icy-read-specific-access-123456780123 。</p> <p data-bbox="912 340 1490 520">以下は、指定した Amazon S3 ロケーションへの読み取り専用アクセスを許可するように AWS Glue がアタッチするインラインポリシーの例です。</p> <pre data-bbox="912 550 1507 1234">{ "Version": "2012-10-17", "Statement": [{ "Effect": "Allow", "Action": ["s3:Get*", "s3:List*"], "Resource": ["arn:aws:s3:::DOC-EXAMPLE-BUCKET/*"] }] }</pre>

選択内容	AWS Glue がアタッチする内容
<p>特定の Amazon S3 ロケーションへのアクセスを許可する (読み取りと書き込み)</p>	<p>選択した IAM ID に埋め込まれたインラインポリシー。詳細については、「IAM ユーザーガイド」の「インラインポリシー」を参照してください。</p> <p>AWS Glue は、次の規則を使用してポリシーに名前を付けます: <code>AWSGlueConsole <Role/User> InlinePolicy-read -and-write-specific-access- <UUID></code>。例: <code>AWSGlueConsoleRoleInlinePolicy-read-and-write-specific-access-123456780123</code>。</p> <p>以下は、指定した Amazon S3 ロケーションへの読み取りと書き込みアクセスを許可するように AWS Glue がアタッチするインラインポリシーの例です。</p> <pre data-bbox="917 1081 1502 1795"> { "Version": "2012-10-17", "Statement": [{ "Effect": "Allow", "Action": ["s3:Get*", "s3:List*", "s3:*Object*"], "Resource": ["arn:aws:s3:::DOC-EXAMPLE-BUCKET1/*", "arn:aws:s3:::DOC-EXAMPLE-BUCKET2/*"] }] } </pre>

選択内容	AWS Glue がアタッチする内容
	}
すべての Amazon S3 へのアクセスを許可する (読み取り専用)	AmazonS3ReadOnlyAccess マネージド IAM ポリシー。詳細については、「 AWS マネージドポリシー: AmazonS3ReadOnlyAccess 」を参照してください。
すべての Amazon S3 へのアクセスを許可する (読み取りと書き込み)	AmazonS3FullAccess マネージド IAM ポリシー。詳細については、「 AWS マネージドポリシー: AmazonS3FullAccess 」を参照してください。

7. [Next] を選択します。

8. アカウントのデフォルトの AWS Glue サービスロールを選択します。サービスロールは、AWS Glue がユーザーに代わって他の AWS のサービスのリソースにアクセスするのに使用する IAM ロールです。詳細については、「[AWS Glue のサービスロール](#)」を参照してください。

- 標準の AWS Glue サービスロールを選択すると、AWS Glue は AWS アカウント に新しい IAM ロールを `AWSGlueServiceRole` という名前で作成して、次のマネージドポリシーをアタッチします。アカウントに既に `AWSGlueServiceRole` という名前の IAM ロールがある場合は、AWS Glue はこれらのポリシーをその既存のロールにアタッチします。

- [AWSGlueServiceRole](#)
- [AmazonS3FullAccess](#)

- 既存の IAM ロールを選択すると、AWS Glue はそのロールをデフォルトとして設定しますが、アクセス許可は何も追加しません。ロールが AWS Glue のサービスロールとして使用されるように設定したことを確認してください。詳細については、[ステップ 1: AWS Glue サービスの IAM ポリシーを作成する](#) および [ステップ 2: AWS Glue 用の IAM ロールを作成する](#) を参照してください。

9. [Next] を選択します。

10. 最後に、選択したアクセス許可を確認し、[変更を適用] を選択します。変更を適用すると、AWS Glue は選択した ID に IAM アクセス許可を追加します。IAM コンソール (<https://console.aws.amazon.com/iam/>) で新しいアクセス許可の表示や変更ができます。

これで、AWS Glue の最小限の IAM アクセス許可のセットアップが完了しました。実稼働環境では、「[AWS Glue のセキュリティ](#)」と「[AWS Glue の Identity and Access Management](#)」をよく理解し、ユースケースに対して AWS リソースを確保できるようにすることをお勧めします。

次のステップ

IAM アクセス許可の設定が完了したので、次のトピックを確認して AWS Glue の使用を開始できます。

- [「AWS スキルビルダー」の「AWS Glue の開始方法」](#)
- [AWS Glue Data Catalog の開始方法](#)

AWS Glue Studio の設定

ビジュアル ETL 用に AWS Glue を初めて使用する場合は、このセクションのタスクを完了してください。

トピック

- [AWS Glue Studio ユーザーに必要な IAM アクセス許可を確認します](#)
- [ETL ジョブに必要な IAM アクセス許可を確認する](#)
- [AWS Glue Studio に対する IAM のアクセス許可の設定](#)
- [ETL ジョブの VPC を設定します](#)

AWS Glue Studio ユーザーに必要な IAM アクセス許可を確認します

AWS Glue Studio を使用するには、ユーザーはさまざまな AWS リソースにアクセスする必要があります。ユーザーは、Amazon S3 バケット、IAM ポリシーとロール、および AWS Glue Data Catalog オブジェクト。

AWS Glue サービスのアクセス許可

AWS Glue Studio は、AWS Glue サービスのアクションとリソースを使用します。AWS Glue Studio を効果的に使用するには、ユーザーはこれらのアクションやリソースに対するアクセス許可が必要です。AWS Glue Studio ユーザーに `AWSGlueConsoleFullAccess` マネージドポリシーを付与するか、より少ないアクセス許可の組み合わせでカスタムポリシーを作成できます。

⚠ Important

セキュリティのベストプラクティスに従って、ポリシーを強化して、Amazon S3 バケットおよび Amazon CloudWatch ロググループへのアクセスをさらに制限することが推奨されます。Amazon S3 ポリシーの例については、「[Writing IAM Policies: How to Grant Access to an Amazon S3 Bucket](#)」を参照してください。

AWS Glue Studio のためのカスタム IAM ポリシーの作成

AWS Glue Studio では、カスタムポリシーを作成するためのアクセス許可セットを、より小規模にできます。カスタムポリシーは、オブジェクトまたはアクションのサブセットに対しアクセス許可を付与できます。カスタムポリシーを作成する際には、以下の情報を参考にしてください。

AWS Glue Studio API を使用するには、IAM アクセス許可のアクションポリシーに `glue:UseGlueStudio` を含めます。`glue:UseGlueStudio` を使用することで、時間の経過とともに API にさらにアクションが追加されても、すべての AWS Glue Studio アクションにアクセスできます。

有向非巡回 (DAG) に関するアクション

- CreateDag
- UpdateDag
- GetDag
- DeleteDag

ジョブに関するアクション

- SaveJob
- GetJob
- CreateJob
- DeleteJob
- GetJobs
- UpdateJob

ジョブ実行に関するアクション

- StartJobRun
- GetJobRuns
- BatchStopJobRun
- GetJobRun
- QueryJobRuns
- QueryJobs
- QueryJobRunsAggregated

スキーマに関するアクション

- GetSchema
- GetInferredSchema

データベースに関するアクション

- GetDatabases

プランに関するアクション

- GetPlan

テーブルに関するアクション

- SearchTables
- GetTables
- GetTable

接続に関するアクション

- CreateConnection
- DeleteConnection
- UpdateConnection
- GetConnections
- GetConnection

マッピングに関するアクション

- GetMapping

S3 プロキシに関するアクション

- ListBuckets
- ListObjectsV2
- GetBucketLocation

セキュリティ設定に関するアクション

- GetSecurityConfigurations

スクリプトに関するアクション

- CreateScript (AWS Glue での同じ名前の API とは異なります)

AWS Glue Studio API へのアクセス

AWS Glue Studio にアクセスするには、IAM アクセス許可のアクションポリシーリストに `glue:UseGlueStudio` を追加します。

以下の例では、`glue:UseGlueStudio` はアクションポリシーに含まれますが、AWS Glue Studio API は個別に識別されません。これは、`glue:UseGlueStudio` を含めることで、IAM アクセス許可で個別に AWS Glue Studio API を指定する必要なしに、内部 API へのアクセス許可が自動的に付与されるためです。

この例では、リストされた追加のアクションポリシー (`glue:SearchTables` など) が AWS Glue Studio API ではないため、場合によっては IAM アクセス許可に含める必要があります。Amazon S3 プロキシアクションを含めて、付与する Amazon S3 アクセス許可のレベルを指定することもできます。以下のポリシーの例では、選択した IAM ロールに十分なアクセス許可がある場合に、AWS Glue Studio を開く、ビジュアルジョブを作成する、ビジュアルジョブを保存/実行するためのアクセス許可が付与されます。

```
{
  "Version": "2012-10-17",
  "Statement": [
```

```
{
  "Sid": "VisualEditor0",
  "Effect": "Allow",
  "Action": [
    "glue:UseGlueStudio",
    "iam:ListRoles",
    "iam:ListUsers",
    "iam:ListGroups",
    "iam:ListRolePolicies",
    "iam:GetRole",
    "iam:GetRolePolicy",
    "glue:SearchTables",
    "glue:GetConnections",
    "glue:GetJobs",
    "glue:GetTables",
    "glue:BatchStopJobRun",
    "glue:GetSecurityConfigurations",
    "glue>DeleteJob",
    "glue:GetDatabases",
    "glue>CreateConnection",
    "glue:GetSchema",
    "glue:GetTable",
    "glue:GetMapping",
    "glue>CreateJob",
    "glue>DeleteConnection",
    "glue>CreateScript",
    "glue:UpdateConnection",
    "glue:GetConnection",
    "glue:StartJobRun",
    "glue:GetJobRun",
    "glue:UpdateJob",
    "glue:GetPlan",
    "glue:GetJobRuns",
    "glue:GetTags",
    "glue:GetJob",
    "glue:QueryJobRuns",
    "glue:QueryJobs",
    "glue:QueryJobRunsAggregated"
  ],
  "Resource": "*"
},
{
  "Action": [
    "iam:PassRole"
  ]
}
```

```
    ],
    "Effect": "Allow",
    "Resource": "arn:aws:iam::*:role/AWSGlueServiceRole*",
    "Condition": {
      "StringLike": {
        "iam:PassedToService": [
          "glue.amazonaws.com"
        ]
      }
    }
  }
]
```

ノートブックとデータのプレビューアクセス許可

データプレビューとノートブックを使用すると、ジョブを実行することなく、ジョブの任意の段階（読み取り、変換、書き込み）でデータのサンプルを表示できます。データにアクセスするときに使用する AWS Identity and Access Management の AWS Glue Studio (IAM) ロールを指定します。IAM ロールは想定可能であることを意図しており、標準の長期認証情報（パスワードやアクセスキーなど）を関連付けることはありません。その代わりに、AWS Glue Studio がそのロールを引き受けて、IAM は一時的なセキュリティ認証情報を提供します。

データプレビューとノートブックコマンドが正しく動作するようにするには、文字列 `AWSGlueServiceRole` で始まる名前のロールを使用します。ロールに別の名前を使用する場合は、`iam:passrole` アクセス許可を追加し、IAM のロールに対するポリシーを設定する必要があります。詳しくは、「[「AWSGlueServiceRole*」という名前ではないロールの IAM ポリシーを作成します。](#)」を参照してください。

Warning

ロールがノートブックに対する `iam:passrole` アクセス許可を与えた場合、ロールチェーンを実装すると、あるユーザーが意図せずにノートブックにアクセスする可能性があります。現在、ノートブックへアクセス許可されているユーザーをモニタリングできる監査は実装されていません。

IAM ID によるデータプレビューセッションの作成を認めない場合は、「[the section called “ID によるデータプレビューセッションの作成を拒否する”](#)」の例を参照してください。

Amazon CloudWatch のアクセス許可

AWS Glue Studio を使用して Amazon CloudWatch ジョブをモニタリングできます。これは、AWS Glue から raw データを収集して読み取り可能なほぼリアルタイムのメトリックに処理します。デフォルトでは、AWS Glue メトリクスデータは CloudWatch に自動的に送信されます。詳細については、Amazon CloudWatch ユーザーガイドの「[Amazon CloudWatch とは何ですか?](#)」およびAWS Glue デベロッパーガイドの「[AWS Glue メトリクス](#)」を参照してください。

CloudWatch ダッシュボードにアクセスするには、AWS Glue Studio では以下のいずれかが必要です。

- AdministratorAccess ポリシー
- CloudWatchFullAccess ポリシー
- これらの特定のアクセス許可の 1 つ以上を含むカスタムポリシー。
 - cloudwatch:GetDashboard および cloudwatch:ListDashboards でダッシュボードを表示する
 - cloudwatch:PutDashboard でダッシュボードを作成または変更する
 - cloudwatch>DeleteDashboards でダッシュボードを削除する

ポリシーを使用してIAMユーザーのアクセス許可を変更する方法の詳細については、IAM ユーザーガイドの[IAM ユーザーのアクセス許可の変更](#)を参照してください。

ETL ジョブに必要な IAM アクセス許可を確認する

AWS Glue Studio を使用してジョブを作成する場合、ジョブは作成時に指定する IAM ロールのアクセス許可を引き継ぎます。この IAM ロールには、データソースからデータを抽出し、ターゲットにデータを書き込み、AWS Glue リソースにアクセスするためのアクセス許可が必要です。

AWSGlueServiceRole で正しく使用するには、ジョブ用に作成するロールの名前が「AWS Glue Studio」の文字列で始まる必要があります。例えば、ロール AWSGlueServiceRole-FlightDataJob に名前を付けることができます。

データソースとデータターゲットのアクセス許可

AWS Glue Studio ジョブでは使用するすべてのソース、ターゲット、スクリプト、および一時ディレクトリに対して Amazon S3 にアクセスできる必要があります。特定の Amazon S3 リソースにきめ細かいアクセス許可を付与するポリシーを作成できます。

- データソースには、`s3:ListBucket` および `s3:GetObject` アクセス許可が必要です。
- データターゲットには、`s3:ListBucket`、`s3:PutObject`、および `s3>DeleteObject` アクセス許可が必要です。

データソースとして Amazon Redshift を選択した場合は、クラスターのアクセス許可のロールを指定できます。Amazon Redshift クラスターに対して実行されるジョブが、一時的な認証情報を使用して一時的なストレージとして Amazon S3 にアクセスするコマンドを発行します。ジョブが 1 時間を超えて実行されると、これらの認証情報が期限切れになり、ジョブが失敗します。この問題を回避するために、一時的な認証情報を使用してジョブに必要なアクセス許可を付与するロールを Amazon Redshift クラスター自体に割り当てることができます。詳細については、AWS Glue データベースデベロッパーガイドから「[Moving Data to and from Amazon Redshift](#)」を参照してください。

ジョブで Amazon S3 以外のデータソースまたはターゲットを使用する場合は、これらのデータソースおよびターゲットにアクセスするために、ジョブで使用される IAM ロールに必要なアクセス許可をアタッチする必要があります。詳細については、AWS Glue デベロッパーガイドの「[Setting Up Your Environment to Access Data Stores](#)」を参照してください。

データストアにコネクタと接続を使用している場合は、追加のアクセス許可が必要です。詳細については、「[the section called “コネクタの使用に必要なアクセス許可”](#)」。

ジョブの削除に必要なアクセス許可

AWS Glue Studio では、コンソールで複数のジョブを選択して削除できます。このアクションを実行するには、`glue:BatchDeleteJob` アクセス許可が必要です。これは、ジョブを削除するために `glue>DeleteJob` のアクセス許可を必要とする AWS Glue コンソールとは異なります。

AWS Key Management Service アクセス許可

AWS Key Management Service (AWS KMS) でサーバー側の暗号化を使用する Amazon S3 ソースとターゲットにアクセスする場合は、ジョブが使用する AWS Glue Studio ロールにポリシーをアタッチして、ジョブがデータを復号化できるようにします。ジョブロールには、`kms:ReEncrypt`、`kms:GenerateDataKey`、および `kms:DescribeKey` のアクセス許可が必要です。さらに、ジョブロールには、AWS KMS カスタマーマスターキー (CMK) で暗号化された

Amazon S3 オブジェクトをアップロードまたはダウンロードするための `kms:Decrypt` アクセス許可が必要です。

AWS KMS CMK を使用するための追加料金がかかります。詳細については、[AWS Key Management Service デベロッパーガイド](#)の「[AWS Key Management Service の概念 - カスタマーマスターキー \(CMK\)](#)」および「[AWS Key Management Service の料金](#)」を参照してください。

コネクタの使用に必要なアクセス許可

AWS Glue カスタムコネクタと接続を使用してデータストアにアクセスしている場合、AWS Glue ETL ジョブの実行に使用されるロールには、追加のアクセス許可がアタッチされている必要があります。

- AWS Marketplace から購入したコネクタにアクセスするための AWS マネージドポリシー `AmazonEC2ContainerRegistryReadOnly`。
- `glue:GetJob` および `glue:GetJobs` アクセス許可。
- 接続で使用されるシークレットにアクセスするための AWS Secrets Manager アクセス許可。IAM ポリシーのサンプルについては、「[例: シークレット値を取得するアクセス許可](#)」を参照してください。

AWS Glue ETL ジョブが Amazon VPC を実行中の VPC 内で実行される場合、VPC は [the section called “ETL ジョブの VPC を設定します”](#) で説明されているように設定する必要があります。

AWS Glue Studio に対する IAM のアクセス許可の設定

AWS 管理者ユーザーを使用して、ロールを作成し、ユーザーとジョブロールにポリシーを割り当てることができます。

`AWSGlueConsoleFullAccess` AWS マネージドポリシーを使用すると、AWS Glue Studio コンソールを使用するために必要なパーミッションを提供することができます。

独自のポリシーを作成するには、AWS Glueデベロッパーガイド の「[Create an IAM Policy for the AWS Glue Service](#)」に従ってください。[AWS Glue Studio ユーザーに必要な IAM アクセス許可を確認します](#) で前述の IAM アクセス許可を含めます。

トピック

- [AWS Glue Studio ユーザーにポリシーを添付する。](#)
- [「AWSGlueServiceRole*」という名前ではないロールの IAM ポリシーを作成します。](#)

AWS Glue Studio ユーザーにポリシーを添付する。

AWS コンソールにサインインする AWS Glue Studio ユーザーは、特定のリソースへのアクセス権限が必要です。IAM ポリシーをユーザーに割り当てることで、これらの権限を提供します。

ユーザーに AWSGlueConsoleFullAccess マネージドポリシーをアタッチするには

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. ナビゲーションペインで、[ポリシー] を選択します。
3. ポリシーのリストで、[AWSGlueConsoleFullAccess] ポリシーの横にあるチェックボックスを選択します。[Filter (フィルター)] メニューと検索ボックスを使用して、ポリシーのリストをフィルタリングできます。
4. [Policy actions] を選択して、[Attach] を選択します。
5. ポリシーをアタッチするユーザーを選択します。[Filter] メニューと検索ボックスを使用して、プリンシパルエンティティのリストをフィルタリングできます。ポリシーをアタッチするユーザーを選択し、[ポリシーのアタッチ] を選択します。
6. 必要に応じて前の手順を繰り返し、ユーザーに追加のポリシーをアタッチします。

「AWSGlueServiceRole*」という名前ではないロールの IAM ポリシーを作成します。

AWS Glue Studio で使用されるロールの IAM ポリシーを設定する方法

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. 新しい IAM ポリシーを追加します。既存のポリシーに追加するか、新しい IAM インラインポリシーを作成できます。IAM ポリシーを作成するには
 1. [Policies] を選択してから、[Create Policy] を選択します [ご利用開始にあたって] ボタンが表示されたら、そのボタンを選択して、[ポリシーの作成] を選択します。
 2. [独自のポリシーを作成] の横で、[選択] を選択します。
 3. [ポリシー名] に、後で参照しやすい任意の値を入力します。オプションとして、[説明] に説明のテキストを入力します。
 4. [ポリシードキュメント] に、以下の形式のポリシーステートメントを入力してから、[ポリシーの作成] を選択します。

3. 次のブロックをコピーして、[構文] 配列の下ポリシーに貼り付け、*my-interactive-session-role-prefix* を、AWS Glue のアクセス許可に関連付けるすべての一般的なロールのプレフィックスに置き換えます。

```
{
  "Action": [
    "iam:PassRole"
  ],
  "Effect": "Allow",
  "Resource": "arn:aws:iam::*:role/my-interactive-session-role-prefix",
  "Condition": {
    "StringLike": {
      "iam:PassedToService": [
        "glue.amazonaws.com "
      ]
    }
  }
}
```

ポリシーに含まれる、[Version] (バージョン) 配列と [Statement] (ステートメント) 配列の完全な例を以下に示します。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Action": [
        "iam:PassRole"
      ],
      "Effect": "Allow",
      "Resource": "arn:aws:iam::*:role/my-interactive-session-role-prefix",
      "Condition": {
        "StringLike": {
          "iam:PassedToService": [
            "glue.amazonaws.com "
          ]
        }
      }
    }
  ]
}
```

4. ユーザーのポリシーを有効にするには、[ユーザー] を選択します。
5. ポリシーをアタッチする ユーザーを選択します。

ETL ジョブの VPC を設定します

Amazon Virtual Private Cloud (Amazon VPC) を使用すると、AWS クラウド 内の論理的に分離された独自の領域に仮想プライベートクラウド (VPC) と呼ばれる仮想ネットワークを定義できます。AWS のリソース (インスタンスなど) を VPC 内部で起動できます。VPC は、お客様自身のデータセンターで運用されている従来のネットワークによく似ていますが、 のスケーラブルなインフラストラクチャを使用できるというメリットがあります。AWS お客様の VPC はお客様が設定できます。IP アドレスレンジの選択、サブネットの作成、ルートテーブル、ネットワークゲートウェイ、セキュリティの設定ができます。VPC のインスタンスをインターネットに接続できます。VPC を自社のデータセンターに接続し、AWS クラウド をデータセンターの拡張部分として使用できます。各サブネットでのリソースの保護には、セキュリティグループ、ネットワークアクセスコントロールリストなど、複数のセキュリティレイヤーを使用できます。詳細については、[Amazon VPC ユーザーガイド](#)を参照してください。

コネクタを使用するとき、AWS Glue ETL ジョブをVPC 内で実行するように設定できます。必要に応じて、次のように VPC を設定する必要があります。

- AWS にはないデータストア用のパブリックネットワークアクセス。ジョブによってアクセスされるすべてのデータストアは、VPC サブネットから使用できる必要があります。
- ジョブが VPC リソースとパブリックインターネットの両方にアクセスする必要がある場合は、VPC 内にネットワークアドレス変換 (NAT) ゲートウェイが必要になります。

詳細については、AWS Glue デベロッパーガイドの「[Setting Up Your Environment to Access Data Stores](#)」を参照してください。

AWS Glue Studio 中でのノートブックの使用開始

AWS Glue Studio を通してノートブックをスタートすると、ほんの数秒後にデータを調査し、ジョブスクリプトの開発を開始できるように、すべての設定ステップが実行されます。

以下のセクションでは、AWS Glue Studio で ETL ジョブのノートブックを使用するために、ロールを作成し、適切なアクセス許可を付与する方法について説明します。

トピック

- [IAM ロールに対するアクセス許可の付与](#)

IAM ロールに対するアクセス許可の付与

ノートブックを使用するには、AWS Glue Studio のセットアップが事前に必要です。

AWS Glue でノートブックを使用するには、ロールに次のものがが必要です。

- `sts:AssumeRole` アクションに対する AWS Glue との信頼関係。タグ付けが必要な場合は `sts:TagSession`。
- ノートブック、AWS Glue、インタラクティブセッションのすべての API オペレーションを含む IAM ポリシー。
- PassRole の IAM ポリシー (ロールがノートブックからインタラクティブセッションに渡されるようにする必要があるため)。

例えば、新しいロールを作成するときに、`AWSGlueConsoleFullAccessRole` などの標準の AWS 管理ポリシーをロールに追加して、ノートブックのオペレーション用と IAM PassRole ポリシー用にそれぞれ新しいポリシーを追加できます。

AWS Glue との信頼関係に必要なアクション

ノートブックセッションを開始する際に、ノートブックに渡されるロールの信頼関係に `sts:AssumeRole` を追加する必要があります。セッションにタグが含まれている場合は、`sts:TagSession` アクションも渡す必要があります。これらのアクションがないと、ノートブックセッションを開始できません。

例:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "glue.amazonaws.com"
      },
      "Action": "sts:AssumeRole"
    }
  ]
}
```

ノートブックの API オペレーションを含むポリシー

次のサンプルポリシーで、ノートブックに必要な AWS IAM 権限について説明します。新しいロールを作成する場合は、以下を含むポリシーを作成します。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "glue:StartNotebook",
        "glue:TerminateNotebook",
        "glue:GlueNotebookRefreshCredentials",
        "glue:DeregisterDataPreview",
        "glue:GetNotebookInstanceStatus",
        "glue:GlueNotebookAuthorize"
      ],
      "Resource": "*"
    }
  ]
}
```

次の IAM ポリシーを使用して、特定のリソースへのアクセスを許可できます。

- `AwsGlueSessionUserRestrictedNotebookServiceRole`: セッションを除くすべての AWS Glue リソースへのフルアクセス権限を提供します。ユーザーが、ユーザーに関連付けられているノートブックセッションのみを作成して使用できるようにします。このポリシーには、AWS Glue が他の AWS サービスで AWS Glue リソースを管理するために必要な他のアクセス許可も含まれています。
- `AwsGlueSessionUserRestrictedNotebookPolicy`: ユーザー自身に関連付けられているノートブックセッションのみを作成および使用できるようにするアクセス許可を提供します。このポリシーには、ユーザーが制限付き AWS Glue セッションロールを渡すことを明示的に許可するアクセス許可も含まれます。

ロールを渡すための IAM ポリシー

ロールを使用してノートブックを作成すると、そのロールがインタラクティブセッションに渡され、その両方で同じロールを使用できるようになります。このように、`iam:PassRole` のアクセス許可がロールのポリシーの一部として必要です。

次の例を使用して、ロール用の新しいポリシーを作成します。アカウント番号とロール名を自分のものに置き換えます。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": "iam:PassRole",
      "Resource": "arn:aws:iam::090000000210:role/<role_name>"
    }
  ]
}
```

AWS Glue 使用状況プロファイルの設定

クラウドプラットフォームを使用する主な利点の1つは、その柔軟性です。ただし、コンピューティングリソースの作成が容易であるため、管理対象外のままにしてガードレールを使用しないと、クラウドコストが急増するリスクがあります。その結果、管理者は、インフラストラクチャコストの増加を回避しながら、ユーザーが不要な摩擦なしで作業できるようにバランスを取る必要があります。

AWS Glue 使用プロファイルを使用すると、管理者はデベロッパー、テスター、製品チームなど、アカウント内のさまざまなクラスのユーザーに対して異なるプロファイルを作成できます。各プロファイルは、さまざまなタイプのユーザーに割り当てることができる一意のパラメータセットです。例えば、開発者はより多くのワーカーを必要とし、最大ワーカー数を増やすことができますが、製品チームはより少ないワーカーと低いタイムアウトまたはアイドルタイムアウト値を必要とする場合があります。

ジョブとジョブ実行の動作の例

プロファイル A を持つユーザー A によってジョブが作成されるとします。ジョブは特定のパラメータ値で保存されます。プロファイル B を持つユーザー B はジョブの実行を試みます。

ユーザー A がジョブを作成したときに、特定の数のワーカーを設定しなかった場合、ユーザー A のプロファイルのデフォルトセットが適用され、ジョブの定義とともに保存されました。

ユーザー B がジョブを実行すると、保存された任意の値で実行されます。ユーザー B の独自のプロフィールがより制限されており、その数のワーカーでの実行が許可されていない場合、ジョブの実行は失敗します。

リソースとしての使用状況プロフィール

AWS Glue 使用状況プロフィールは、Amazon リソースネーム (ARN) によって識別されるリソースです。アクションベースおよびリソースベースの承認を含む、すべてのデフォルトの IAM (Identity and Access Management) コントロールが適用されます。管理者は、AWS Glue リソースを作成するユーザーの IAM ポリシーを更新し、プロフィールを使用するためのアクセス権を付与する必要があります。

The screenshot shows the AWS Glue console interface for managing usage profiles. The main content area is titled "Usage profiles (1/9)" and includes a search bar and a table of profiles. The table has columns for Name, Status, Description, and Created on (UTC).

Name	Status	Description	Created on (UTC)
dev-profile-1	Assigned	-	April 30, 2024, 02:19:53
dev-profile-2	Not assigned	I edited the description and default workers	April 25, 2024, 22:10:17
product-profile-1	Not assigned	-	April 30, 2024, 02:19:02
product-profile-2	Assigned	-	May 7, 2024, 20:39:18
tester-profile-1	Assigned	test description has been edited	May 7, 2024, 20:55:25
tester-profile-2	Assigned	glue testing profile	May 7, 2024, 21:20:13
test	Assigned	I edited this successfully again	April 25, 2024, 20:28:48
test profile	Not assigned	Description I edited this	April 30, 2024, 17:17:53

トピック

- [使用状況プロフィールの作成と管理](#)
- [使用状況プロフィールとジョブ](#)

使用状況プロフィールの作成と管理

AWS Glue 使用状況プロフィールの作成

管理者は使用状況プロフィールを作成し、さまざまなユーザーに割り当てます。使用状況プロフィールを作成するときは、デフォルト値と、さまざまなジョブおよびセッションパラメータに許可される値の範囲を指定します。ジョブまたはインタラクティブセッションには、少なくとも1つのパラメータを設定する必要があります。ジョブにパラメータ値が指定されていない場合に使用するデフォ

ルト値をカスタマイズしたり、このプロファイルを使用するときにユーザーがパラメータ値を指定した場合に、範囲制限または検証に許可される値のセットを設定したりできます。

デフォルトは、管理者がジョブの作成者を支援するために設定したベストプラクティスです。ユーザーが新しいジョブを作成し、タイムアウト値を設定しない場合、使用プロファイルのデフォルトのタイムアウトが適用されます。作成者にプロファイルがない場合、AWS Glue サービスのデフォルトが適用され、ジョブの定義に保存されます。実行時に、はプロファイルに設定された制限 (最小、最大、許可されたワーカー) AWS Glue を適用します。

パラメータを設定すると、他のすべてのパラメータはオプションになります。ジョブまたはインタラクティブセッション用にカスタマイズできるパラメータは次のとおりです。

- **ワーカーの数** — コンピューティングリソースの過剰な使用を避けるため、ワーカーの数を制限します。デフォルト値、最小値、最大値を設定できます。最小値は 1 です。
- **ワーカータイプ** — ワークロードに関連するワーカータイプを制限します。デフォルトのタイプを設定し、ユーザープロファイルのワーカータイプを許可できます。
- **タイムアウト** — ジョブまたはインタラクティブセッションが終了する前にリソースを実行および消費できる最大時間を定義します。長時間実行されるジョブを避けるためにタイムアウト値をセットアップします。

デフォルト値、最小値、最大値を分単位で設定できます。最小値は 1 (分) です。AWS Glue デフォルトのタイムアウトは 2,880 分ですが、使用状況プロファイルで任意のデフォルト値を設定できます。

「デフォルト」の値を設定するのがベストプラクティスです。この値は、ユーザーが値を設定しなかった場合、ジョブまたはセッションの作成に使用されます。

- **アイドルタイムアウト** — セルの実行後にタイムアウトするまでにインタラクティブセッションが非アクティブになる分数を定義します。作業の完了後に終了するインタラクティブセッションのアイドルタイムアウトを定義します。アイドルタイムアウト範囲はタイムアウトの制限内である必要があります。

デフォルト値、最小値、最大値を分単位で設定できます。最小値は 1 (分) です。AWS Glue デフォルトのタイムアウトは 2,880 分ですが、使用状況プロファイルで任意のデフォルト値を設定できます。

「デフォルト」の値を設定するのがベストプラクティスです。この値は、ユーザーが値を設定しなかった場合に、セッションの作成に使用されます。

管理者として AWS Glue 使用状況プロファイルを作成するには (コンソール)

1. 左側のナビゲーションメニューで、コスト管理 を選択します。
2. 使用プロファイルの作成 を選択します。
3. 使用プロファイルの使用プロファイル名を入力します。
4. 使用プロファイルの目的を他のユーザーが認識するのに役立つオプションの説明を入力します。
5. プロファイルに少なくとも 1 つのパラメータを定義します。フォーム内のフィールドはパラメータです。例えば、セッションアイドルタイムアウトの最小値などです。
6. 使用プロファイルに適用される任意のタグを定義します。
7. [保存] を選択します。

The screenshot shows the AWS Glue console interface for creating a usage profile. The left sidebar contains a navigation menu with categories like 'Data Catalog', 'Data Integration and ETL', and 'Legacy pages'. The main content area is titled 'Create usage profile' and includes sections for 'Name and description', 'Parameter configurations', and 'Timeout'. The 'Parameter configurations' section contains a warning box and sections for 'Number of workers' and 'Worker type'. The 'Timeout' section contains fields for default, minimum, and maximum values.

使用状況プロファイルを作成するには (AWS CLI)

1. 次のコマンドを入力します。

```
aws glue create-usage-profile --name profile-name --configuration file://  
config.json --tags list-of-tags
```

ここで、config.json はインタラクティブセッション (SessionConfiguration) とジョブ () のパラメータ値を定義できます JobConfiguration。

```
//config.json (There is a separate blob for session/job configuration
{
  "SessionConfiguration": {
    "timeout": {
      "DefaultValue": "2880",
      "MinValue": "100",
      "MaxValue": "4000"
    },
    "idleTimeout": {
      "DefaultValue": "30",
      "MinValue": "10",
      "MaxValue": "4000"
    },
    "workerType": {
      "DefaultValue": "G.2X",
      "AllowedValues": [
        "G.2X",
        "G.4X",
        "G.8X"
      ]
    },
    "numberOfWorkers": {
      "DefaultValue": "10",
      "MinValue": "1",
      "MaxValue": "10"
    }
  },
  "JobConfiguration": {
    "timeout": {
      "DefaultValue": "2880",
      "MinValue": "100",
      "MaxValue": "4000"
    },
    "workerType": {
      "DefaultValue": "G.2X",
      "AllowedValues": [
        "G.2X",
        "G.4X",
        "G.8X"
      ]
    }
  }
}
```

```
    ],
  },
  "numberOfWorkers": {
    "DefaultValue": "10",
    "MinValue": "1",
    "MaxValue": "10"
  }
}
```

2. 次のコマンドを入力して、作成された使用プロファイルを確認します。

```
aws glue get-usage-profile --name profile-name
```

レスポンス :

```
{
  "ProfileName": "foo",
  "Configuration": {
    "SessionConfiguration": {
      "numberOfWorkers": {
        "DefaultValue": "10",
        "MinValue": "1",
        "MaxValue": "10"
      },
      "workerType": {
        "DefaultValue": "G.2X",
        "AllowedValues": [
          "G.2X",
          "G.4X",
          "G.8X"
        ]
      },
      "timeout": {
        "DefaultValue": "2880",
        "MinValue": "100",
        "MaxValue": "4000"
      },
      "idleTimeout": {
        "DefaultValue": "30",
        "MinValue": "10",
        "MaxValue": "4000"
      }
    }
  }
}
```

```
    },
    "JobConfiguration": {
      "numberOfWorkers": {
        "DefaultValue": "10",
        "MinValue": "1",
        "MaxValue": "10"
      },
      "workerType": {
        "DefaultValue": "G.2X",
        "AllowedValues": [
          "G.2X",
          "G.4X",
          "G.8X"
        ]
      },
      "timeout": {
        "DefaultValue": "2880",
        "MinValue": "100",
        "MaxValue": "4000"
      }
    },
    "CreatedOn": "2024-01-19T23:15:24.542000+00:00"
  }
}
```

使用状況プロファイルの管理に使用される追加の CLI コマンド :

- `aws Glue list-usage-profiles`
- `aws glue update-usage-profile --name profile-name --configuration file://config.json`
- `aws glue delete-usage-profile --name profile-name`

使用状況プロファイルの編集

管理者は、作成した使用状況プロファイルを編集して、ジョブとインタラクティブセッションのプロファイルパラメータ値を変更できます。

使用状況プロファイルを編集するには :

管理者として AWS Glue 使用状況プロファイルを編集するには (コンソール)

1. 左側のナビゲーションメニューで、**コスト管理** を選択します。

2. 編集する権限を持つ使用プロファイルを選択し、編集を選択します。
3. 必要に応じてプロファイルを変更します。デフォルトでは、すでに値を持つパラメータが展開されます。
4. 編集の保存 を選択します。

aws Services Search for services, features, blogs, docs, and more [Alt+S] N. Virginia MyRole/AWSUser @ 0123-4567-8901

AWS Glue > Usage profiles > dev-profile-1 > Edit

Edit dev-profile-1

Name and description

Usage profile name

Usage profile description - optional

Write any details that will help you or others recognize the purpose of this configuration.

Descriptions can be up to 2048 characters long.

▼ Parameter configurations for jobs Info

Configure usage restrictions for AWS Glue jobs. Each parameter has a default value preconfigured for different types of jobs.

▼ Number of workers

The number of workers of a defined worker_type that are allocated. Customize number of workers to avoid excessive use of compute resources.

Default	Minimum	Maximum
<input type="text" value="10"/>	<input type="text" value="1"/>	<input type="text" value="20"/>

Between minimum and maximum Minimum allowed value: 1

▼ Worker type

The type of a unit capable of performing operational processes dictated by its fleet management system. Select the relevant worker types for your wo

Default worker type

Allowed worker types

Choose one or more worker types

▶ Timeout

The maximum time in minutes that a job run can consume resources before it is terminated and. Setup timeout values to avoid long running jobs.

▼ Parmeter configurations for sessions Info

Configure usage restrictions for AWS Glue interactive sessions. Each parameter has a default value preconfigured for different types of interactive sessions.

▶ Number of workers

The number of workers of a defined worker_type that are allocated. Customize number of workers to avoid excessive use of compute resources.

▶ Worker type

The type of a unit capable of performing operational processes dictated by its fleet management system. Select the relevant worker types for your workloads.

▼ Idle timeout

The number of minutes of inactivity after which an interactive session will timeout after a cell has been executed. Define idle-timeout for sessions to terminate after the work completed.

Default (minutes)	Minimum (minutes)	Maximum (minutes)
<input type="text" value="2880"/>	<input type="text" value="1"/>	<input type="text" value="4000"/>

Between minimum and maximum Minimum allowed value: 1

▶ Timeout

The maximum time in minutes that an interactive session run can consume resources before it is terminated. Setup timeout values to avoid long running sessions.

▶ Tags - optional

Tags are user-defined key-value pairs that provide metadata to organize and classify your AWS resources.

CloudShell Feedback Language © 2023, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

使用状況プロファイルを編集するには (AWS CLI)

- 次のコマンドを入力します。上記の create コマンドと同じ `--configuration` ファイル構文が使用されます。

```
aws glue update-usage-profile --name profile-name --configuration file://  
config.json
```

ここで、`config.json` はインタラクティブセッション (SessionConfiguration) とジョブ () のパラメータ値を定義します JobConfiguration。

使用プロファイルの割り当て

使用状況プロファイルページの使用状況ステータス列には、使用状況プロファイルがユーザーに割り当てられているかどうかが表示されます。ステータスにカーソルを合わせると、割り当てられた IAM エンティティが表示されます。

管理者は、AWS Glue リソースを作成するユーザー/ロールに AWS Glue 使用状況プロファイルを割り当てることができます。プロファイルの割り当ては、次の 2 つのアクションの組み合わせです。

- IAM ユーザー/ロールタグを `glue:UsageProfile` キーで更新してから、
- ユーザー/ロールの IAM ポリシーの更新。

AWS Glue Studio を使用してジョブ/インタラクティブセッションを作成するユーザーの場合、管理者は次のロールにタグを付けます。

- ジョブの制限については、管理者はログインしたコンソールロールにタグを付けます。
- インタラクティブセッションの制限については、管理者がノートブックの作成時にユーザーが提供するロールにタグを付けます。

AWS Glue リソースを作成する IAM ユーザー/ロールで管理者が更新する必要があるポリシーの例を次に示します。

```
{  
  "Effect": "Allow",  
  "Action": [  
    "glue:GetUsageProfile"  
  ],  
}
```

```
"Resource": [  
  "arn:aws:glue:us-east-1:123456789012:usageProfile/foo"  
]  
}
```

AWS Glue は、AWS Glue 使用プロファイルで指定された値に基づいてジョブ、ジョブ実行、およびセッションリクエストを検証し、リクエストが許可されていない場合は例外を発生させます。同期 APIs の場合、エラーがユーザーにスローされます。非同期パスの場合、失敗したジョブ実行は、入力パラメータがユーザー/ロールに割り当てられたプロファイルの許容範囲外であるというエラーメッセージとともに作成されます。

使用状況プロファイルをユーザー/ロールに割り当てるには：

1. (Identity and Access Management) IAM コンソールを開きます。
2. 左側のナビゲーションで、ユーザー または ロール を選択します。
3. ユーザーまたはロールを選択します。
4. [タグ] タブを選択します。
5. [新しいタグを追加] をクリックします。
6. キー `glue:UsageProfile` と使用プロファイルの名前の値を含むタグを追加します。
7. [Save changes] (変更の保存) を選択します。

The screenshot shows the AWS IAM console interface for the `AWSGlueServiceRole` role. The `Tags (1)` tab is selected, displaying a table with one tag. The tag key is `glue:UsageProfile` and the value is `foo`. The console also shows summary information for the role, including its creation date (June 28, 2023, 12:35 UTC-07:00), last activity (27 days ago), ARN (`arn:aws:iam:::role/service-role/AWSGlueServiceRole`), and maximum session duration (1 hour).

Key	Value
<code>glue:UsageProfile</code>	<code>foo</code>

割り当てられた使用状況プロファイルの表示

ユーザーは、割り当てられた使用状況プロファイルを表示し、API コールを実行して AWS Glue ジョブとセッションリソースを作成するとき、またはジョブを開始するとき使用できます。

プロファイルのアクセス許可は IAM ポリシーで提供されます。発信者ポリシーに `アクセスglue:UsageProfile` 許可がある限り、ユーザーはプロファイルを表示できます。そうしないと、アクセス拒否エラーが発生します。

割り当てられた使用状況プロファイルを表示するには：

1. 左側のナビゲーションメニューで、コスト管理 を選択します。
2. 表示する権限を持つ使用プロファイルを選択します。

Usage profile "dev-provile-1" successfully updated. Usage profile "dev-provile-1" successfully updated. To assign it to IAM roles or users, go to AWS IAM service through the "Open AWS IAM" button and tag the IAM role or user with key: glue:UsageProfile and value: dev-profile-1.

[Open AWS IAM](#)

AWS Glue > Usage profiles > dev-profile-1

dev-profile-1

[Edit](#) [Delete](#)

Usage profile details

Usage profile name dev-profile-1	Status Assigned	Created on October 18, 2023, 14:32 (UTC+3:30)
-------------------------------------	--------------------	--

Usage profile description
A long description of the flow. Long description of the flow.

Assigned IAM roles (8)

Find IAM roles

- AmazonSageMakerServiceCatalogProductsCloudformationRole
- GlueRedshiftDevRole
- GlueRedshiftTestRole
- GlueRedshiftTestRole-2
- GlueEMRRole
- GlueEMRDevRole
- GlueTestRole
- GlueAppFlowRole

Assigned IAM users (100)

Find IAM users

- glue-dev-user-1
- glue-dev-user-2
- glue-dev-user-3
- glue-test-user-1
- glue-test-user-2
- glue-test-user-3
- glue-product-user-1
- glue-product-user-1

Parameter configurations for jobs

Number of workers			Worker type	
Default	Minimum	Maximum	Default type	Allowed types
10	1	20	G.2X	-

Timeout (minutes)		
Default	Minimum	Maximum
2880	100	4000

Parameter configurations for sessions

Number of workers			Worker type	
Default	Minimum	Maximum	Default type	Allowed types
10	1	20	G.1X	G.1X, G.4X, G.8X

Timeout (minutes)			Idle timeout (minutes)		
Default	Minimum	Maximum	Default	Minimum	Maximum
2880	100	4000	30	10	200

Tags (3)

A tag is a label that you assign to an AWS resource. Each tag consists of a key and an optional value. You can use tags to search and filter your resources or track your AWS costs.

Find tags

Key	Value
Key-1	Value-1
Key-2	Value-2
Key-3	Value-3

Use status of the file management

使用状況プロファイルとジョブ

使用状況プロファイルを使用したジョブの作成

ジョブの作成中に、使用プロファイルで設定された制限とデフォルトが適用されます。プロファイルは保存時にジョブに割り当てられます。

使用状況プロファイルを使用したジョブの実行

ジョブの実行を開始すると、は発信者のプロファイルに設定された制限 AWS Glue を適用します。直接発信者がいない場合、Glue は作成者によってジョブに割り当てられたプロファイルの制限を適用します。

Note

ジョブがスケジュールに従って (AWS Glue ワークフローまたは AWS Glue トリガーによって) 実行されると、作成者が適用するジョブに割り当てられたプロファイル。外部サービス (Step Functions、MWAA) または StartJobRun API によってジョブが実行されると、呼び出し元のプロファイル制限が適用されます。

AWS Glue ワークフローまたは AWS Glue トリガーの場合: 既存のジョブを更新して新しいプロファイル名を保存し、プロファイルの制限 (最小、最大、許可されたワーカー) をスケジュールされた実行の実行時に強制する必要があります。

ジョブに割り当てられた使用状況プロファイルの表示

ジョブに割り当てられたプロファイル (スケジュールされた AWS Glue ワークフローまたは AWS Glue トリガーで実行時に使用される) を表示するには、ジョブの詳細タブを確認します。ジョブ実行の詳細タブで、過去の実行で使用されたプロファイルを確認することもできます。

ジョブにアタッチされた使用状況プロファイルの更新または削除

ジョブに割り当てられたプロファイルは、更新時に変更されます。作成者に使用状況プロファイルが割り当てられていない場合、以前にジョブにアタッチされたプロファイルはそこから削除されます。

AWS Glue Data Catalog の開始方法

AWS Glue Data Catalog は永続的な技術メタデータストアです。AWS クラウドでメタデータを保存、注釈付け、および共有するために使用することができるマネージド型サービスです。詳細については、「[AWS Glue Data Catalog](#)」を参照してください。

AWS Glue コンソールおよび一部のユーザーインターフェイスが最近更新されました。

概要

このチュートリアルを使用して、最初の AWS Glue データカタログを作成できます。このデータカタログは、データソースとして Simple Storage Service (Amazon S3) バケットを使用します。

このチュートリアルでは、AWS Glue コンソールを使用して以下の作業を行います。

1. データベースの作成
2. テーブルを作成する
3. データソースとして Amazon S3 バケットを使用します。

これらのステップを完了すると、データソースとして Amazon S3 バケットを使用して、AWS Glue データカタログにデータを入力できるようになります。

ステップ 1: データベースを作成する

開始するには、AWS Management Console にサインインして [AWS Glue コンソール](#) を開きます。

AWS Glue コンソールを使用してデータベースを作成する:

1. AWS Glue コンソールで、左側のメニューにある [Data catalog] (データカタログ) から、[Databases] (データベース) を選択します。
2. [Add database] (データベースの追加) を選択します。
3. [データベースの作成] ページで、データベースの名前を入力します。[Location - optional] セクションで、データカタログのクライアントが使用する URI の場所を設定します。これがわからない場合は、データベースの作成を続行してください。
4. (オプション)。データベースの説明を入力します。

5. [データベースの作成] を選択します。

これで、AWS Glue コンソールで最初のデータベースの作成が終わりました。新しいデータベースが、使用可能なデータベースのリストに表示されます。データベースは、[Databases] (データベース) ダッシュボードからデータベースの名前を選択することで編集できます。

次のステップ

[Other ways to create a database:] (データベースを作成するその他の方法)

AWS Glue コンソールでデータベースを作成しましたが、データベースを作成する方法は他にもあります。

- クローラーを使用して、データベースとテーブルを自動的に作成することができます。クローラーを使用してデータベースを設定するには、「[AWS Glue コンソールでクローラーを使用する](#)」を参照してください。
- AWS CloudFormation テンプレートを使用できます。「[AWS Glue リソースをAWS Glue Data Catalog テンプレートで作成する](#)」を参照してください。
- AWS Glue データベース API オペレーションを使用してデータベースを作成することもできます。

create オペレーションを使用してデータベースを作成するには、DatabaseInput (必須) パラメータを含めてリクエストを構成します。

例:

CLI、Boto3、または DDL を使用して、チュートリアルで使用した S3 バケットからの同じ flights_data.csv ファイルに基づいてテーブルを定義する方法の一例を以下に示します。

CLI

```
aws glue create-database --database-input "{\"Name\":\"clidb\"}"
```

Boto3

```
glueClient = boto3.client('glue')  
  
response = glueClient.create_database(  

```

```
DatabaseInput={  
    'Name': 'boto3db'  
}  
)
```

Database API のデータタイプ、構造、およびオペレーションの詳細については、「[データベース API](#)」を参照してください。

次のステップ

次のセクションでは、テーブルを作成し、そのテーブルをデータベースに追加します。

データカタログの設定とアクセス許可を確認することもできます。「[AWS Glue コンソールでデータカタログ設定を使用する](#)」を参照してください。

ステップ 2。テーブルを作成する

このステップでは、AWS Glue コンソールを使用してテーブルを作成します。

1. AWS Glue コンソールで、左側のメニューにある [Tables] (テーブル) を選択します。
2. [Add table (テーブルの追加)] を選択します。
3. [Table details] (テーブルの詳細) でテーブルの名前を入力して、テーブルのプロパティを設定します。
4. [Databases] (データベース) セクションのドロップダウンメニューから、ステップ 1 で作成したデータベースを選択します。
5. [Add a data store] (データストアの追加) セクションのデフォルト状態では、ソースのタイプとして [S3] が選択されています。
6. [Data located in] (データがある場所) で、[Specified path in another account] (他のアカウントの指定したパス) を選択します。
7. パスをコピーして、[Include path] (パスを含める) の入力フィールドに貼り付けます。

```
s3://crawler-public-us-west-2/flight/2016/csv/
```

8. [Data format] (データ形式) セクションにある [Classification] (分類) で [CSV] を選択し、[Delimiter] (区切り記号) では [comma (,)] (カンマ (,)) を選択します。[Next] を選択します。
9. スキーマを定義するように求められます。スキーマは、データレコードの構造と形式を定義します。[Add column] (列の追加) を選択します。(詳細については、「[スキーマレジストリ](#)」を参照してください)。

10. 列のプロパティを指定します。
 - a. 列名を入力します。
 - b. [Column type] (列タイプ) の場合、デフォルトで 'string' が既に選択されています。
 - c. [Column number] (列番号) の場合、デフォルトで '1' が既に選択されています。
 - d. 追加] を選択します。
11. パーティションインデックスを追加するように求められます。これはオプションです。このステップをスキップするには、[Next] (次へ) を選択します。
12. テーブルプロパティのサマリーが表示されます。すべてが期待とおりに動作している場合、[作成] を選択します。それ以外の場合は、[Back] (戻る) を選択して、必要に応じて編集します。

これで、テーブルを手動で作成し、データベースに関連付ける作業が完了しました。新しく作成したテーブルは [Tables] (テーブル) ダッシュボードに表示されます。ダッシュボードからは、すべてのテーブルの編集と管理を行うことができます。

詳細については、「[AWS Glue コンソールでテーブルを操作する](#)」を参照してください。

次のステップ

次のステップ

データカタログを入力した後は、AWS Glue でジョブの作成を開始できます。「[AWS Glue Studio でビジュアル ETL ジョブの作成](#)」を参照してください。

コンソールを使用する以外にも、データカタログでテーブルを定義する方法として、次のような方法があります。

- [クローラーを作成して実行する](#)
- [AWS Glue のクローラーに分類子の追加](#)
- [AWS Glue テーブル API を使用する](#)
- [AWS Glue Data Catalog テンプレートを使用する](#)
- [Apache Hive メタストアを移行する](#)
- [AWS CLI](#)、Boto3、またはデータ定義言語 (DDL) を使用する

CLI、Boto3、または DDL を使用して、チュートリアルで使用した S3 バケットからの同じ flights_data.csv ファイルに基づいてテーブルを定義する方法の一例を以下に示します。

AWS CLI コマンドを構造化する方法については、ドキュメントを参照してください。CLI の例には、「aws glue create-table --table-input」値の JSON 構文が含まれています。

CLI

```
{
  "Name": "flights_data_cli",
  "StorageDescriptor": {
    "Columns": [
      {
        "Name": "year",
        "Type": "bigint"
      },
      {
        "Name": "quarter",
        "Type": "bigint"
      }
    ],
    "Location": "s3://crawler-public-us-west-2/flight/2016/csv",
    "InputFormat": "org.apache.hadoop.mapred.TextInputFormat",
    "OutputFormat":
"org.apache.hadoop.hive.q1.io.HiveIgnoreKeyTextOutputFormat",
    "Compressed": false,
    "NumberOfBuckets": -1,
    "SerdeInfo": {
      "SerializationLibrary":
"org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe",
      "Parameters": {
        "field.delim": ",",
        "serialization.format": ","
      }
    }
  },
  "PartitionKeys": [
    {
      "Name": "mon",
      "Type": "string"
    }
  ],
  "TableType": "EXTERNAL_TABLE",
  "Parameters": {
    "EXTERNAL": "TRUE",
    "classification": "csv",
  }
}
```

```
        "columnsOrdered": "true",
        "compressionType": "none",
        "delimiter": ",",
        "skip.header.line.count": "1",
        "typeOfData": "file"
    }
}
```

Boto3

```
import boto3

glue_client = boto3.client("glue")

response = glue_client.create_table(
    DatabaseName='sampledb',
    TableInput={
        'Name': 'flights_data_manual',
        'StorageDescriptor': {
            'Columns': [{
                'Name': 'year',
                'Type': 'bigint'
            }, {
                'Name': 'quarter',
                'Type': 'bigint'
            }
        ],
            'Location': 's3://crawler-public-us-west-2/flight/2016/csv',
            'InputFormat': 'org.apache.hadoop.mapred.TextInputFormat',
            'OutputFormat':
'org.apache.hadoop.hive.q1.io.HiveIgnoreKeyTextOutputFormat',
            'Compressed': False,
            'NumberOfBuckets': -1,
            'SerdeInfo': {
                'SerializationLibrary':
'org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe',
                'Parameters': {
                    'field.delim': ',',
                    'serialization.format': ','
                }
            }
        },
        'PartitionKeys': [{
```

```
    'Name': 'mon',
    'Type': 'string'
  ]],
  'TableType': 'EXTERNAL_TABLE',
  'Parameters': {
    'EXTERNAL': 'TRUE',
    'classification': 'csv',
    'columnsOrdered': 'true',
    'compressionType': 'none',
    'delimiter': ',',
    'skip.header.line.count': '1',
    'typeOfData': 'file'
  }
}
```

DDL

```
CREATE EXTERNAL TABLE `sampledb`.`flights_data` (
  `year` bigint,
  `quarter` bigint)
PARTITIONED BY (
  `mon` string)
ROW FORMAT DELIMITED
  FIELDS TERMINATED BY ','
STORED AS INPUTFORMAT
  'org.apache.hadoop.mapred.TextInputFormat'
OUTPUTFORMAT
  'org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat'
LOCATION
  's3://crawler-public-us-west-2/flight/2016/csv/'
TBLPROPERTIES (
  'classification'='csv',
  'columnsOrdered'='true',
  'compressionType'='none',
  'delimiter'=',',
  'skip.header.line.count'='1',
  'typeOfData'='file')
```

データストアへのネットワークアクセスを設定する

抽出、変換、ロード (ETL) ジョブを実行するには、AWS Glue がデータストアにアクセスできる必要があります。ジョブを Virtual Private Cloud (VPC) サブネットで行う必要がない場合 (例えば、Amazon S3 から Amazon S3 へのデータ変換など)、追加の設定は必要ありません。

ジョブを VPC サブネットで行う必要がある場合 (プライベートサブネット内の JDBC データストアからのデータの変換など)、AWS Glue は、ジョブを VPC 内の他のリソースに安全に接続できるようにする [Elastic Network Interfaces](#) をセットアップします。それぞれの Elastic Network Interface には、指定したサブネット内の IP アドレス範囲からプライベート IP アドレスが割り当てられます。パブリック IP アドレスは割り当てられません。AWS Glue 接続で指定されたセキュリティグループは、各 Elastic Network Interface に適用されます。詳しくは、「[AWS Glue から Amazon RDS データストアに JDBC 接続するための Amazon VPC の設定](#)」を参照してください。

ジョブによってアクセスされるすべての JDBC データストアは、VPC サブネットから使用できる必要があります。VPC 内から Amazon S3 にアクセスするには [VPC エンドポイント](#) が必須です。ジョブが VPC リソースとパブリックインターネットの両方にアクセスする必要がある場合は、VPC 内にネットワークアドレス変換 (NAT) ゲートウェイが必要になります。

ジョブまたは開発エンドポイントは、一度に 1 つの VPC (およびサブネット) にのみアクセスできます。異なる VPC のデータストアにアクセスする必要がある場合は、次のオプションがあります。

- VPC ピア接続を使用してデータストアにアクセスします。VPC ピア接続の詳細については、「[VPC ピア接続の基本](#)」を参照してください。
- 中間的なストレージとして Amazon S3 バケットを使用します。ジョブ 1 での Amazon S3 出力をジョブ 2 への入力としながら、作業を 2 つのジョブに分割します。

Amazon VPC を使用して Amazon Redshift データストアに接続する方法の詳細については、「[the section called “Redshift の設定”](#)」を参照してください。

Amazon VPC を使用して Amazon RDS データストアに接続する方法の詳細については、「[the section called “Amazon RDS データストアに接続するための VPC の設定”](#)」を参照してください。

Amazon VPC で必要なルールを設定したら、データストアに接続するために必要なプロパティを使用して、AWS Glue で接続を作成します。接続の詳細については、「[データへの接続](#)」を参照してください

Note

AWS Glue の DNS 環境を必ず設定してください。詳細については、「[VPC での DNS のセットアップ](#)」を参照してください。

トピック

- [AWS Glue の PyPI に接続する VPC のセットアップ](#)
- [VPC での DNS のセットアップ](#)

AWS Glue の PyPI に接続する VPC のセットアップ

Python Package Index (PyPI) は、Python プログラミング言語用のソフトウェアのリポジトリです。このトピックでは、(セッションの作成者が `--additional-python-modules` フラグを使用して指定した) pip インストールパッケージの使用をサポートするために必要な詳細について説明します。

コネクタで AWS Glue インタラクティブセッションを使用すると、コネクタに指定されたサブネット経由で VPC ネットワークが使用されます。そのため、特別な設定をセットアップしない限り、AWS サービスやその他のネットワーク送信先は使用できません。

この問題には、次のような解決策があります。

- セッションからアクセス可能なインターネットゲートウェイの使用。
- パッケージセットの依存関係の推移閉包を含む PyPI/simple リポジトリを使用した S3 バケットのセットアップと使用。
- PyPI をミラーリングし、VPC にアタッチされている CodeArtifact リポジトリの使用。

インターネットゲートウェイをセットアップする

技術面については、「[NAT ゲートウェイのユースケース](#)」で詳しく説明していますが、`--additional-python-modules` を使用するには以下の要件に注意してください。具体的には、`--additional-python-modules` には、VPC の設定によって決まる `pypi.org` へのアクセス権が必要です。次の要件に注意してください。

1. ユーザーのセッションに追加の Python モジュールを pip インストール経由でインストールする要件。セッションがコネクタを使用する場合、設定が影響を受ける可能性があります。

2. `--additional-python-modules` でコネクタを使用している場合、セッションが開始されると、`pypi.org` に到達するためのネットワークパスがコネクタの `PhysicalConnectionRequirements` に関連付けられたサブネットに指定される必要があります。
3. 設定が正しいかどうかを判断する必要があります。

ターゲット PyPI/simple リポジトリをホストする Amazon S3 バケットのセットアップ

この例では、Amazon S3 に一連のパッケージとその依存関係について PyPI のミラーリングをセットアップします。

一連のパッケージに PyPI のミラーリングを設定するには:

```
# pip download all the dependencies
pip download -d s3pypi --only-binary :all: plotly ggplot
pip download -d s3pypi --platform manylinux_2_17_x86_64 --only-binary :all: psycpg2-binary
# create and upload the pypi/simple index and wheel files to the s3 bucket
s3pypi -b test-domain-name --put-root-index -v s3pypi/*
```

既存のアーティファクトリポジトリがある場合は、pip 用のインデックス URL があります。これを、上記の Amazon S3 バケットのサンプル URL の代わりに指定できます。

カスタム `index-url` をいくつかのサンプルパッケージで使用するには:

```
%%configure
{
  "--additional-python-modules": "psycpg2_binary==2.9.5",
  "python-modules-installer-option": "--no-cache-dir --verbose --index-url https://test-domain-name.s3.amazonaws.com/ --trusted-host test-domain-name.s3.amazonaws.com"
}
```

VPC にアタッチされた PyPI の CodeArtifact ミラーリングのセットアップ

ミラーリングをセットアップするには:

1. コネクタが使用するサブネットと同じリージョンにリポジトリを作成します。

[Public upstream repositories]、[pypi-store] の順に選択します。

2. サブネットの VPC からリポジトリにアクセスできるようにします。
3. `python-modules-installer-option` を使用して正しい `--index-url` を指定します。

```
%%configure
{
  "--additional-python-modules": "psycopg2_binary==2.9.5",
  "python-modules-installer-option": "--no-cache-dir --verbose --index-url https://
test-domain-name.s3.amazonaws.com/ --trusted-host test-domain-name.s3.amazonaws.com"
}
```

詳細については、「[VPC から CodeArtifact を使用する](#)」を参照してください。

VPC での DNS のセットアップ

ドメインネームシステム (DNS) は、インターネットで使用する名前を対応する IP アドレスに解決するための標準です。DNS ホスト名は、ホスト名とドメイン名で構成され、コンピュータに一意の名前を付けます。DNS サーバーは DNS ホスト名を対応する IP アドレスに解決します。

VPC で DNS をセットアップするには、DNS ホスト名と DNS 解決の両方が VPC で有効になっていることを確認します。VPC ネットワーク属性の `enableDnsHostnames` と `enableDnsSupport` を `true` に設定する必要があります。これらの属性を表示および変更するには、VPC コンソール (<https://console.aws.amazon.com/vpc/>) を開きます。

詳細については、「[Using DNS with Your VPC](#)」を参照してください。また、AWS CLI を使用して `modify-vpc-attribute` コマンドを実行し、VPC ネットワーク属性を設定することもできます。

Note

Route 53 を使用している場合は、設定によって DNS ネットワーク属性が上書きされないようにする必要があります。

AWS Glue での暗号化のセットアップ

次のワークフローの例では、AWS Glue で暗号化を使用するときに設定するオプションに注目します。この例では、特定の AWS Key Management Service (AWS KMS) キーの使用法を示していますが、特定のニーズに基づいて他の設定を選択することもできます。このワークフローでは、AWS Glue の設定時に暗号化に関連するオプションのみをハイライト表示しています。

1. AWS Glue コンソールのユーザーがすべての AWS Glue API オペレーションを許可するアクセス許可ポリシー (例: "glue:*") を使用しない場合は、次のアクションが許可されていることを確認してください。

- "glue:GetDataCatalogEncryptionSettings"
- "glue:PutDataCatalogEncryptionSettings"
- "glue:CreateSecurityConfiguration"
- "glue:GetSecurityConfiguration"
- "glue:GetSecurityConfigurations"
- "glue>DeleteSecurityConfiguration"

2. 暗号化されたカタログにアクセスまたは書き込みするクライアント (つまりコンソールユーザー、クローラー、ジョブ、開発エンドポイントなど) には、以下のアクセス許可が必要です。

```
{
  "Version": "2012-10-17",
  "Statement": {
    "Effect": "Allow",
    "Action": [
      "kms:GenerateDataKey",
      "kms:Decrypt",
      "kms:Encrypt"
    ],
    "Resource": "<key-arns-used-for-data-catalog>"
  }
}
```

3. 暗号化された接続パスワードにアクセスするユーザーまたはロールには、次のアクセス許可が必要です。

```
{
  "Version": "2012-10-17",
  "Statement": {
    "Effect": "Allow",
    "Action": [
      "kms:Decrypt"
    ],
    "Resource": "<key-arns-used-for-password-encryption>"
  }
}
```

4. 暗号化されたデータを Amazon S3 に書き込む抽出、変換、ロード (ETL) ジョブのロールには、以下のアクセス許可が必要です。

```
{
  "Version": "2012-10-17",
  "Statement": {
    "Effect": "Allow",
    "Action": [
      "kms:Decrypt",
      "kms:Encrypt",
      "kms:GenerateDataKey"
    ],
    "Resource": "<key-arns-used-for-s3>"
  }
}
```

5. 暗号化された Amazon CloudWatch Logs を書き込む ETL ジョブまたはクローラーでは、キーポリシーおよび IAM ポリシーに以下のアクセス許可が必要です。

キーポリシー (IAM ポリシーではない) の場合:

```
{
  "Effect": "Allow",
  "Principal": {
    "Service": "logs.region.amazonaws.com"
  },
  "Action": [
    "kms:Encrypt*",
    "kms:Decrypt*",
    "kms:ReEncrypt*",
    "kms:GenerateDataKey*",
    "kms:Describe*"
  ],
  "Resource": "<arn of key used for ETL/crawler cloudwatch encryption>"
}
```

キーポリシーの詳細については、AWS Key Management Service デベロッパーガイドの「[Using Key Policies in AWS KMS](#)」を参照してください。

logs:AssociateKmsKey 許可にアタッチする IAM ポリシーの場合:

```
{
```

```
"Effect": "Allow",
"Principal": {
  "Service": "logs.region.amazonaws.com"
},
"Action": [
  "logs:AssociateKmsKey"
],
"Resource": "<arn of key used for ETL/crawler cloudwatch encryption>"
}
```

6. 暗号化されたジョブのブックマークを使用する ETL ジョブには、次のアクセス許可が必要です。

```
{
  "Version": "2012-10-17",
  "Statement": {
    "Effect": "Allow",
    "Action": [
      "kms:Decrypt",
      "kms:Encrypt"
    ],
    "Resource": "<key-arns-used-for-job-bookmark-encryption>"
  }
}
```

7. AWS Glue コンソールで、ナビゲーションペインの [設定] を選択します。

- a. [Data catalog settings] (データカタログ設定) ページで [Metadata encryption] (メタデータの暗号化) を選択し、Data Catalog を暗号化します。このオプションは、選択した AWS KMS キーを使用して、Data Catalog のすべてのオブジェクトを暗号化します。
- b. [AWS KMS key] で、[aws/glue] を選択します。自身で作成した AWS KMS キーを選択することもできます。

Important

AWS Glue は、対称カスターマスターキー (CMK) のみをサポートしています。[AWS KMS key] (KMS キー) リストには、対称キーのみが表示されます。ただし、[Choose a AWS KMS key ARN] (KMS キー ARN を選択します) を選択した場合、任意のキータイプの ARN をコンソールで入力します。対称キーには ARN だけを入力するようにしてください。

暗号化が有効になっている場合、Data Catalog にアクセスするクライアントには、AWS KMS のアクセス許可が付与されている必要があります。

8. ナビゲーションペインで、[セキュリティ設定] を選択します。セキュリティ設定は、AWS Glue プロセスを設定するために使用できるセキュリティプロパティのセットです。次に [セキュリティ設定の追加] を選択します。設定で、次のオプションのいずれかを選択します。
 - a. [S3 暗号化] を選択します。[暗号化モード] では、SSE-KMSを選択します。[AWS KMS キー] については、[aws/s3] を選択します (ユーザーがこのキーを使用するアクセス権限を持っていることを確認してください)。これにより、ジョブによって Amazon S3 に書き込まれたデータは、AWS 管理の AWS Glue AWS KMS キーを使用できるようになります。
 - b. [CloudWatch Logs の暗号化] を選択し、CMK を選択します (ユーザーがこのキーを使用するアクセス許可があることを確認してください)。詳しくは、AWS Key Management Service デベロッパーガイドの「[Encrypt Log Data in CloudWatch Logs Using AWS KMS](#)」を参照してください。

 Important

AWS Glue は、対称カスターマスターキー (CMK) のみをサポートしています。[AWS KMS key] (KMS キー) リストには、対称キーのみが表示されます。ただし、[Choose a AWS KMS key ARN] (KMS キー ARN を選択します) を選択した場合、任意のキータイプの ARN をコンソールで入力します。対称キーには ARN だけを入力するようにしてください。

- c. [詳細プロパティ] を選択し、[ジョブのブックマーク暗号化モード] を選択します。[AWS KMS キー] には、[aws/glue] を選択します (ユーザーがこのキーを使用するアクセス権限を持っていることを確認してください)。これにより、AWS Glue AWS KMS キーを使用して Amazon S3 に書き込まれたジョブブックマークの暗号化が可能になります。
9. ナビゲーションペインで [Connections (接続)] を選択します。
 - a. [接続の追加] を選択して、ETL ジョブのターゲットとなる Java Database Connectivity (JDBC) データ・ストアへの接続を作成します。
 - b. Secure Sockets Layer (SSL) 暗号化が使用されるようにするには、[Require SSL connection (SSL接続を要求)] を選択し、接続をテストします。
10. ナビゲーションペインで ジョブを選択します。
 - a. [ジョブを追加] を選択して、データを変換するジョブを作成します。
 - b. ジョブ定義で、作成したセキュリティ構成を選択します。

11 AWS Glue コンソールで、オンデマンドでジョブを実行します。ジョブによって書き込まれた Amazon S3 データ、ジョブによって書き込まれた CloudWatch Logs、およびジョブブックマークがすべて暗号化されていることを確認します。

AWS Glue のための開発用ネットワークの設定

AWS Glue で抽出、変換、ロード (ETL) スクリプトを実行するには、開発エンドポイントを使用してスクリプトを開発し、テストすることができます。開発エンドポイントは、AWS Glue バージョン 2.0 ジョブでの使用がサポートされていません。バージョン 2.0 以降では、Jupyter Notebook といずれかの AWS Glue カーネルを使用する開発方法が推奨されています。詳しくは、「[the section called “AWS Glue のインタラクティブセッションの開始方法”](#)」を参照してください。

開発エンドポイント用にネットワークを設定する

開発エンドポイントを設定するときは、Virtual Private Cloud (VPC)、サブネット、およびセキュリティグループを指定します。

Note

AWS Glue の DNS 環境を必ず設定してください。詳しくは、「[VPC での DNS のセットアップ](#)」を参照してください。

AWS Glue による必要なリソースに対するアクセスを許可するには、サブネットのルートテーブルに行を追加して、Amazon S3 のプレフィックスリストを VPC エンドポイントに関連付けます。プレフィックスリスト ID は、VPC からのトラフィックが VPC エンドポイント経由で AWS のサービスにアクセスできるようにする、アウトバウンドセキュリティグループのルールを作成するために必要です。この開発エンドポイントに関連付けられているノートブックサーバーへの接続をローカルマシンから簡単に行うには、ルートテーブルに行を追加してインターネットゲートウェイ ID を追加します。詳細については、「[VPC エンドポイント](#)」を参照してください。サブネットのルートテーブルを更新すると次の表のようになります。

デスティネーション	ターゲット		
10.0.0.0/16	local		
Amazon S3 の pl-id	vpce-id		

デスティネーション	ターゲット		
0.0.0.0/0	igw-xxxx		

AWS Glue がコンポーネント間で通信できるようにするには、すべての TCP ポートに対して自己参照のインバウンドルールを持つセキュリティグループを指定します。自己参照ルールを作成することで、ソースを VPC 内の同じセキュリティグループに制限することができ、ネットワーク全体には公開されません。VPC のデフォルトのセキュリティグループには、すでに ALL Traffic (すべてのトラフィック) の自己参照インバウンドルールがある場合があります。

セキュリティグループを設定するには

1. AWS Management Console にサインインし、Amazon EC2 コンソール (<https://console.aws.amazon.com/ec2/>) を開きます。
2. 左のナビゲーションペインで セキュリティグループを選択します。
3. リストから既存のセキュリティグループを選択するか、[Create Security Group] (セキュリティグループの作成) を選択して、開発エンドポイントで使用します。
4. セキュリティグループペインで、[Inbound] (インバウンド) タブに移動します。
5. 自己参照ルールを追加して、AWS Glue コンポーネントが通信できるようにします。具体的には、[Type (タイプ)] All TCP、[Protocol (プロトコル)] は TCP、[Port Range (ポート範囲)] にはすべてのポートが含まれ、[Source (ソース)] は [Group ID (グループ ID)] と同じセキュリティグループ名であるというルールを追加または確認します。

インバウンドルールは次のようになります。

タイプ	プロトコル	ポート範囲	ソース
すべての TCP	TCP	0 ~ 65535	<i>security-group</i>

次に、自己参照インバウンドルールの例を示します。

6. アウトバウンドトラフィックのルールも追加します。すべてのポートへのアウトバウンドトラフィックを開くか、または [Type (タイプ)] All TCP、[Protocol (プロトコル)] は TCP、[Port Range (ポート範囲)] にすべてのポートが含まれ、[Source (ソース)] は [Group ID (グループ ID)] と同じセキュリティグループ名の自己参照ルールを作成します。

アウトバウンドルールは、次のいずれかのルールのようになります。

タイプ	プロトコル	ポート範囲	デスティネーション
すべての TCP	TCP	0 ~ 65535	<i>security-group</i>
すべてのトラフィック	すべて	すべて	0.0.0.0/0

ノートブックサーバー用の Amazon EC2 のセットアップ

開発エンドポイントでは、Jupyter Notebook で ETL スクリプトをテストするためのノートブックサーバーを作成できます。ノートブックとの通信を可能にするには、HTTPS (ポート 443) と SSH (ポート 22) の両方のインバウンドルールを持つセキュリティグループを指定します。ルールのソースが 0.0.0.0/0 か、ノートブックに接続しているマシンの IP アドレスであることを確認してください。

セキュリティグループを設定するには

1. AWS Management Console にサインインし、Amazon EC2 コンソール (<https://console.aws.amazon.com/ec2/>) を開きます。
2. 左のナビゲーションペインで セキュリティグループを選択します。
3. リストから既存のセキュリティグループを選択するか、[Create Security Group] (セキュリティグループの作成) を選択して、ノートブックサーバーで使用します。開発エンドポイントに関連付けられているセキュリティグループは、ノートブックサーバーの作成にも使用されます。
4. セキュリティグループペインで、[Inbound] (インバウンド) タブに移動します。
5. 次のようなインバウンドルールを追加します。

タイプ	プロトコル	ポート範囲	ソース
SSH	TCP	22	0.0.0.0/0
HTTPS	TCP	443	0.0.0.0/0

次に、セキュリティグループのインバウンドルールの例を示します。

Security Group: sg-19e1b768

...

Description

Inbound

Outbound

Tags

Edit

Type ⓘ	Protocol ⓘ	Port Range ⓘ	Source ⓘ
SSH	TCP	22	0.0.0.0/0
HTTPS	TCP	443	0.0.0.0/0

AWS Glue でのデータ検出とカタログ化

AWS Glue Data Catalog は、組織のデータセットに関するメタデータを保存する一元化されたリポジトリです。データソースの場所、スキーマ、およびランタイムメトリクスへのインデックスとして機能します。メタデータはメタデータテーブルに保存され、ここでは各テーブルが 1 つのデータストアを表します。

データソースを自動的にスキャンしてメタデータを抽出するクローラーを使用してデータカタログに入力できます。クローラーは、内部 (AWS ベース) および AWS 外部のデータソースに接続できます。

サポートされるデータソースについては、「[クロール可能なデータストア](#)」を参照してください。

特定の要件に応じてテーブル構造、スキーマ、パーティション構造を定義することで、データカタログにテーブルを手動で作成することもできます。

メタデータテーブルの手動作成については、「[メタデータの手動定義](#)」を参照してください。

データカタログ内の情報を使用して、ETL ジョブを作成し、モニタリングできます。データカタログは他の AWS 分析サービスと統合され、データソースの統合ビューを提供するため、データの管理と分析が容易になります。

- Amazon Athena – SQL を使用して、Amazon S3 データのデータカタログにテーブルメタデータを保存し、クエリします。
- AWS Lake Formation – きめ細かなデータアクセスポリシーを一元的に定義および管理し、データアクセスを監査します。
- Amazon EMR – ビッグデータ処理のためにデータカタログで定義されたデータソースにアクセスします。
- Amazon SageMaker – 機械学習モデルを迅速かつ確実に構築、トレーニング、デプロイします。

データカタログの主な機能

データカタログの主な側面を次に示します。

メタデータリポジトリ

データカタログは中央メタデータリポジトリとして機能し、データソースの場所、スキーマ、プロパティに関する情報を保存します。このメタデータは、従来のリレーショナルデータベースカタログと同様に、データベースとテーブルにまとめられます。

自動データ検出可能性

AWS Glue クローラー では、新規または更新されたデータソースを自動的に検出してカタログ化できるため、手動メタデータ管理のオーバーヘッドが軽減され、データカタログが最新の状態を維持できます。データソースをカタログ化することで、データカタログでは、ユーザーやアプリケーションが組織内で利用可能なデータアセットを簡単に検出して理解できるようになり、データの再利用とコラボレーションが促進されます。

データカタログは、Amazon S3、Amazon RDS、Amazon Redshift、Apache Hive など、さまざまなデータソースをサポートしています。AWS Glue クローラー を使用して、これらのソースからメタデータを自動的に推測して保存できます。

詳細については、「[クローラーを使用したデータカタログへの入力](#)」を参照してください。

スキーマ管理

データカタログは、スキーマの推論、進化、バージョニングなど、データソースのスキーマを自動的にキャプチャして管理します。AWS Glue ETL ジョブを使用して、データカタログ内のスキーマとパーティションを更新できます。

テーブル最適化

Amazon Athena、Amazon EMR、AWS Glue ETL ジョブなどの AWS 分析サービスによる読み取りパフォーマンスを向上させるために、データカタログは、データカタログ内の Iceberg テーブル用にマネージド圧縮 (小さな Amazon S3 オブジェクトを圧縮してより大きなオブジェクトにコンパクト化するプロセス) を提供しています。AWS Glue コンソール、AWS Lake Formation コンソール、AWS CLI API、または AWS API を使用して、データカタログ内の個々の Iceberg テーブルの圧縮を有効または無効にすることができます。

詳細については、[Iceberg テーブルの最適化](#) を参照してください。

列統計

追加のデータパイプラインを設定することなく、Parquet、ORC、JSON、ION、CSV、XML などのデータ形式でデータカタログテーブルの列レベルの統計を計算できます。列統計は、列内の値に関するインサイトを得ることで、データプロファイルを理解するのに役立ちます。データカタログは、最小値、最大値、null 値の合計、個別の値の合計、値の平均長、true 値の合計出現数などの列の値の統計の生成をサポートしています。

詳細については、[列統計を使用したクエリのパフォーマンスの最適化](#) を参照してください。

データリネージュ

データカタログは、データに対して実行された変換とオペレーションの記録を維持し、データ系統情報を提供します。この系統情報は、データ出所の監査、コンプライアンス、理解に役立ちます。

AWS の他のサービスとの統合

データカタログは、AWS Lake Formation、Amazon Athena、Amazon Redshift Spectrum、Amazon EMR などの他の AWS サービスとシームレスに統合されます。この統合により、1 つの一貫性のあるメタデータレイヤーを使用して、さまざまなデータストアのデータをクエリおよび分析できます。

セキュリティとアクセスコントロール

AWS Glue は AWS Lake Formation と統合して、データカタログリソースのきめ細かなアクセスコントロールをサポートすることで、組織のポリシーと要件に基づいてアクセス許可を管理し、データアセットに安全にアクセスできます。AWS Glue は AWS Key Management Service (AWS KMS) と統合して、データカタログに保存されているメタデータを暗号化します。

トピック

- [AWS Glue データカタログの入力](#)
- [トランザクションテーブルへのデータ入力と管理](#)
- [データカタログの管理](#)
- [データカタログにアクセスする](#)
- [AWS Glue データカタログのベストプラクティス](#)
- [AWS Glue スキーマレジストリ](#)

AWS Glue データカタログの入力

次の方法を使用して AWS Glue Data Catalog にデータを入力できます。

- AWS Glue クローラー – AWS Glue クローラー は、データベース、データレイク、ストリーミングデータなどのデータソースを自動的に検出してカタログ化できます。クローラーは、さまざまなデータソースのメタデータを自動的に検出して推測できるため、データカタログにデータを入力する最も一般的な推奨方法です。
- メタデータの手動追加 – AWS Glue コンソール、Lake Formation コンソール、AWS CLI API、または AWS Glue API を使用して、データベース、テーブル、および接続の詳細を手動で定義し、

これらをデータカタログに追加できます。手動入力は、クローラできないデータソースをカタログ化する場合に便利です。

- 他の AWS サービスとの統合 — AWS Lake Formation や Amazon Athena などのサービスからのメタデータをデータカタログに入力できます。これらのサービスは、データカタログでデータソースを検出して登録できます。
- 既存のメタデータリポジトリからの入力 – Apache Hive Metastore などの既存のメタデータストアがある場合は、AWS Glue を使用してそのメタデータをデータカタログにインポートできます。詳細については、GitHub で「[Migration between the Hive Metastore and the AWS Glue Data Catalog](#)」を参照してください。

トピック

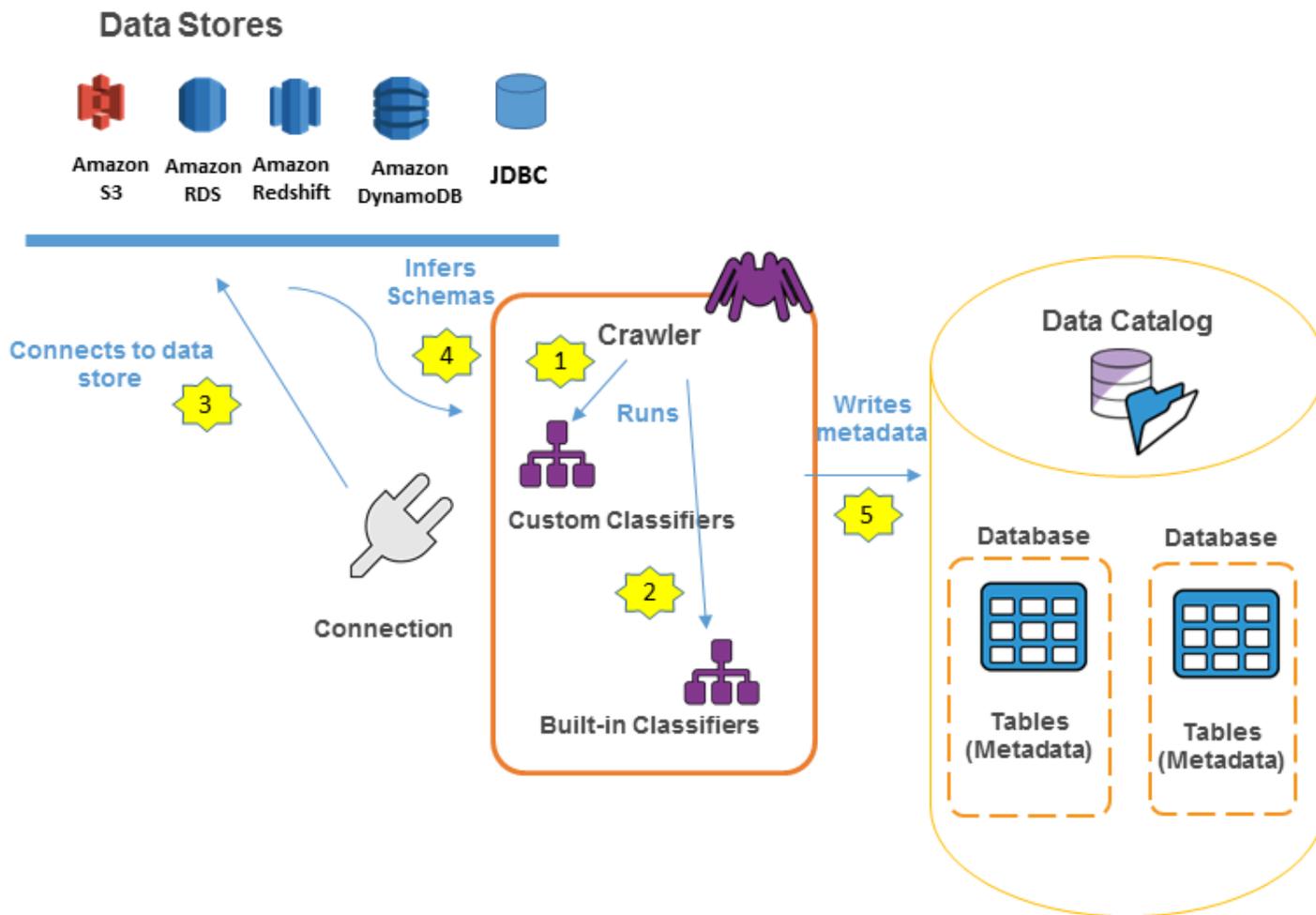
- [クローラーを使用したデータカタログへの入力](#)
- [メタデータの手動定義](#)
- [他の AWS のサービスとの統合](#)
- [データカタログの設定](#)

クローラーを使用したデータカタログへの入力

AWS Glue クローラー を使用して、AWS Glue Data Catalog にデータベースとテーブルを入力できます。これは、AWS Glue ユーザーが最もよく使用する基本的な方法です。クローラーは 1 回の実行で複数のデータストアをクローラできます。完了すると、クローラーはデータカタログで 1 つ以上のテーブルを作成または更新します。AWS Glue で定義した抽出、変換、ロード (ETL) ジョブは、これらのデータカタログテーブルをソースおよびターゲットとして使用します。ETL ジョブは、ソースおよびターゲットのデータカタログテーブルで指定されているデータストアに対して読み取りと書き込みを行います。

ワークフロー

次のワークフロー図は、AWS Glue クローラーがデータストアや他の要素とやり取りしてデータカタログに入力する方法を示しています。



クローラーが AWS Glue Data Catalog に入力する一般的なワークフローを以下に示します。

1. クローラーが選択した任意のカスタム分類子を実行し、データの形式とスキーマを推論します。カスタム分類子のコードを提供すると、指定した順序で実行されます。

データの構造を正常に認識した最初のカスタム分類子がスキーマを作成するために使用されます。リスト内の下位のカスタム分類子はスキップされます。

2. カスタム分類子と一致するデータのスキーマがない場合は、組み込み分類子がデータのスキーマを認識します。組み込み分類子の例に、JSON を認識する分類子があります。
3. クローラーがデータストアに接続します。一部のデータストアでは、クローラーがアクセスするために接続プロパティを必要とします。
4. データの推測されたスキーマが作成されます。

5. クローラーはデータカタログにメタデータを書き込みます。テーブル定義にはデータストア内のデータに関するメタデータが含まれています。テーブルは、Data Catalog でテーブルのコンテナとなるデータベースに書き込まれます。テーブルの属性には分類が含まれます。これは、テーブルのスキーマを推測した分類子により作成されるラベルです。

トピック

- [クローラーの仕組み](#)
- [クローラー可能なデータストア](#)
- [クローラーは、どのようにパーティションを作成するタイミングを判断していますか？](#)
- [クローラーの前提条件](#)
- [クローラーの設定](#)
- [AWS Glue でのクローラーへの分類子の追加](#)
- [AWS Glue クローラーのスケジュール](#)
- [クローラーの結果と詳細の表示](#)
- [クローラーの動作のカスタマイズ](#)
- [チュートリアル: AWS Glue クローラーの追加](#)

クローラーの仕組み

クローラーを実行すると、クローラーは以下のアクションを使用してデータストアを調査します。

- 生データの形式、スキーマ、および関連プロパティを確認するためにデータを分類する – カスタム分類子を作成して分類の結果を設定できます。
- データをテーブルまたはパーティションにグループ化する – データはクローラーのヒューリスティックに基づいてグループ化されます。
- メタデータをデータカタログに書き込む – クローラーでテーブルやパーティションを追加、更新、削除する方法を設定できます。

クローラーを定義する場合、データの形式を評価してスキーマを推測する分類子を 1 つ以上選択します。クローラーを実行すると、リストで最初にデータストアの認識に成功した分類子を使用してテーブルのスキーマが作成されます。組み込み分類子を使用するか、独自に定義することができます。カスタム分類子は、クローラーを定義する前に別のオペレーションで定義します。AWS Glue には、JSON、CSV、Apache Avro などの形式の共通ファイルからスキーマを推論するための組み込み

分類子が用意されています。AWS Glue の組み込み分類子の最新のリストについては、「[AWS Glue の組み込み分類子](#)」を参照してください。

クローラーで作成するメタデータテーブルは、クローラーの定義時にデータベースに含まれます。クローラーがデータベースを指定しない場合、テーブルはデフォルトのデータベースに配置されます。さらに、各テーブルには、最初にデータストアの認識に成功した分類子により入力された分類子の列があります。

クローラーするファイルが圧縮されている場合、クローラーはダウンロードして処理する必要があります。クローラーを実行すると、ファイルを調査して形式と圧縮タイプを判定し、これらのプロパティをデータカタログに書き込みます。一部のファイル形式 (Apache Parquet など) では、ファイルの書き込み時にファイルの部分を圧縮できます。これらのファイルでは、圧縮されたデータはファイルの内部コンポーネントであり、AWS Glue はテーブルをデータカタログ内に書き込むときに `compressionType` プロパティを事前設定しません。一方、ファイル全体を圧縮アルゴリズム (gzip など) で圧縮する場合は、テーブルをデータカタログ内に書き込むときに `compressionType` プロパティが事前設定されます。

クローラーは、作成するテーブルの名前を生成します。AWS Glue Data Catalog に保存されるテーブルの名前は、以下のルールに従います。

- 英数字とアンダースコア (`_`) のみを使用できます。
- カスタムプレフィックスは 64 文字より長くすることはできません。
- 名前の最大長は 128 文字より長くすることはできません。クローラーは、生成した名前が制限内に収まるように切り詰めます。
- 重複するテーブル名が発生した場合、クローラーは名前にハッシュ文字列のサフィックスを追加します。

クローラーが複数回実行される場合 (おそらくスケジュールに基づいて)、データストア内の新規または変更されたファイルやテーブルが検索されます。クローラーの出力には、前回の実行以降に検索された、新しいテーブルとパーティションが含まれています。

クローラー可能なデータストア

クローラーは、次に示すファイルベースおよびテーブルベースのデータストアをクローラーできます。

クローラーが使用するアクセスタイプ	データストア
ネイティブクライアント	<ul style="list-style-type: none"> • Amazon Simple Storage Service (Amazon S3) • Amazon DynamoDB • Delta Lake 2.0.x • Apache Iceberg 1.5 • Apache Hudi 0.14
JDBC	<p>Amazon Redshift</p> <p>Snowflake</p> <p>Amazon Relational Database Service (Amazon RDS) 内、または Amazon RDS の外部:</p> <ul style="list-style-type: none"> • Amazon Aurora • MariaDB • Microsoft SQL Server • MySQL • Oracle • PostgreSQL
MongoDB クライアント	<ul style="list-style-type: none"> • MongoDB • MongoDB Atlas • Amazon DocumentDB (MongoDB 互換性)

 Note

現在、AWS Glue ではデータストリームのクローラーがサポートされていません。

JDBC、MongoDB、MongoDB Atlas、Amazon DocumentDB (MongoDB 互換) データストアでは、クローラーがデータストアに接続するときに使用可能な AWS Glue 接続 を指定する必要があります。Amazon S3 の場合、オプションで Network タイプの接続を指定できます。接続とは、認証情

報、URL、Amazon Virtual Private Cloud 情報などの接続情報を保存するデータカタログオブジェクトです。詳細については、「[データへの接続](#)」を参照してください。

クローラーがサポートするドライバーのバージョンは次のとおりです。

製品	クローラーがサポートするドライバー
PostgreSQL	42.2.1
Amazon Aurora	ネイティブクローラードライバーと同じ
MariaDB	8.0.13
Microsoft SQL Server	6.1.0
MySQL	8.0.13
Oracle	11.2.2
Amazon Redshift	4.1
Snowflake	3.13.20
MongoDB	4.7.2
MongoDB Atlas	4.7.2

次に、さまざまなデータストアに関する注意事項を示します。

Amazon S3

自分のアカウントのパスをクロールするか、または別のアカウントのパスをクロールするかを選択できます。フォルダにあるすべての Amazon S3 ファイルが同じスキーマを持つ場合、クローラーはテーブルを 1 つ作成します。また、Amazon S3 オブジェクトが分割されている場合、メタデータテーブルは 1 つしか作成されず、そのテーブルのパーティション情報がデータカタログに追加されます。

Amazon S3 と Amazon DynamoDB

クローラーは、AWS Identity and Access Management (IAM) ロールをアクセス許可のために使用して、データストアにアクセスします。クローラーに渡すロールは、クロールされる Amazon S3

パスと Amazon DynamoDB テーブルにアクセスするためのアクセス許可を持っている必要があります。

Amazon DynamoDB

AWS Glue コンソールを使用してクローラーを定義するとき、1つのDynamoDB テーブルを指定します。AWS Glue API を使用している場合、テーブルのリストを指定できます。クローラーの実行時間を短縮するために、データの小さなサンプルのみをクロールするように選択できます。

Delta Lake

各 Delta Lake データストアで、Delta テーブルの作成方法を指定します。

- [Create Native tables] (ネイティブテーブルの作成): Delta トランザクションログの直接クエリをサポートするクエリエンジンとの統合を可能にします。詳細については、「[Delta Lake テーブルのクエリを実行する](#)」を参照してください。
- [Create Symlink tables] (Symlink テーブルの作成): 指定された設定パラメータに基づいて、マニフェストファイルを分割キーで分割して `_symlink_manifest` フォルダを作成します。

Iceberg

Iceberg データストアごとに、Iceberg テーブルのメタデータを含む Amazon S3 パスを指定します。クローラーは Iceberg テーブルのメタデータを検出すると、それをデータカタログに登録します。クローラーがテーブルを更新し続けるように、スケジュールを設定できます。

データストアに次のパラメータを定義できます。

- Exclusions: 特定のフォルダをスキップできます。
- Maximum Traversal Depth: Amazon S3 バケットでクローラーがクロールできる深度の制限を設定します。デフォルトの最大トラバーサル深度は 10 で、設定できる最大深度は 20 です。

Hudi

Hudi データストアごとに、Hudi テーブルのメタデータを含む Amazon S3 パスを指定します。クローラーは Hudi テーブルのメタデータを検出すると、それをデータカタログに登録します。クローラーがテーブルを更新し続けるように、スケジュールを設定できます。

データストアに次のパラメータを定義できます。

- Exclusions: 特定のフォルダをスキップできます。
- Maximum Traversal Depth: Amazon S3 バケットでクローラーがクロールできる深度の制限を設定します。デフォルトの最大トラバーサル深度は 10 で、設定できる最大深度は 20 です。

Note

Hudi 0.13.1 およびタイムスタンプ型と互換性がないため、論理型として `millis` を含むタイムスタンプ列は、`bigint` として解釈されます。解決策は、今後の Hudi リリースで提供される可能性があります。

Hudi テーブルは次のように分類され、それぞれが固有の意味を持っています。

- **書き込み時コピー (CoW):** データは列形式 (Parquet) で保存され、更新ごとに書き込み中にファイルの新しいバージョンが作成されます。
- **読み取り時マージ (MoR):** データは、列形式 (Parquet) と行形式 (Avro) の組み合わせを使用して保存されます。更新は、行形式の差分ファイルに記録され、必要に応じて圧縮されて、新しいバージョンの列形式のファイルが作成されます。

CoW データセットでは、レコードが更新されるたびに、そのレコードを含むファイルが更新された値で書き換えられます。MoR データセットでは、更新があるたびに、Hudi によって変更されたレコードの行のみが書き込まれます。MoR は、読み取りが少なく書き込みまたは変更が多いワークロードに適しています。CoW は、頻繁に変更されないデータの読み取りが多いワークロードに適しています。

Hudi にはデータにアクセスするためのクエリタイプが 3 種類用意されています。

- **スナップショットクエリ:** 指定されたコミットまたは圧縮アクションの時点で最新のテーブルスナップショットを参照するクエリです。MoR テーブルの場合、スナップショットクエリは、クエリ時に最新のファイルスライスのベースファイルとデルタファイルをマージして、テーブルの最新の状態を提示します。
- **増分クエリ:** クエリは、指定されたコミット/圧縮以降にテーブルに書き込まれた新しいデータのみを参照します。これにより、増分データパイプラインを有効にするための変更ストリームが効果的に提供されます。
- **読み取り最適化クエリ:** MoR テーブルの場合、クエリは圧縮された最新のデータを参照します。CoW テーブルの場合、クエリはコミットされた最新のデータを参照します。

書き込み時コピー (CoW) テーブルの場合、クローラーは `ReadOptimized serde org.apache.hudi.hadoop.HoodieParquetInputFormat` を使用して、データカタログに 1 つのテーブルを作成します。

読み取り時マージ (MoR) テーブルの場合、クローラーはデータカタログ内の同じテーブルの場所に次の 2 つのテーブルを作成します。

- サフィックス `_ro` を持ち、ReadOptimized serde `org.apache.hudi.hadoop.HoodieParquetInputFormat` を使用するテーブル。
- サフィックス `_rt` を持ち、スナップショットクエリを可能にする RealTime Serde `org.apache.hudi.hadoop.realtime.HoodieParquetRealtimeInputFormat` を使用するテーブル。

MongoDB と Amazon DocumentDB (MongoDB 互換)

MongoDB バージョン 3.2 以降がサポートされています。クローラーの実行時間を短縮するために、データの小さなサンプルのみをクロールするように選択できます。

リレーショナルデータベース

データベースのユーザー名とパスワードを使用して認証します。データベースエンジンのタイプに応じて、どのオブジェクト (データベース、スキーマ、テーブルなど) をクロールするかを選択できます。

Snowflake

Snowflake JDBC クローラーは、テーブル、外部テーブル、ビュー、マテリアライズドビューのクローリングをサポートしています。マテリアライズドビュー定義は入力されません。

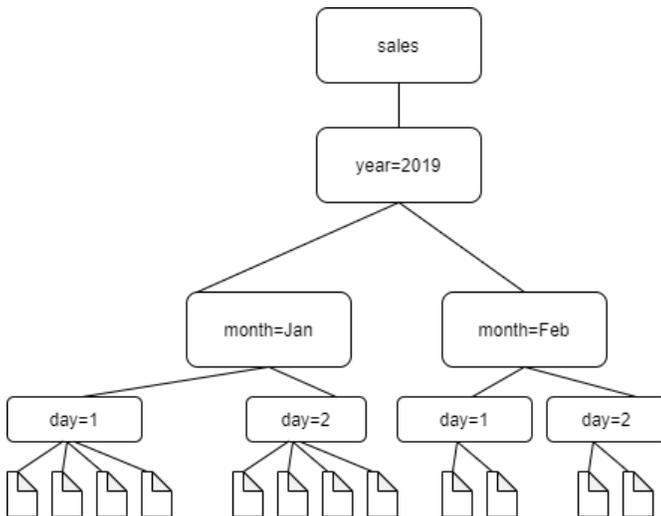
Snowflake 外部テーブルについては、Amazon S3 の場所を指している場合のみクローラーによってクローリングされます。クローラーは、テーブルスキーマの他に、Amazon S3 の場所、ファイル形式、および出力をデータカタログテーブルのテーブルパラメータとしてクローリングします。分割された外部テーブルのパーティション情報は入力されないことに注意してください。

Snowflake クローラーを使用して作成されたデータカタログテーブルでは、ETL は現在サポートされていません。

クローラーは、どのようにパーティションを作成するタイミングを判断していますか？

AWS Glue クローラーは、Amazon S3 をスキャンしてバケット内に複数のフォルダを検出すると、フォルダ構造のテーブルのルート、およびどのフォルダがテーブルのパーティションであるかを確認します。テーブルの名前は Amazon S3 プレフィックスまたはフォルダ名に基づいています。クロールするフォルダレベルを指すインクルードパスはユーザーが指定します。フォルダレベルの大半のスキーマが類似している場合、クローラーはテーブルを別個に作成せずに、テーブルのパーティションを作成します。クローラーで別個のテーブルを作成するには、クローラーを定義するときに各テーブルのルートフォルダを別個のデータストアとして追加します。

例えば、次の Amazon S3 フォルダ構造を考えてみます。



4つの最下位レベルのフォルダへのパスは次のとおりです。

```

S3://sales/year=2019/month=Jan/day=1
S3://sales/year=2019/month=Jan/day=2
S3://sales/year=2019/month=Feb/day=1
S3://sales/year=2019/month=Feb/day=2
  
```

クローラターゲットが Sales に設定されており、 $day=n$ フォルダ内のすべてのファイルが同じ形式 (暗号化されていない JSON など) で、同じか非常によく似たスキーマを持っているとします。クローラは、パーティションキー $year$ 、 $month$ 、および day で、4つのパーティションを持つ単一のテーブルを作成します。

次の例では、次の Amazon S3 構造を想定します。

```

s3://bucket01/folder1/table1/partition1/file.txt
s3://bucket01/folder1/table1/partition2/file.txt
s3://bucket01/folder1/table1/partition3/file.txt
s3://bucket01/folder1/table2/partition4/file.txt
s3://bucket01/folder1/table2/partition5/file.txt
  
```

$table1$ と $table2$ のファイルのスキーマが類似しており、クローラにインクルードパス $s3://bucket01/folder1/$ で定義されているデータストアが 1つの場合、クローラは 2つのパーティションキー列を持つ 1つのテーブルを作成します。最初のパーティションキー列には $table1$ および $table2$ が含まれ、2番目のパーティションキー列には、 $table1$ パーティションに対して $partition1$ から $partition3$ まで、および $table2$ パーティションに対して $partition4$ と $partition5$ が含まれます。2つの個別のテーブルを作成するには、2つのデータストアを持つ

クローラーを定義します。この例では、最初のインクルードパスを `s3://bucket01/folder1/table1/` として定義し、2 番目を `s3://bucket01/folder1/table2` として定義します。

Note

Amazon Athena の場合、各テーブルは Amazon S3 プレフィックス (すべてのオブジェクトを含む) に対応します。オブジェクト別にスキーマが異なる場合、Athena では同じプレフィックス内の異なるオブジェクトを別個のテーブルとして認識しません。これは、クローラーで同じ Amazon S3 プレフィックスから複数のテーブルを作成する場合に発生することがあります。そのため、Athena のクエリで結果が何も返されない場合があります。Athena でテーブルを適切に認識してクエリを実行するには、Amazon S3 フォルダ構造内の異なるテーブルスキーマごとに別個のインクルードパスを持つクローラーを作成します。詳細については、「[Athena で AWS Glue を使用するときのベストプラクティス](#)」およびこの [AWS ナレッジセンターの記事](#) を参照してください。

クローラーの前提条件

クローラーは、定義する時に指定する AWS Identity and Access Management (IAM) ロールのアクセス許可を取得します。この IAM ロールには、データストアからデータを抽出してデータカタログに書き込むためのアクセス許可が必要です。AWS Glue コンソールには、AWS Glue プリンシパルサービスの信頼ポリシーがアタッチされた IAM ロールだけがリスト表示されています。コンソールから、クローラーがアクセスする Amazon S3 データストアにアクセスするための IAM ポリシーを持つ IAM ロールを作成できます。AWS Glue のロールの指定の詳細については、「[AWS Glue のアイデンティティベースのポリシー](#)」を参照してください。

Note

Delta Lake データストアをクローラーする場合、Amazon S3 の場所に対する読み取り/書き込み権限が必要です。

クローラーには、ロールを作成して次のポリシーをアタッチできるようになりました。

- `AWSGlueServiceRole` AWS 管理ポリシー。データカタログに必要なアクセス許可を付与します。
- データソースに対するアクセス許可を付与するインラインポリシー。
- ロールに対する `iam:PassRole` アクセス許可を付与するインラインポリシー。

AWS Glue コンソールクローラーウィザードにロールを作成させるのが、より迅速なアプローチです。これによって作成されるロールは、クローラー専用であり、AWSGlueServiceRole AWS 管理ポリシーと、指定したデータソースに必要なインラインポリシーを含んでいます。

クローラーに既存のロールを指定する場合は、そのロールに AWSGlueServiceRole ポリシーまたは同等のポリシー (またはこのポリシーのスコープダウンバージョン) と、必要なインラインポリシーが含まれていることを確認します。例えば、Amazon S3 データストアの場合、インラインポリシーは少なくとも次のようになります。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "s3:PutObject"
      ],
      "Resource": [
        "arn:aws:s3:::bucket/object*"
      ]
    }
  ]
}
```

Amazon DynamoDB データストアの場合、ポリシーは、少なくとも次のようになります。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "dynamodb:DescribeTable",
        "dynamodb:Scan"
      ],
      "Resource": [
        "arn:aws:dynamodb:region:account-id:table/table-name*"
      ]
    }
  ]
}
```

また、クローラーが AWS Key Management Service (AWS KMS) 暗号化 Amazon S3 データを読み取る場合、IAM ロールには AWS KMS キーの復号化のアクセス許可が必要です。詳細については、[ステップ 2: AWS Glue 用の IAM ロールを作成する](#) を参照してください。

クローラーの設定

クローラーはデータストアにアクセスし、メタデータを抽出して、テーブル定義を AWS Glue Data Catalog に作成します。コンソールの [Crawlers (クローラー)AWS Glue] ペインには、作成したクローラーがすべて一覧表示されます。リストには、クローラーの最後の実行のステータスとメトリクスが表示されます。

Note

独自の JDBC ドライバーバージョンを導入する場合、AWS Glue クローラーは AWS Glue ジョブと Amazon S3 バケットのリソースを使用して、用意したドライバーが自分の環境で実行されるようにします。リソースの追加使用量はアカウントに反映されます。さらに、独自の JDBC ドライバーを用意しても、クローラーがドライバーの機能をすべて活用できるわけではありません。ドライバーは、「[AWS Glue 接続の追加](#)」に記載されているプロパティに限定されます。

クローラーを設定するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。ナビゲーションペインで [Crawlers (クローラー)] を選択します。
2. [クローラーの作成] を選択して、[クローラーの追加] ウィザードの手順に従います。ウィザードは、クローラーの作成に必要なステップをガイドします。スキーマを定義するためにカスタム分類子を追加する場合は、「[AWS Glue でのクローラーへの分類子の追加](#)」を参照してください。

ステップ 1: クローラープロパティを設定する

クローラーの名前と説明 (オプション) を入力します。必要に応じて、クローラーに [タグキー] とオプションの [タグ値] でタグを付加できます。作成後、タグキーは読み取り専用になります。リソースを整理、識別しやすいように、いくつかのリソースでタグを使用します。詳細については、「AWS Glue の AWS タグ」を参照してください。

名前

名前に使用できるのは、文字 (A~Z)、数字 (0~9)、ハイフン (-)、またはアンダースコア (_) で、長さは 255 文字までです。

説明

説明の長さは、最大 2048 文字です。

タグ

タグを使用して、リソースの整理と識別を行います。詳細については、次を参照してください:

- [AWSAWS Glue のタグ](#)

ステップ 2: データソースと分類子を選択する

データソースの設定

[データは AWS Glue テーブルにマッピングされていますか] で、適切なオプション (「未実施」または「はい」) を選択します。デフォルトでは、「未実施」が選択されています。

クローラーは、クローリングのソースとして直接データストアにアクセスでき、また既存のデータカタログのテーブルをソースとして使用することもできます。クローラーが既存のカタログテーブルを使用する場合、それらのカタログのテーブルで指定されたデータストアをクローリングします。

- [Not yet] (未実施): クローリングする 1 つまたは複数のデータソースを選択します。クローラーは、さまざまなタイプ (Amazon S3、JDBC など) の複数のデータストアをクローリングできます。

一度に設定できるデータストアは 1 つだけです。接続情報とインクルードパスと除外パターンを指定すると、別のデータストアを追加することもできます。

- [Yes] (はい): AWS Glue データカタログから既存のテーブルを選択します。カタログテーブルは、クローリングするデータストアを指定します。クローラーは 1 回の実行でのカタログテーブルのみをクローリングできます。他のソースタイプを混在させることはできません。

ソースとしてカタログテーブルを指定する一般的な理由は、(データストアの構造をすでに知っていたため) テーブルを手動で作成した場合に、新しいパーティションの追加を含め、クローラーにテーブルを更新された状態に保持させるためです。他の理由の説明については、「[クローラーを使用して手動で作成されたデータカタログテーブルを更新する](#)」を参照してください。

既存のテーブルをクローラーソースのタイプとして指定する場合は、次の条件が適用されます。

- データベース名はオプションです。
- Amazon S3 または Amazon DynamoDB データストアを指定するカタログテーブルのみが使用できます。
- クローラーが実行されるとき、新しいカタログテーブルは作成されません。既存のテーブルは、新しいパーティションの追加を含め、必要に応じて更新されます。
- データストアで見つかる削除されたオブジェクトは無視されます。カタログテーブルが削除されません。代わりに、クローラーによりログメッセージが書き込まれます (SchemaChangePolicy.DeleteBehavior=LOG)
- Amazon S3 パスごとに単一のスキーマを作成するクローラー設定オプションはデフォルトで有効になっており、無効にすることはできません。(TableGroupingPolicy=CombineCompatibleSchemas) 詳しくは、「[各 Amazon S3 インクルードパスの単一のスキーマを作成する方法](#)」を参照してください。
- カタログテーブルをソースとして他のソースタイプ (例: Amazon S3 または Amazon DynamoDB) と混在させることはできません。

データソース

クローラーでスキャンするデータソースのリストを選択または追加します。

(オプション) データソースとして JDBC を選択した場合、ドライバー情報が保存されている Connection アクセスを指定するときに、独自の JDBC ドライバーを使用できます。

インクルードパス

クローラーで何を含め、何を除外するかを評価する際、クローラーは必要なインクルードパスを評価することから始めます。Amazon S3、MongoDB、MongoDB Atlas、Amazon DocumentDB (MongoDB 互換)、リレーショナルデータストアの場合は、インクルードパスを指定する必要があります。

Amazon S3 データストアの場合

このアカウントのパスを指定するか、別のアカウントのパスを指定するかを選択し、Amazon S3 パスを参照して選択します。

Amazon S3 データストアの場合、インクルードパスの構文は bucket-name/folder-name/file-name.ext です。バケット内のすべてのオブジェクトをクローラーするには、インクルードパスにバケット名のみ指定します。除外パターンは、インクルードパスを基準とする相対パスです。

Delta Lake データストアの場合

Delta テーブルへの Amazon S3 パスを `s3://bucket/prefix/object` として 1 つ以上指定します。

Iceberg または Hudi データストアの場合

Iceberg または Hudi テーブルのメタデータを持つフォルダを含む Amazon S3 パスを `s3://bucket/prefix` として 1 つ以上指定します。

Hudi データストアの場合、Hudi フォルダはルートフォルダの子フォルダ内に存在する場合があります。クローラーは、Hudi フォルダのパス以下にあるすべてのフォルダをスキャンします。

JDBC データストアの場合

データベース製品に応じて、`<database>/<schema>/<table>` または `<database>/<table>` を入力します。Oracle Database と MySQL は、パス内のスキーマをサポートしません。`<schema>` または `<table>` は、パーセント (%) 文字に置き換えることができます。たとえば、システム識別子 (SID) が `orcl` の Oracle データベースの場合、`orcl/%` を入力して、接続で指定されたユーザーがアクセスできるすべてのテーブルをインポートします。

Important

このフィールドでは、大文字と小文字が区別されます。

MongoDB、MongoDB Atlas、Amazon DocumentDB データストアの場合

`database/collection` と入力します。

MongoDB、MongoDB Atlas、Amazon DocumentDB (MongoDB 互換) の場合、構文は `database/collection` です。

JDBC データストアの場合、構文は `database-name/schema-name/table-name` または `database-name/table-name` です。構文は、データベースエンジンでデータベース内のスキーマがサポートされているかどうかによって依存します。たとえば、MySQL や Oracle などのデータベースエンジンの場合は、インクルードパスに `schema-name` を指定しません。インクルードパスでスキーマやテーブルの代わりにパーセント記号 (%) を使用することで、データベース内のすべてのスキーマやテーブルを表すことができます。インクルードパスでデータベースの代わりにパーセント記号 (%) を使用することはできません。

最大トラバーサル深度 (Iceberg または Hudi データストアのみ)

クローラーが Amazon S3 パス内の Iceberg または Hudi メタデータフォルダを検出するために通過できる Amazon S3 パスの最大深度を定義します。このパラメータの目的は、クローラーの実行時間を制限することです。デフォルト値は 10 で、最大値は 20 です。

除外パターン

この機能を利用すると、特定のファイルまたはテーブルをクロールから除外できます。エクスクルードパスは、インクルードパスを基準とする相対パスです。たとえば、JDBC データストア内のテーブルを除外するには、エクスクルードパスにテーブル名を入力します。

クローラーでは、JDBC データストアに接続するために、JDBC URI 接続文字列を含む AWS Glue 接続を使用します。クローラーは、データベースエンジン内のオブジェクトにのみアクセスできます。そのために、AWS Glue 接続で JDBC ユーザー名とパスワードを使用します。クローラーは、JDBC 接続を介してアクセスできるテーブルのみ作成できます。クローラーは、JDBC URI を使用してデータベースエンジンにアクセスした後で、インクルードパスを使用してデータカタログで作成するデータベースエンジン内のテーブルを決定します。例えば、MySQL の場合、MyDatabase/% というインクルードパスを指定すると、MyDatabase 内のすべてのテーブルがデータカタログで作成されます。Amazon Redshift にアクセスする場合、MyDatabase/% というインクルードパスを指定すると、データベース MyDatabase の各スキーマ内のすべてのテーブルがデータカタログで作成されます。MyDatabase/MySchema/% のインクルードパスを指定すると、データベース MyDatabase のすべてのテーブルとスキーマ MySchema が作成されます。

インクルードパスの指定後に、1 つ以上の Unix 形式の glob 除外パターンを指定することで、インクルードパスに含まれる予定であったオブジェクトをクロールから除外できます。これらのパターンはインクルードパスに適用されて、どのオブジェクトを除外するか決定します。また、これらのパターンはクローラーによって作成されるテーブルのプロパティとして保存されます。AWS GluePySpark 拡張機能 (`create_dynamic_frame.from_catalog` など) は、テーブルプロパティを読み取り、除外パターンによって定義されたオブジェクトを除外します。

AWS Glue は、除外パターンで次の種類の glob パターンをサポートしています。

除外パターン	説明
*.csv	.csv で終わる現在のフォルダ内のオブジェクト名を表す Amazon S3 パスと一致する

除外パターン	説明
.	ドットを含むオブジェクト名すべてと一致する
*.{csv,avro}	.csv か .avro で終わるオブジェクト名と一致する
foo.?	foo. で始まり、その後に 1 文字の拡張子が続くオブジェクト名と一致する
myfolder/*	/myfolder/mysource など、myfolder のサブフォルダの 1 つのレベルにあるオブジェクトと一致する
myfolder/**	/myfolder/mysource/data など、myfolder のサブフォルダの 2 つのレベルにあるオブジェクトと一致する
myfolder/**	myfolder のすべてのサブフォルダにあるオブジェクト (/myfolder/mysource/mydata や /myfolder/mysource/data など) と一致する
myfolder**	/myfolder や /myfolder/mydata.txt などの myfolder 以下のファイルと同様に、サブフォルダ myfolder と一致する
Market*	JDBC データベースの Market で始まる名前のテーブル (Market_us や Market_fr など) と一致する

AWS Glue は、glob 除外パターンを次のように解釈します。

- スラッシュ (/) 文字は、Amazon S3 キーをフォルダ階層に区切る区切り記号です。
- アスタリスク (*) 記号は、フォルダの境界を超えない、0 文字以上の名前の要素と一致します。

- 二重アスタリスク (**) は、フォルダやスキーマの境界を越える 0 個以上の文字に相当します。
- 疑問符 (?) 記号は、名前の要素のちょうど 1 文字に相当します。
- バックスラッシュ (\) 文字は、本来ならば特殊文字として解釈される文字をエスケープ処理するために使用されます。\\ 式はバックスラッシュ 1 つに相当し、\{ は左括弧に相当します。
- 角括弧 ([]) は、一連の文字の中から、名前の要素の 1 文字に相当する角括弧式を作成します。たとえば、[abc] は a、b、または c に一致します。ハイフン (-) は、範囲を指定するために使用されます。つまり、[a-z] は a から z (この値を含みます) までに相当する範囲を指定します。これらのフォームは組み合わせることができます。そのため、[abce-g] は a、b、c、e、f、または g に一致します。角括弧 ([]) の後の文字が感嘆符 (!) の場合、角括弧式は否定の意味になります。たとえば、[!a-c] は a、b、または c 以外のすべての文字に一致します。

角括弧式内では、*、?、および \ 文字は、文字通りの意味です。ハイフン (-) 文字は、角括弧内で最初の文字だった場合、または式を否定する ! の次の文字だった場合は、文字通りの意味です。

- 中括弧 ({ }) は、グループ内のサブパターンが一致する場合にグループが一致するサブパターンのグループを囲みます。カンマ (,) 文字は、サブパターンを分割するために使用されます。グループはネストできません。
- ファイル名の先頭のピリオドまたはドット文字は、マッチ操作では通常の文字として扱われます。たとえば、* 除外パターンは、ファイル名 .hidden に一致します。

Example Amazon S3 除外パターンの例

各除外パターンは、インクルードパスに対して評価されます。例えば、次の Amazon S3 デイレクトリ構造があるとします。

```
/mybucket/myfolder/  
  departments/  
    finance.json  
    market-us.json  
    market-emea.json  
    market-ap.json  
  employees/  
    hr.json  
    john.csv  
    jane.csv  
    juan.txt
```

インクルードパスが `s3://mybucket/myfolder/` の場合、以下は除外パターンのサンプル結果の一部です。

除外パターン	結果
<code>departments/**</code>	<code>departments</code> フォルダ内のすべてのファイルとフォルダを除外し、 <code>employees</code> フォルダとそのファイルを含めます。
<code>departments/market*</code>	<code>market-us.json</code> 、 <code>market-emea.json</code> 、および <code>market-ap.json</code> を除外します
<code>** .csv</code>	名前が <code>.csv</code> で終わる <code>myfolder</code> 以下のすべてのオブジェクトを除外します
<code>employees/*.csv</code>	<code>employees</code> フォルダ内のすべての <code>.csv</code> ファイルを除外します。

Example Amazon S3 パーティションのサブセットの除外の例

データが日別に分割されていて、1年の日別に異なる Amazon S3 パーティションに入っているとします。2015年1月には、31のパーティションがあります。ここで、1月の第1週のみデータをクローリングするには、1日〜7日を除くすべてのパーティションを除外する必要があります。

```
2015/01/{[!0],0[8-9]}**, 2015/0[2-9]**, 2015/1[0-2]**
```

この glob パターンの各パートを見てみます。最初のパートは `2015/01/{[!0],0[8-9]}**` で、"0" で始まらないすべての日付、および 2015年01月の08日目および09日目を除外しています。"" を日数パターンのサフィックスとして使用すると、下位レベルフォルダへのフォルダ境界を越えることに注意してください。"" を使用すると、下位レベルフォルダは除外されます。

2番目のパートは `2015/0[2-9]**` で、2015年02月から09月までの日を除外します。

3番目のパートは `2015/1[0-2]**` で、2015年10、11、12月の日を除外します。

Example JDBC の除外パターン

次のスキーマ構造を使用して JDBC データベースをクローリングしているとします。

```
MyDatabase/MySchema/
  HR_us
  HR_fr
  Employees_Table
  Finance
  Market_US_Table
  Market_EMEA_Table
  Market_AP_Table
```

インクルードパスが MyDatabase/MySchema/% の場合、以下は除外パターンのサンプル結果の一部です。

除外パターン	結果
HR*	HR で始まる名前を持つテーブルを除外
Market_*	Market_ で始まる名前を持つテーブルを除外
**_Table	_Table で終わる名前を持つテーブルをすべて除外

追加のクローラーソースパラメータ

各ソースのタイプには、異なる追加パラメータのセットが必要です。次のリストはその一部です。

Connection

AWS Glue 接続を選択または追加します。接続の詳細については、「[データへの接続](#)」を参照してください

追加メタデータ - オプション (JDBC データストア用)

クローラーがクローリングするための追加のメタデータプロパティを選択します。

- [Comments] (コメント): 関連するテーブルレベルと列レベルのコメントをクローリングします。

- [Raw types] (未加工型): 表の列の未加工のデータ型を追加のメタデータに保持します。デフォルトの動作として、クローラーは未処理のデータ型を Hive 互換の型に変換します。

JDBC ドライバークラス名 - オプション (JDBC データストア用)

クローラーがデータソースに接続するためのカスタム JDBC ドライバークラス名を入力します。

- Postgres: org.postgresql.Driver
- MySQL: com.mysql.jdbc.Driver、com.mysql.cj.jdbc.Driver
- Redshift: com.amazon.redshift.jdbc.Driver、com.amazon.redshift.jdbc42.Driver
- Oracle: oracle.jdbc.driver.OracleDriver
- SQL Server: com.microsoft.sqlserver.jdbc.SQLServerDriver

JDBC ドライバース3パス - オプション (JDBC データストア用)

.jar ファイルへの既存の Amazon S3 パスを選択します。クローラーがデータソースに接続するためのカスタム JDBC ドライバーを使用するときに、この場所に .jar ファイルが保存されません。

データサンプリングを有効にする (Amazon DynamoDB、MongoDB、MongoDB Atlas、Amazon DocumentDB データストアのみ)

データサンプルのみをクロールするかどうかを選択します。選択しない場合は、テーブル全体がクロールされます。テーブルが高スループットテーブルではない場合、すべてのレコードのスキャンには時間がかかることがあります。

クエリ用のテーブルを作成する (Delta Lake データストア用のみ)

Delta Lake テーブルの作成方法を次の中から選択します。

- [Create Native tables] (ネイティブテーブルの作成): Delta トランザクションログの直接クエリをサポートするクエリエンジンとの統合を可能にします。
- [Create Symlink tables] (Symlink テーブルの作成): 指定された設定パラメータに基づいて、マニフェストファイルを分割キーで分割して Symlink マニフェストを作成します。

スキャンレート - オプション (DynamoDB データストア用のみ)

クローラーで使用する DynamoDB テーブル読み取りキャパシティーユニットの割合を指定します。読み取りキャパシティーユニットは、DynamoDB で定義されている用語で、テーブルに対して実行できる読み取り回数/秒のレート制限として機能する数値です。0.1~1.5 の値を入力します。指定しない場合、プロビジョニングされたテーブルの場合は 0.5%、オンデマンドテーブ

ルの場合は最大の設定済みキャパシティーの 1/4 にデフォルト設定されます。AWS Glue クローラーは、プロビジョニング容量モードでのみ使用できることにご注意ください。

 Note

DynamoDB データストアの場合、テーブルの読み取りと書き込みを処理するプロビジョニング容量モードを設定します。AWS Glue クローラーは、オンデマンド容量モードでは使用しないでください。

ネットワーク接続 - オプション (Amazon S3 データストア用のみ)

必要に応じて、この Amazon S3 ターゲットで使用するネットワーク接続を含めます。各クローラーは 1 つのネットワーク接続に制限されているため、他の Amazon S3 ターゲットも同じ接続を使用します (空欄の場合は接続なし)。

接続の詳細については、「[データへの接続](#)」を参照してください

ファイルのサブセットのみのサンプリングとサンプルサイズ (Amazon S3 データストア用のみ)

データセット内のサンプルファイルをクローリングするとき、クローリングされる各リーフフォルダ内のファイル数を指定します。この機能をオンにすると、このデータセット内のすべてのファイルをクローリングする代わりに、クローラーはクローリングする各リーフフォルダ内の一部のファイルをランダムに選択します。

サンプリングするクローラーは、データ形式について事前に知識があり、フォルダ内のスキーマが変更されないことを知っているお客様に最適です。この機能をオンにすると、クローラーの実行時間が大幅に短縮されます。

有効な値は、1 から 249 までの整数です。指定しない場合、すべてのファイルがクローリングされます。

以降のクローラー実行

このフィールドは、すべての Amazon S3 データソースに影響するグローバルフィールドです。

- [Crawl all sub-folders] (すべてのサブフォルダをクローリング): 以降クローリングするたびに、すべてのフォルダを再度クローリングします。
- [Crawl new sub-folders only] (新しいサブフォルダのみをクローリング): 最近のクローリング以降に追加された Amazon S3 フォルダのみをクローリングします。スキーマに互換性がある場

合、新しいパーティションが既存のテーブルに追加されます。詳細については、「[the section called “新しいパーティションを追加するための増分クローリング”](#)」を参照してください。

- [Crawl based on events] (イベントに基づいてクローリング): Amazon S3 イベントを利用して、クローリングするフォルダを制御します。詳細については、「[the section called “Amazon S3 イベント通知を使用した加速クローリング”](#)」を参照してください。

カスタム分類子 - オプション

クローラーを定義する前にカスタム分類子を定義します。分類子は、指定されたファイルがクローラーで処理可能な形式であるかどうかをチェックします。処理可能な場合、分類子は、そのデータ形式と一致する StructType オブジェクトの形式でスキーマを作成します。

詳細については、「[AWS Glue でのクローラーへの分類子の追加](#)」を参照してください。

ステップ 3: セキュリティ設定を構成する

IAM ロール

クローラーはこのロールを取得します。これは、AWS 管理ポリシー `AWSGlueServiceRole` と同様のアクセス許可を持っている必要があります。Amazon S3 および DynamoDB ソースの場合、データストアへのアクセス許可も必要です。クローラーが AWS Key Management Service (AWS KMS) で暗号化された Amazon S3 データを読み取る場合は、ロールには AWS KMS キーでの復号化のアクセス許可が必要です。

Amazon S3 データストアの場合、ロールにアタッチされた追加のアクセス許可は次のようになります。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "s3:PutObject"
      ],
      "Resource": [
        "arn:aws:s3:::bucket/object*"
      ]
    }
  ]
}
```

```
}

```

Amazon DynamoDB データストアの場合、ロールにアタッチされた追加のアクセス許可は次のようになります。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "dynamodb:DescribeTable",
        "dynamodb:Scan"
      ],
      "Resource": [
        "arn:aws:dynamodb:region:account-id:table/table-name*"
      ]
    }
  ]
}
```

自身の JDBC ドライバーを追加するには、アクセス許可を新たに追加する必要があります。

- ジョブアクション CreateJob、DeleteJob、GetJob、GetJobRun、StartJobRun にアクセス権限を付与します。
- Amazon S3 アクション
s3:DeleteObjects、s3:GetObject、s3:ListBucket、s3:PutObject にアクセス権限を付与します。

Note

Amazon S3 バケットポリシーが無効になっている場合、s3:ListBucket は必要ありません。

- サービスプリンシパルに Amazon S3 ポリシーのバケット/フォルダへのアクセス権限を付与します。

Amazon S3 ポリシーの例:

```
{
  "Version": "2012-10-17",
```

```
"Statement": [
  {
    "Sid": "VisualEditor0",
    "Effect": "Allow",
    "Action": [
      "s3:PutObject",
      "s3:GetObject",
      "s3:ListBucket",
      "s3:DeleteObject"
    ],
    "Resource": [
      "arn:aws:s3:::bucket-name/driver-parent-folder/driver.jar",
      "arn:aws:s3:::bucket-name"
    ]
  }
]
```

AWS Glue は以下のフォルダ (`_crawler` と `_glue_job_crawler`) を Amazon S3 バケットに JDBC ドライバーと同じレベルで作成します。例えば、ドライバーパスが `<s3-path/driver_folder/driver.jar>` の場合、以下のフォルダがまだ存在しない場合は作成されません。

- `<s3-path/driver_folder/_crawler>`
- `<s3-path/driver_folder/_glue_job_crawler>`

必要に応じて、クローラーにセキュリティ設定を追加して、保管時の暗号化オプションを指定することができます。

詳細については、[ステップ 2: AWS Glue 用の IAM ロールを作成する](#) および [AWS Glue の Identity and Access Management](#) を参照してください。

Lake Formation 設定 - オプション

クローラーが Lake Formation の認証情報をデータソースのクローリングに使用できるようにします。

[Use Lake Formation credentials for crawling S3 data source] (S3 データソースのクローリングに Lake Formation 認証情報を使用する) をオンにすると、クローラーは Lake Formation 認証情報を使用してデータソースをクローリングできるようになります。データソースが別のアカウントに属する場合は、登録されているアカウント ID を指定する必要があります。そうしないと、ク

クローラーはアカウントに関連付けられているデータソースのみをクローリングします。Amazon S3 とデータカタログのデータソースにのみ適用されます。

セキュリティ設定 - オプション

設定にはセキュリティ設定が含まれます。詳細については、次を参照してください:

- [AWS Glue によって書き込まれたデータの暗号化](#)

Note

クローラーにセキュリティ設定を設定すると、変更できますが、削除することはできません。クローラーのセキュリティレベルを下げるには、設定内でセキュリティ機能を明示的に DISABLED に設定するか、新しいクローラーを作成します。

ステップ 4: 出力とスケジュールを設定する

出力設定

オプションには、検出されたスキーマの変更、データストア内の削除されたオブジェクトをクローラーが処理する方法があります。詳細については、「[クローラーの動作のカスタマイズ](#)」を参照してください。

クローラースケジュール

AWS Glue では、オンデマンドでクローラーを実行することも、クローラーとジョブの時間ベースのスケジュールを定義することもできます。これらのスケジュールの定義は、Unix と同様の cron 構文を使用します。詳細については、「[AWS Glue クローラのスケジュール](#)」を参照してください。

ステップ 5: 確認して作成する

設定したクローラー設定を確認して、クローラーを作成します。

AWS Glue でのクローラーへの分類子の追加

分類子はデータをデータストアに読み取ります。データの形式を認識すると、スキーマが生成されます。分類子も、形式の認識がどれほど確実かを示す確信度数を返します。

AWS Glue では一連の組み込み分類子が用意されていますが、カスタム分類子を作成することもできます。AWS Glue は、クローラ定義で指定した順序で、カスタム分類子を最初に呼び出します。カス

タム分類子から返された結果に応じて、AWS Glue が組み込みの分類子を呼び出す場合もあります。処理中に分類子が `certainty=1.0` を返した場合、正しいスキーマを 100% 確実に作成できることを示しています。次に、AWS Glue はその分類子の出力を使用します。

分類子が `certainty=1.0` を返さない場合、AWS Glue は最も高い確実性を持つ分類子の出力を使用します。どの分類子からも `0.0` 以上の確実性が返されない場合、AWS Glue は UNKNOWN のデフォルト分類文字列を返します。

分類子を使用するタイミング

データストアをクローリングして AWS Glue Data Catalog でメタデータテーブルを定義する際に分類子を使用します。順序が設定された一連の分類子を使用してクローラをセットアップできます。クローラが分類子を呼び出す際、分類子はデータが認識されるかどうかを判断します。分類子でデータを認識できないか、100% 確実ではない場合、クローラはリストにある次の分類子を呼び出して、データを認識できるかどうか判断します。

AWS Glue コンソールを使用して分類子を作成する方法の詳細については、「[AWS Glue コンソールでの分類子の操作](#)」を参照してください。

カスタム分類子

分類子の出力には、ファイルの分類や形式 (たとえば、`json`)、およびファイルのスキーマを示す文字列が含まれます。カスタム分類子の場合、分類子のタイプに基づいてスキーマを作成するためのロジックを定義します。分類子のタイプには、`grok` パターン、XML タグ、および JSON パスに基づくスキーマの定義が含まれています。

分類子の定義を変更すると、変更前の分類子を使用してクローリングしたデータは再分類されません。クローラは、以前にクローリングしたデータを追跡します。新しいデータは、更新された分類子で分類されるため、スキーマが更新される場合があります。データのスキーマが更新された場合は、クローラの実行時に分類子を更新してスキーマの変更を反映してください。データを再分類して不正な分類子を修正するには、更新された分類子を使用して新しいクローラを作成します。

AWS Glue でカスタム分類子を作成する方法については、「[カスタム分類子の書き込み](#)」を参照してください。

Note

組み込み分類子のいずれかでデータ形式が認識される場合、カスタム分類子を作成する必要はありません。

AWS Glue の組み込み分類子

AWS Glue は、JSON、CSV、ウェブログ、および多くのデータベースシステムを含む、さまざまな形式の組み込み分類子を提供します。

AWS Glue が入力データ形式に適したカスタム分類子を 100% の確実度で検出できない場合、次の表に示すような順番で組み込み分類子を呼び出します。組み込み分類子は、形式が一致するか (certainty=1.0)、または一致しないか (certainty=0.0) どうかを示す結果を返します。certainty=1.0 を持つ最初の分類子は、Data Catalog での分類文字列とメタデータテーブルのスキーマを提供します。

分類子タイプ	分類文字列	メモ
Apache Avro	avro	ファイルの先頭からスキーマを読み取って形式を判断します。
Apache ORC	orc	ファイルのメタデータを読み取って形式を判断します。
Apache Parquet	parquet	ファイルの末尾からスキーマを読み取って形式を判断します。
JSON	json	ファイルの先頭から読み取って形式を判断します。
バイナリ JSON	bson	ファイルの先頭から読み取って形式を判断します。
XML	xml	<p>ファイルの先頭から読み取って形式を判断します。AWS Glue は、ドキュメントの XML タグに基づいてテーブルスキーマを判定します。</p> <p>カスタム XML を作成してドキュメントの行を指定するには、「XML カスタム分類子の書き込み」を参照してください。</p>
Amazon Ion	ion	ファイルの先頭から読み取って形式を判断します。
Combined Apache ログ	combined_apache	grok パターンを通じてログ形式を判断します。
Apache ログ	apache	grok パターンを通じてログ形式を判断します。

分類子タイプ	分類文字列	メモ
Linux カーネルログ	linux_kernel	grok パターンを通じてログ形式を判断します。
Microsoft ログ	microsoft_log	grok パターンを通じてログ形式を判断します。
Ruby ログ	ruby_logger	ファイルの先頭から読み取って形式を判断します。
Squid 3.x ログ	squid	ファイルの先頭から読み取って形式を判断します。
Redis 監視ログ	redismonlog	ファイルの先頭から読み取って形式を判断します。
Redis ログ	redislog	ファイルの先頭から読み取って形式を判断します。
CSV	csv	次の区切り記号をチェックします。カンマ (,)、パイプ ()、タブ (\t)、セミコロン (;)、および Ctrl-A (\u0001)。Ctrl-A は Start Of Heading の Unicode 制御文字です。
Amazon Redshift	redshift	JDBC 接続を使用してメタデータをインポートします。
MySQL	mysql	JDBC 接続を使用してメタデータをインポートします。
PostgreSQL	postgresql	JDBC 接続を使用してメタデータをインポートします。
Oracle データベース	oracle	JDBC 接続を使用してメタデータをインポートします。
Microsoft SQL Server	sqlserver	JDBC 接続を使用してメタデータをインポートします。
Amazon DynamoDB	dynamodb	DynamoDB テーブルからデータを読み取ります。

以下の圧縮形式のファイルは分類できます。

- ZIP (1つのファイルのみを含むアーカイブでサポートされています)。Zip は他のサービスで十分にサポートされていないことに注意してください (アーカイブのため)。
- BZIP
- GZIP
- LZ4
- Snappy (標準および Hadoop ネイティブの両方の Snappy フォーマットでサポート)

組み込みの CSV 分類子

組み込みの CSV 分類子では、CSV ファイルの内容を解析して AWS Glue テーブルのスキーマを判断します。この分類子は以下の区切り記号を確認します。

- カンマ (,)
- パイプ (|)
- タブ (\t)
- セミコロン (;)
- Ctrl-A (\u0001)

Ctrl-A は Start Of Heading の Unicode 制御文字です。

CSV として分類されるためには、テーブルのスキーマに少なくとも 2 つのデータ列と 2 つのデータ行が必要です。CSV 分類子では、いくつかのヒューリスティックを使用して特定のファイルにヘッダーがあるかどうかを判断します。分類子で最初のデータ行にヘッダーを確認できない場合は、列のヘッダーが col1、col2、col3 のように表示されます。組み込みの CSV 分類子では、以下のファイルの特性を評価することで、ヘッダーを推測するかどうかを決めます。

- ヘッダー候補の各列が STRING データ型として解析されます。
- 最後の列を除き、ヘッダー候補の列ごとに 150 文字未満のコンテンツがあります。末尾の区切り記号を許可するには、ファイル全体で最後の列を空にすることができます。
- ヘッダー候補の各列が、列名に関する AWS Glue regex 要件を満たす必要があります。
- ヘッダー行は、データ行と十分に異なっている必要があります。これを判断するには、1 つ以上の行が STRING 型以外として解析されることを確認します。すべての列が STRING 型である場合、最初のデータ行は以降の行と十分に異なっていないため、ヘッダーとして使用できません。

Note

組み込みの CSV 分類子で必要な AWS Glue テーブルが作成されない場合は、以下のいずれかの代替方法を使用できます。

- Data Catalog で列名を変更し、SchemaChangePolicy を LOG に設定して、将来のクローラ実行に関してパーティションの出力設定を InheritFromTable に設定します。
- データを分類するためのカスタム grok 分類子を作成し、必要な列を割り当てます。
- 組み込みの CSV 分類子では、LazySimpleSerDe をシリアル化ライブラリとして参照するテーブルを作成します。これは、型の推定に適しています。ただし、CSV データ内に引用符で囲まれた文字列がある場合は、テーブル定義を編集して SerDe ライブラリを OpenCSVSerDe に変更します。推定した型を STRING に調整し、SchemaChangePolicy を LOG に設定して、将来のクローラ実行に関してパーティションの出力設定を InheritFromTable に設定します。SerDe ライブラリの詳細については、Amazon Athena ユーザーガイドの「[SerDe リファレンス](#)」を参照してください。

カスタム分類子の書き込み

AWS Glue でデータを分類するためのカスタム分類子を提供できます。grok パターン、XML タグ、JavaScript Object Notation (JSON)、またはカンマ区切り値 (CSV) を使用してカスタム分類子を作成できます。AWS Glue クローラがカスタム分類子を呼び出します。分類子がデータを認識すると、データの分類とスキーマがクローラに返されます。組み込みの分類子にデータが一致しない場合、または、クローラにより作成されたテーブルをカスタマイズする場合は、カスタム分類子を定義する必要があるかもしれません。

AWS Glue コンソールを使用して分類子を作成する方法の詳細については、「[AWS Glue コンソールでの分類子の操作](#)」を参照してください。

AWS Glue は、組み込みの分類子の前に、指定した順序でカスタム分類子を実行します。クローラがデータに一致する分類子を検出すると、分類の文字列とスキーマが AWS Glue Data Catalog に書き込まれるテーブルの定義で使用されます。

トピック

- [Grok カスタム分類子の書き込み](#)
- [XML カスタム分類子の書き込み](#)
- [JSON カスタム分類子の書き込み](#)

• CSV カスタム分類子の書き込み

Grok カスタム分類子の書き込み

Grok は、一致するパターンによりテキストデータを解析するために使用するツールです。grok パターンは名前のついた一連の正規表現 (regex) で、一度に 1 行のデータごとに一致させるために使用されます。AWS Glue は grok パターンを使用してデータのスキーマを推測します。grok パターンがデータと一致すると、AWS Glue はそのパターンを使用してデータの構造を判断し、フィールドにマッピングします。

AWS Glue は数多くの組み込みパターンを提供します。または、独自のパターンを定義することもできます。組み込みパターンとカスタム分類子の定義にあるカスタムパターンを使用して grok パターンを作成できます。grok パターンをカスタマイズしてカスタムテキストファイル形式を分類できます。

Note

AWS Glue grok カスタム分類子は、AWS Glue Data Catalog で作成されたテーブルの GrokSerDe シリアル化ライブラリを使用します。AWS Glue Data Catalog で Amazon Athena、Amazon EMR、または Redshift Spectrum を使用している場合、GrokSerDe のサポートに関する情報については、これらのサービスに関するドキュメントを確認してください。現在、Amazon EMR と Redshift Spectrum の GrokSerDe を使用して作成されたテーブルのクエリを実行する際に問題が発生することがあります。

grok パターンのコンポーネントの基本的な構文を以下に示します。

```
%{PATTERN:field-name}
```

名付けられた PATTERN に一致するデータは、スキーマの field-name 列にデフォルトのデータ型 string でマッピングされます。必要に応じて、フィールドのデータ型は結果のスキーマの byte、boolean、double、short、int、long、または float にキャストできます。

```
%{PATTERN:field-name:data-type}
```

たとえば、num フィールドを int データ型にキャストするには、以下のパターンを使用することができます。

```
%{NUMBER:num:int}
```

パターンは他のパターンで構成できます。たとえば、SYSLOG タイムスタンプのパターンを月、日、時間のパターン (Feb 1 06:25:43 など) で定義できます。このデータの場合、次のパターンを定義できます。

```
SYSLOGTIMESTAMP %{MONTH} +%{MONTHDAY} %{TIME}
```

Note

grok パターンは、一度に 1 行ずつしか処理できません。複数行のパターンはサポートされていません。また、パターン内の改行はサポートされていません。

AWS Glue のカスタム分類子の値

grok 分類子を定義する場合、AWS Glue に以下の値を指定しカスタム分類子を作成します。

名前

分類子の名前。

分類

分類されたデータの形式を説明するために記述されたテキスト文字列 (例: special-logs)。

Grok パターン

データストアに適用される一連のパターンで、一致があるかどうかを決定します。これらのパターンは AWS Glue の [組み込みパターン](#) と定義されたカスタムパターンによるものです。

grok パターンのシンプルな例を次に示します。

```
%{TIMESTAMP_IS08601:timestamp} \[%{MESSAGEPREFIX:message_prefix}\]  
%{CRAWLERLOGLEVEL:loglevel} : %{GREEDYDATA:message}
```

データが TIMESTAMP_IS08601 と一致すると、スキーマの列 timestamp が作成されます。動作は、例にある他の名前付きパターンに似ています。

カスタムパターン

独自に定義するオプションのカスタムパターン。これらのパターンは、データを分類する grok パターンにより参照されます。データに適用される grok パターンでこれらのカスタムパターンを参照できます。各カスタムコンポーネントパターンは別々の行にある必要があります。[正規表現 \(regex\)](#) 構文は、パターンを定義するために使用されます。

以下は、カスタムパターンを使用する例です。

```
CRAWLERLOGLEVEL (BENCHMARK|ERROR|WARN|INFO|TRACE)
MESSAGEPREFIX .*-.*-.*-.*-.*
```

最初の名前付きカスタムパターンである CRAWLERLOGLEVEL は、列挙された文字列の 1 つとデータが一致するとき一致となります。2 番目のカスタムパターン MESSAGEPREFIX は、メッセージのプレフィックス文字列との一致を試みます。

AWS Glue は、作成日時、最終更新時間、分類子のバージョンを追跡します。

AWS Glue 組み込みパターン

AWS Glue は、カスタム分類子を構築するために使用できる多くの一般的なパターンを提供します。分類子の定義の grok pattern に名前付きパターンを追加します。

次のリストは、各パターンの行です。各行で、パターン名の後に定義があります。[正規表現 \(regex\)](#) 構文は、パターンを定義するために使用されます。

```
#<noLOC>&GLU;</noLOC> Built-in patterns
USERNAME [a-zA-Z0-9._-]+
USER %{USERNAME:UNWANTED}
INT (?:[+-]?(?:[0-9]+))
BASE10NUM (?![0-9.+])(?>[+-]?(?:[0-9]+(?:\.[0-9]+)?)|(?:\.[0-9]+))
NUMBER (?:%{BASE10NUM:UNWANTED})
BASE16NUM (?![0-9A-Fa-f])(?:[+-]?(?:0x)?(?:[0-9A-Fa-f]+))
BASE16FLOAT \b(?![0-9A-Fa-f.])?(?:[+-]?(?:0x)?(?:[0-9A-Fa-f]+(?:\.[0-9A-Fa-f]*)?)|
(?:\.[0-9A-Fa-f]+))\b
BOOLEAN (?i)(true|false)

POSINT \b(?:[1-9][0-9]*)\b
NONNEGINT \b(?:[0-9]+)\b
WORD \b\w+\b
NOTSPACE \S+
```

```

SPACE \s*
DATA .*?
GREEDYDATA .*
#QUOTEDSTRING (?: (?<!\\)(?: "(?: \\.| [^\\""])*" | '(?: \\.| [^\\"'])*') | (? `(?: \\.| [^\\"`])* `) )
QUOTEDSTRING (?: (?<!\\)(?: "(?: \\.| [^\\""])+)" | "'(?: \\.| [^\\"'])+'" | '(?: `(?: \\.| [^\\"`])+`)' )
UUID [A-Fa-f0-9]{8}-(?:[A-Fa-f0-9]{4}-){3}[A-Fa-f0-9]{12}

# Networking
MAC (?:%{CISCOMAC:UNWANTED}|%{WINDOWSMAC:UNWANTED}|%{COMMONMAC:UNWANTED})
CISCOMAC (?: (?: [A-Fa-f0-9]{4}\.){2}[A-Fa-f0-9]{4} )
WINDOWSMAC (?: (?: [A-Fa-f0-9]{2}-){5}[A-Fa-f0-9]{2} )
COMMONMAC (?: (?: [A-Fa-f0-9]{2}:){5}[A-Fa-f0-9]{2} )
IPV6 ( ( ( [0-9A-Fa-f]{1,4}: ){7} ( [0-9A-Fa-f]{1,4}|: ) ) | ( ( [0-9A-Fa-f]{1,4}: ){6} ( : [0-9A-Fa-f]{1,4}| ( (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) (\.(25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) ){3} ) | : ) ) | ( ( [0-9A-Fa-f]{1,4}: ){5} ( ( ( : [0-9A-Fa-f]{1,4} ){1,2} ) | : ( (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) (\.(25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) ){3} ) | : ) ) | ( ( [0-9A-Fa-f]{1,4}: ){4} ( ( ( : [0-9A-Fa-f]{1,4} ){1,3} ) | ( ( : [0-9A-Fa-f]{1,4} )? : ( (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) (\.(25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) ){3} ) ) | : ) ) | ( ( [0-9A-Fa-f]{1,4}: ){3} ( ( ( : [0-9A-Fa-f]{1,4} ){1,4} ) | ( ( : [0-9A-Fa-f]{1,4} ){0,2} : ( (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) (\.(25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) ){3} ) ) | : ) ) | ( ( [0-9A-Fa-f]{1,4}: ){2} ( ( ( : [0-9A-Fa-f]{1,4} ){1,5} ) | ( ( : [0-9A-Fa-f]{1,4} ){0,3} : ( (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) (\.(25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) ){3} ) ) | : ) ) | ( ( [0-9A-Fa-f]{1,4}: ){1} ( ( ( : [0-9A-Fa-f]{1,4} ){1,6} ) | ( ( : [0-9A-Fa-f]{1,4} ){0,4} : ( (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) (\.(25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) ){3} ) ) | : ) ) | ( ( ( ( : [0-9A-Fa-f]{1,4} ){1,7} ) | ( ( : [0-9A-Fa-f]{1,4} ){0,5} : ( (25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) (\.(25[0-5]|2[0-4]\d|1\d\d|1-9)?\d) ){3} ) ) | : ) ) ) ( % . + ) ?
IPV4 ( ?<! [0-9] ) ( ? : ( ? : 25 [0-5] | 2 [0-4] [0-9] | [0-1] ? [0-9] { 1, 2 } ) [ . ] ( ? : 25 [0-5] | 2 [0-4] [0-9] | [0-1] ? [0-9] { 1, 2 } ) [ . ] ( ? : 25 [0-5] | 2 [0-4] [0-9] | [0-1] ? [0-9] { 1, 2 } ) [ . ] ( ? : 25 [0-5] | 2 [0-4] [0-9] | [0-1] ? [0-9] { 1, 2 } ) ) ( ? ! [0-9] )
IP ( ? : % { IPV6 : UNWANTED } | % { IPV4 : UNWANTED } )
HOSTNAME \b ( ? : [0-9A-Za-z][0-9A-Za-z-_] { 0, 62 } ) ( ? : \. ( ? : [0-9A-Za-z][0-9A-Za-z-_] { 0, 62 } ) ) * ( \. ? | \b )
HOST % { HOSTNAME : UNWANTED }
IPORHOST ( ? : % { HOSTNAME : UNWANTED } | % { IP : UNWANTED } )
HOSTPORT ( ? : % { IPORHOST } : % { POSINT : PORT } )

# paths
PATH ( ? : % { UNIXPATH } | % { WINPATH } )
UNIXPATH ( ?> / ( ?> [ \w _ % ! $ @ : . , ~ - ] + | \\ . ) * ) +
#UNIXPATH ( ?<! [ \w \ / ] ) ( ? : / [ ^ \ \ s ? * ] * ) +
TTY ( ? : / dev / ( pts | tty ( [ pq ] ) ? ) ( \w + ) ? / ( ? : [ 0-9 ] + ) )
WINPATH ( ?> [ A-Za-z ] + : | \\ ) ( ? : \\ [ ^ \ \ ? * ] * ) +
URIPROTO [ A-Za-z ] + ( \ + [ A-Za-z ] + ) ?

```

```

URIHOST %{IPORHOST}(?:%{POSINT:port})?
# uripath comes loosely from RFC1738, but mostly from what Firefox
# doesn't turn into %XX
URIPATH (?:/[A-Za-z0-9$.+!*'(){}~:;=@#%_\-]*)+
#URIPARAM \?(?:[A-Za-z0-9]+(?:=(?:[^\&]*))?(?:&(?:[A-Za-z0-9]+(?:=(?:[^\&]*)))?)*)?
URIPARAM \?[A-Za-z0-9$.+!*'|(){}~@#%&/=:;_?\-\[\]]*
URIPATHPARAM %{URIPATH}(?:%{URIPARAM})?
URI %{URIPROTO}://(?:%{USER}(?:[^\@]*)?@)?(?:%{URIHOST})?(?:%{URIPATHPARAM})?

# Months: January, Feb, 3, 03, 12, December
MONTH \b(?:Jan(?:uary)?|Feb(?:ruary)?|Mar(?:ch)?|Apr(?:il)?|May|Jun(?:e)?|Jul(?:y)?|
Aug(?:ust)?|Sep(?:tember)?|Oct(?:ober)?|Nov(?:ember)?|Dec(?:ember)?)\b
MONTHNUM (?:\d?[1-9]|1[0-2])
MONTHNUM2 (?:\d[1-9]|1[0-2])
MONTHDAY (?:\d?[1-9])|(?:[12][0-9])|(?:3[01])|[1-9])

# Days: Monday, Tue, Thu, etc...
DAY (?:Mon(?:day)?|Tue(?:sday)?|Wed(?:nesday)?|Thu(?:rsday)?|Fri(?:day)?|
Sat(?:urday)?|Sun(?:day)?)

# Years?
YEAR (?:>\d\d){1,2}
# Time: HH:MM:SS
#TIME \d{2}:\d{2}(?::\d{2}(?:\.\d+)?)?
# TIME %{POSINT<24}:%{POSINT<60}(?::%{POSINT<60}(?:\.%{POSINT})?)?
HOUR (?:2[0123]|[01]?[0-9])
MINUTE (?:[0-5][0-9])
# '60' is a leap second in most time standards and thus is valid.
SECOND (?:\d?[0-5]?[0-9]|60)(?:[:.,][0-9]+)?
TIME (?!<[0-9])%{HOUR}:%{MINUTE}(?::%{SECOND})(?![0-9])
# datestamp is YYYY/MM/DD-HH:MM:SS.UUUU (or something like it)
DATE_US %{MONTHNUM}[/-]%{MONTHDAY}[/-]%{YEAR}
DATE_EU %{MONTHDAY}[./-]%{MONTHNUM}[./-]%{YEAR}
DATESTAMP_US %{DATE_US}[- ]%{TIME}
DATESTAMP_EU %{DATE_EU}[- ]%{TIME}
ISO8601_TIMEZONE (?:Z|[+-%]{HOUR}(?::%{MINUTE}))
ISO8601_SECOND (?:%{SECOND}|60)
TIMESTAMP_ISO8601 %{YEAR}-%{MONTHNUM}-%{MONTHDAY}[T ]%{HOUR}:%{MINUTE}(?::%{
%{SECOND}})?%{ISO8601_TIMEZONE}?
TZ (?:[PMCE][SD]T|UTC)
DATESTAMP_RFC822 %{DAY} %{MONTH} %{MONTHDAY} %{YEAR} %{TIME} %{TZ}
DATESTAMP_RFC2822 %{DAY}, %{MONTHDAY} %{MONTH} %{YEAR} %{TIME} %{ISO8601_TIMEZONE}
DATESTAMP_OTHER %{DAY} %{MONTH} %{MONTHDAY} %{TIME} %{TZ} %{YEAR}
DATESTAMP_EVENTLOG %{YEAR}%{MONTHNUM2}%{MONTHDAY}%{HOUR}%{MINUTE}%{SECOND}

```

```

CISCOTIMESTAMP %{{MONTH}} %{{MONTHDAY}} %{{TIME}}

# Syslog Dates: Month Day HH:MM:SS
SYSLOGTIMESTAMP %{{MONTH}} +%{{MONTHDAY}} %{{TIME}}
PROG (?:[\w._/%-]+)
SYSLOGPROG %{{PROG:program}}(?:\[{{POSINT:pid}}\])?
SYSLOGHOST %{{IPORHOST}}
SYSLOGFACILITY <{{NONNEGINT:facility}}.{{NONNEGINT:priority}}>
HTTPDATE %{{MONTHDAY}}/%{{MONTH}}/%{{YEAR}}:%{{TIME}} %{{INT}}

# Shortcuts
QS %{{QUOTEDSTRING:UNWANTED}}

# Log formats
SYSLOGBASE %{{SYSLOGTIMESTAMP:timestamp}} (?:%{{SYSLOGFACILITY}} )?%{{SYSLOGHOST:logsource}}
%{{SYSLOGPROG}}:

MESSAGESLOG %{{SYSLOGBASE}} %{{DATA}}

COMMONAPACHELOG %{{IPORHOST:clientip}} %{{USER:ident}} %{{USER:auth}}
\[{{HTTPDATE:timestamp}}\] "(?:%{{WORD:verb}} %{{NOTSPACE:request}}(?: HTTP/
%{{NUMBER:httpversion}})?|%{{DATA:rawrequest}})" %{{NUMBER:response}} (?:%{{Bytes:bytes=
%{{NUMBER}}|-}))
COMBINEDAPACHELOG %{{COMMONAPACHELOG}} %{{QS:referrer}} %{{QS:agent}}
COMMONAPACHELOG_DATATYPED %{{IPORHOST:clientip}} %{{USER:ident;boolean}} %{{USER:auth}}
\[{{HTTPDATE:timestamp;date;dd/MMM/yyyy:HH:mm:ss Z}}\] "(?:%{{WORD:verb;string}}
%{{NOTSPACE:request}}(?: HTTP/%{{NUMBER:httpversion;float}})?|%{{DATA:rawrequest}})"
%{{NUMBER:response;int}} (?:%{{NUMBER:bytes;long}}|-)

# Log Levels
LOGLEVEL ([A|a]lert|ALERT|[T|t]race|TRACE|[D|d]ebug|DEBUG|[N|n]otice|NOTICE|[I|i]nfo|
INFO|[W|w]arn(?:ing)?|WARN(?:ING)?|[E|e]rr(?:or)?|ERR(?:OR)?|[C|c]rit(?:ical)?|
CRIT(?:ICAL)?|[F|f]atal|FATAL|[S|s]evere|SEVERE|EMERG(?:ENCY)?|[E|e]merg(?:ency)?)

```

XML カスタム分類子の書き込み

XML は、ファイル内のタグを使用してドキュメントの構造を定義します。XML カスタム分類子で、行の定義に使用されるタグ名を指定できます。

AWS Glue のカスタム分類子の値

XML 分類子を定義する場合、AWS Glue に以下の値を指定し分類子を作成します。この分類子の分類フィールドは `xml` に設定してあります。

名前

分類子の名前。

行タグ

XML ドキュメントでテーブル行を定義する XML タグ名、山括弧 < > なし。名前は XML タグ規則に沿って命名する必要があります。

Note

行データを含む要素は、自動で閉じる空の要素にすることはできません。たとえば、次の空の要素は によって解析されませんAWS Glue。

```
<row att1="xx" att2="yy" />
```

空の要素は次のように記述できます。

```
<row att1="xx" att2="yy"> </row>
```

AWS Glue は、作成日時、最終更新時間、分類子のバージョンを追跡します。

たとえば、次の XML ファイルがあるとします。筆者と役職の列のみを含む AWS Glue テーブルを作成するには、AnyCompany として [行タグ] を使用し、AWS Glue コンソールで分類子を作成します。次に、このカスタム分類子を使用するクローラを追加して実行します。

```
<?xml version="1.0"?>
<catalog>
  <book id="bk101">
    <AnyCompany>
      <author>Rivera, Martha</author>
      <title>AnyCompany Developer Guide</title>
    </AnyCompany>
  </book>
  <book id="bk102">
```

```
<AnyCompany>
  <author>Stiles, John</author>
  <title>Style Guide for AnyCompany</title>
</AnyCompany>
</book>
</catalog>
```

JSON カスタム分類子の書き込み

JSON は、データ交換形式です。名前と値のペア、または順序付きの値のリストでデータ構造を定義します。JSON カスタム分類子では、データ構造への JSON パスを指定し、それを使用してテーブルのスキーマを定義できます。

AWS Glue のカスタム分類子の値

JSON 分類子を定義する場合、AWS Glue に以下の値を指定し分類子を作成します。この分類子の分類フィールドは `json` に設定してあります。

名前

分類子の名前。

JSON パス

テーブルスキーマを定義するために使用されるオブジェクトを指す JSON パス。JSON パスは、ドット表記またはブラケット表記で記述できます。以下の演算子がサポートされています。

説明

JSON オブジェクトのルート要素。すべてのパス式はこれで始まります

ワイルドカード文字。JSON パスで名前または数値が必要な箇所ですべてでも使用可能。

ドット表記の子。JSON オブジェクトの子フィールドを指定します。

ブラケット表記の子。JSON オブジェクトの子フィールドを指定します。1 つの子フィールドのみを指定できます。

配列インデックス。インデックスにより配列の値を指定します。

AWS Glue は、作成日時、最終更新時間、分類子のバージョンを追跡します。

Example JSON 分類子を使用して配列からレコードをプルする

JSON データがレコードの配列だとします。たとえば、ファイルの最初の数行は次のようになります。

```
[
  {
    "type": "constituency",
    "id": "ocd-division\country:us\state:ak",
    "name": "Alaska"
  },
  {
    "type": "constituency",
    "id": "ocd-division\country:us\state:al\cd:1",
    "name": "Alabama's 1st congressional district"
  },
  {
    "type": "constituency",
    "id": "ocd-division\country:us\state:al\cd:2",
    "name": "Alabama's 2nd congressional district"
  },
  {
    "type": "constituency",
    "id": "ocd-division\country:us\state:al\cd:3",
    "name": "Alabama's 3rd congressional district"
  },
  {
    "type": "constituency",
    "id": "ocd-division\country:us\state:al\cd:4",
    "name": "Alabama's 4th congressional district"
  },
  {
    "type": "constituency",
    "id": "ocd-division\country:us\state:al\cd:5",
    "name": "Alabama's 5th congressional district"
  },
  {
    "type": "constituency",
    "id": "ocd-division\country:us\state:al\cd:6",
    "name": "Alabama's 6th congressional district"
  },
]
```

```
{
  "type": "constituency",
  "id": "ocd-division\country:us\state:al\cd:7",
  "name": "Alabama's 7th congressional district"
},
{
  "type": "constituency",
  "id": "ocd-division\country:us\state:ar\cd:1",
  "name": "Arkansas's 1st congressional district"
},
{
  "type": "constituency",
  "id": "ocd-division\country:us\state:ar\cd:2",
  "name": "Arkansas's 2nd congressional district"
},
{
  "type": "constituency",
  "id": "ocd-division\country:us\state:ar\cd:3",
  "name": "Arkansas's 3rd congressional district"
},
{
  "type": "constituency",
  "id": "ocd-division\country:us\state:ar\cd:4",
  "name": "Arkansas's 4th congressional district"
}
]
```

組み込み JSON 分類子を使用してクローラを実行する場合、ファイル全体がスキーマを定義するために使用されます。JSON パスを指定しないので、クローラはデータを1つのオブジェクト、つまり、ただの配列として処理します。たとえば、スキーマは次のようになります。

```
root
|-- record: array
```

ただし、JSON 配列の各レコードに基づいたスキーマを作成するには、カスタム JSON 分類子を作成し、JSON パスを `$[*]` として指定します。この JSON パスを指定すると、分類子は配列内の 12 レコードすべてに問合せスキーマを決定します。結果のスキーマには、各オブジェクトに次の例と同様の個別のフィールドが含まれています。

```
root
```

```
|-- type: string
|-- id: string
|-- name: string
```

Example JSON 分類子を使用してファイルの一部のみを確認する

JSON データが、<http://everypolitician.org/> から取られた JSON ファイルの例 `s3://awsglue-datasets/examples/us-legislators/all/areas.json` のパターンと同様だとします。JSON ファイルのサンプルオブジェクトは、次のようになります。

```
{
  "type": "constituency",
  "id": "ocd-division/country:us/state:ak",
  "name": "Alaska"
}
{
  "type": "constituency",
  "identifiers": [
    {
      "scheme": "dmoz",
      "identifier": "Regional/North_America/United_States/Alaska/"
    },
    {
      "scheme": "freebase",
      "identifier": "\/m\/0hjy"
    },
    {
      "scheme": "fips",
      "identifier": "US02"
    },
    {
      "scheme": "quora",
      "identifier": "Alaska-state"
    },
    {
      "scheme": "britannica",
      "identifier": "place/Alaska"
    },
    {
      "scheme": "wikidata",
      "identifier": "Q797"
    }
  ]
}
```

```
],
  "other_names": [
    {
      "lang": "en",
      "note": "multilingual",
      "name": "Alaska"
    },
    {
      "lang": "fr",
      "note": "multilingual",
      "name": "Alaska"
    },
    {
      "lang": "nov",
      "note": "multilingual",
      "name": "Alaska"
    }
  ],
  "id": "ocd-division/country:us/state:ak",
  "name": "Alaska"
}
```

組み込み JSON 分類子を使用してクローラを実行する場合、ファイル全体がスキーマを作成するために使用されます。最終的に次のようなスキーマになります。

```
root
|-- type: string
|-- id: string
|-- name: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
```

ただし、「id」オブジェクトだけを使用してスキーマを作成するには、カスタム JSON 分類子を作成し、JSON パスを \$.id と指定します。その後、スキーマは「id」フィールドのみに基づくものとなります。

```
root
|-- record: string
```

このスキーマで抽出されたデータの最初の数行は次のようになります。

```
{"record": "ocd-division/country:us/state:ak"}
{"record": "ocd-division/country:us/state:al/cd:1"}
{"record": "ocd-division/country:us/state:al/cd:2"}
{"record": "ocd-division/country:us/state:al/cd:3"}
{"record": "ocd-division/country:us/state:al/cd:4"}
{"record": "ocd-division/country:us/state:al/cd:5"}
{"record": "ocd-division/country:us/state:al/cd:6"}
{"record": "ocd-division/country:us/state:al/cd:7"}
{"record": "ocd-division/country:us/state:ar/cd:1"}
{"record": "ocd-division/country:us/state:ar/cd:2"}
{"record": "ocd-division/country:us/state:ar/cd:3"}
{"record": "ocd-division/country:us/state:ar/cd:4"}
{"record": "ocd-division/country:us/state:as"}
{"record": "ocd-division/country:us/state:az/cd:1"}
{"record": "ocd-division/country:us/state:az/cd:2"}
{"record": "ocd-division/country:us/state:az/cd:3"}
{"record": "ocd-division/country:us/state:az/cd:4"}
{"record": "ocd-division/country:us/state:az/cd:5"}
{"record": "ocd-division/country:us/state:az/cd:6"}
{"record": "ocd-division/country:us/state:az/cd:7"}
```

JSON ファイルの「identifier」のように、深くネストされたオブジェクトに基づいてスキーマを作成するには、カスタム JSON 分類子を作成して JSON パスを \$.identifiers[*].identifier と指定します。スキーマは前の例と似ていますが、JSON ファイル内の別のオブジェクトに基づいています。

スキーマは次のようになります。

```
root
```

```
|-- record: string
```

テーブルからのデータの最初の数行のリストには、スキーマが「identifier」オブジェクトのデータに基づくものであることが示されます。

```
{ "record": "Regional/North_America/United_States/Alaska/" }
{ "record": "/m/0hjy" }
{ "record": "US02" }
{ "record": "5879092" }
{ "record": "4001016-8" }
{ "record": "destination/alaska" }
{ "record": "1116270" }
{ "record": "139487266" }
{ "record": "n79018447" }
{ "record": "01490999-8dec-4129-8254-eef6e80fadc3" }
{ "record": "Alaska-state" }
{ "record": "place/Alaska" }
{ "record": "Q797" }
{ "record": "Regional/North_America/United_States/Alabama/" }
{ "record": "/m/0gyh" }
{ "record": "US01" }
{ "record": "4829764" }
{ "record": "4084839-5" }
{ "record": "161950" }
{ "record": "131885589" }
```

JSON ファイルの「name」配列の「other_names」フィールドのように、別の深くネストされたオブジェクトに基づいてテーブルを作成するには、カスタム JSON 分類子を作成して JSON パスを `$.other_names[*].name` と指定します。スキーマは前の例と似ていますが、JSON ファイル内の別のオブジェクトに基づいています。スキーマは次のようになります。

```
root
|-- record: string
```

テーブルからのデータの最初の数行のリストには、「name」配列内の「other_names」オブジェクトのデータに基づくものであることが示されます。

```
{ "record": "Alaska" }
```

```
{"record": "Alaska"}
{"record": "Аляска"}
{"record": "Alaska"}
{"record": "#####"}
{"record": "#####"}
{"record": "#####"}
{"record": "Alaska"}
{"record": "Alyaska"}
{"record": "Alaska"}
{"record": "Alaska"}
{"record": "Штат Аляска"}
{"record": "Аляска"}
{"record": "Alaska"}
{"record": "#####"}

```

CSV カスタム分類子の書き込み

カスタム CSV 分類子を使用すると、カスタム csv 分類子フィールドの各列のデータ型を指定できます。各列のデータ型をコンマで区切って指定できます。データ型を指定することで、クローラーが推測したデータ型を上書きし、データが適切に分類されるようにすることができます。

分類子で CSV を処理するための SerDe を設定でき、データカタログに適用されます。

カスタム分類子を作成すると、その分類子を別のクローラーで再利用することもできます。

- ヘッダーのみ (データなし) の csv ファイルの場合、十分な情報が提供されないため、これらのファイルは UNKNOWN に分類されます。[Column headings] (列見出し) オプションで CSV の [Has headings](見出しあり) を指定し、データ型を指定すると、これらのファイルを正しく分類できます。

カスタム CSV 分類子を使用して、さまざまな種類の CSV データのスキーマを推測できます。分類子に指定できるカスタム属性には、区切り記号、CSV SerDe オプション、ヘッダーに関するオプション、およびデータに対して特定の検証を実行するかどうかなどがあります。

AWS Glue のカスタム分類子の値

CSV 分類子を定義する場合、AWS Glue に以下の値を入力して分類子を作成します。この分類子の分類フィールドは csv に設定してあります。

分類子名

分類子の名前。

CSV SerDe

分類子で CSV を処理するため、データカタログに適用される SerDe を設定します。オプションは、「Open CSV SerDe」、「Lazy Simple SerDe」、および「None」です。クローラーが検出を行う場合は、None 値を指定できます。

列の区切り文字

行内の各列エントリを区切るものを示すカスタム記号。ユニコード文字を入力します。区切り記号を入力できない場合は、コピーと貼り付けをすることができます。これは、システムでサポートされていない文字 (通常は □ と表示) も含め、印刷可能な文字にも有効です。

引用記号

コンテンツを単一の列の値に結合するものを示すカスタム記号。列の区切り文字とは異なる必要があります。ユニコード文字を入力します。区切り記号を入力できない場合は、コピーと貼り付けをすることができます。これは、システムでサポートされていない文字 (通常は □ と表示) も含め、印刷可能な文字にも有効です。

列見出し

CSV ファイルで列見出しを検出する方法の動作を示します。カスタム CSV ファイルに列見出しがある場合は、列見出しのカンマ区切りリストを入力します。

処理オプション: 単一系列のファイルを許可します

1 つの列のみを含むファイルの処理を有効にします。

処理オプション: 列値を識別する前に空白を削除します

列の値の型を識別する前に値をトリミングするかどうかを指定します。

カスタムデータ型 – オプション

カスタムデータ型をコンマで区切って入力します。CSV ファイル内のカスタムデータ型を指定します。カスタムデータ型は、サポートされているデータ型である必要があります。サ

ポートされているデータ型は、「BINARY」、「BOOLEAN」、「DATE」、「DECIMAL」、「DOUBLE」、「FLOAT」、「INT」、「LONG」、「SHORT」、「STRING」、「TIMESTAMP」です。サポートされていないデータ型の場合、エラーが表示されます。

AWS Glue コンソールでの分類子の操作

分類子は、データのスキーマを決定します。カスタムの分類子を記述し、AWS Glue から指定します。

分類子の表示

作成したすべての分類子のリストを表示するには、<https://console.aws.amazon.com/glue/> で AWS Glue コンソールを開き、[Classifiers] (分類子) タブを選択します。

リストには、各分類子に関する次のプロパティが表示されます。

- 分類子 – 分類子名。分類子を作成するときは、その名前を指定する必要があります。
- 分類 – この分類子によって推測されたテーブルの分類タイプ。
- 最終更新 – 分類子が最後に更新された時刻。

分類子の管理

コンソールの [ClassifiersAWS Glue] (分類子) リストから、分類子の追加、編集、削除ができます。分類子の詳細を表示するには、リスト内の分類子名を選択します。詳細には、分類子を作成したときに定義した情報が含まれます。

分類子の作成

AWS Glue コンソールに分類子を追加するには、[Add classifier (分類子を追加)] を選択します。分類子を定義する場合、以下の値を指定します。

- 分類子名 – 分類子の一意の名前を指定します。
- 分類子タイプ – この分類子によって推測されたテーブルの分類タイプ。
- 最終更新 – 分類子が最後に更新された時刻。

分類子名

分類子の一意の名前を指定します。

分類子タイプ

作成する分類子のタイプを選択します。

選択した分類子のタイプに応じて、分類子の次のプロパティを設定します。

Grok

- 分類

分類されたデータの形式またはタイプを記述、またはカスタムラベルを指定します。

- Grok パターン

これを使用してデータを解析して構造化スキーマにします。grok パターンは、データストアの形式を記述する名前付きパターンで構成されています。この grok パターンは、AWS Glue によって提供された名前付きの組み込みパターンと、[Custom patterns (カスタムパターン)] フィールドに書き込み、含まれるカスタムパターンを使用して書き込みます。grok デバッガーの結果は AWS Glue の結果と正確には一致しませんが、grok デバッガーでサンプルデータを使用してパターンを試すことをお勧めします。ウェブ上で grok デバッガーを見つけることができます。AWS Glue によって提供される名前付き組み込みパターンは、一般にウェブ上で利用可能な grok パターンと互換性があります。

名前付きパターンを反復的に追加して grok パターンを作成し、デバッガーで結果を確認します。このアクティビティを使用すると、AWS Glue クローラが grok パターンを実行するときにデータを解析できるという確信が得られます。

- カスタムパターン

grok 分類子の場合、これらは、記述した [Grok pattern] (Grok パターン) のオプションの構成要素です。組み込みのパターンでデータを解析できない場合は、カスタムパターンを記述する必要があります。これらのカスタムパターンはこのフィールドで定義され、[Grok pattern] (Grok パターン) フィールドで参照されます。各カスタムパターンは個別の行に定義されています。組み込みパターンと同様に、[\[regular expression \(regex\)\]](#) (正規表現) 構文を使用する名前付きパターン定義で構成されています。

たとえば、次の MESSAGEPREFIX という名前は、その後に正規表現の定義が続いてデータに適用され、パターンに従っているかどうか判断されます。

```
MESSAGEPREFIX .*-.*-.*-.*-.*
```

XML

- 行タグ

XML 分類子では、これは XML 文書のテーブル行を定義する XML タグの名前です。山括弧 < > を付けずに名前を入力します。名前は XML タグ規則に沿って命名する必要があります。

詳しくは、「[XML カスタム分類子の書き込み](#)」を参照してください。

JSON

- JSON パス

JSON 分類子の場合、これは、作成するテーブルの行を定義するオブジェクト、配列、または値への JSON パスです。名前をドットで入力するか、AWS Glue でサポートされる演算子を使用して JSON 構文を括弧で囲んでください。

詳細については、「[JSON カスタム分類子の書き込み](#)」の演算子のリストを参照してください。

CSV

- 列の区切り文字

行内の各列エントリの区切りを示す単一の文字または記号。リストから区切り文字または記号を選択するか、Other を選択して、カスタム区切り文字または記号を入力します。

- 引用記号

コンテンツを結合して単一の列の値にすることを示す単一の文字または記号。列の区切り文字とは異なる必要があります。リストから引用記号を選択するか、Other を選択して、カスタム引用文字を入力します。

- 列見出し

CSV ファイルで列見出しを検出する方法の動作を示します。Has headings、No headings または Detect headings を選択できます。カスタム CSV ファイルに列見出しがある場合は、列見出しのカンマ区切りリストを入力します。

- 単一列のファイルを許可

CSV とみなされるためには、データが 2 列以上かつ 2 行以上必要です。このオプションを使用すると、1 つの列のみを含むファイルが処理できます。

- 列の値を識別する前に空白を削除

このオプションによって、列の値のタイプを識別する前に値の空白を削除するかどうかを指定します。

- カスタムデータ型

(オプション) - カスタムデータ型をカンマ区切りリストに入力します。サポートされているデータ型は、「BINARY」、「BOOLEAN」、「DATE」、「DECIMAL」、「DOUBLE」、「FLOAT」、「INT」、「LONG」、「SHORT」、「STRING」、「TIMESTAMP」です。

- CSV Serde

(オプション) - 分類子で CSV を処理するための SerDe を設定でき、データカタログに適用されます。Open CSV SerDe、Lazy Simple SerDe、または None から選択します。クローラーが検出を行う場合は、None 値を指定できます。

詳細については、「[カスタム分類子の書き込み](#)」を参照してください。

AWS Glue クローラーのスケジュール

AWS Glue クローラーは、オンデマンドで、または定期的なスケジュールで実行できます。クローラスケジュールは cron 形式で表すことができます。詳細については、Wikipedia の「[cron](#)」を参照してください。

スケジュールに基づいてクローラーを作成する場合は、クローラーの実行頻度、実行する曜日、実行時間などの特定の制約を指定できます。これらの制約は cron に基づいています。クローラスケジュールを設定するときは、cron の機能と制限を考慮する必要があります。たとえば、毎月 31 日にクローラーを実行することを選択した場合、いくつかの月には 31 日間はないことに注意してください。

各クローラーのクローリングは最大 12 か月間のみ有効です

cron を使用してジョブおよびクローラーをスケジュールする方法の詳細については、「[ジョブとクローラーの時間ベースのスケジュール](#)」を参照してください。

クローラーの結果と詳細の表示

クローラーが正常に実行されると、データカタログにテーブル定義が作成されます。クローラーの実行が完了したら、ナビゲーションペインで [Tables] (テーブル) をクリックして、指定したデータベースにクローラーにより作成されたテーブルを表示します。

クローラー自体に関連する情報は、次のように表示されます。

- AWS Glue コンソールの [Crawlers] (クローラー) ページに、クローラーの次のプロパティが表示されます。

プロパティ	説明
[Name] (名前)	クローラーを作成する場合、一意の名前を付ける必要があります。
[ステータス]	クローラーには、準備完了、開始中、停止中、スケジュールあり、スケジュール停止などの状態があります。実行中のクローラーは、開始中から停止中に向かって処理していきます。クローラーにアタッチされたスケジュールを再開または一時停止できます。
スケジュール	クローラーをオンデマンドで実行するか、または、スケジュールで頻度を選択できます。クローラーのスケジュールの詳細については、「 クローラのスケジュール 」を参照してください。
最後の実行	クローラーを最後に実行した日時。
ログ	クローラーの最後の実行からの使用可能なログにリンクします。
最後の実行以後のテーブルの変更	クローラーの最後の実行により更新された AWS Glue Data Catalog のテーブルの数。

- クローラーの履歴を表示するには、ナビゲーションペインで [Crawlers] (クローラー) を選択して、作成したクローラーを表示します。使用可能なクローラーのリストからクローラーを選択します。クローラーのプロパティとクローラーの履歴は [Crawler runs] (クローラーの実行) タブに表示されます。

[Crawler runs] (クローラーの実行) タブには、クローラーの各実行履歴に対して、[Start time (UTC)] (開始時刻 (UTC))、[End time (UTC)] (終了時刻 (UTC))、[Duration] (所要時間)、[Status] (ステータス)、[DPU hours] (DPU 時間)、[Table changes] (表の変更) などが表示されます。

[クローラーの実行] タブは、クローラー履歴機能の起動日以降に発生したクロールのみを表示し、最大 12 か月分のクロールのみを保持します。古いクロールは返されません。

- 追加情報を表示するには、クローラー詳細ページのタブを選択します。各タブには、クローラーに関連する情報が表示されます。
 - [Schedule] (スケジュール): クローラー用に作成されたすべてのスケジュールがここに表示されます。
 - [Data sources] (データソース): クローラーによってスキャンされたすべてのデータソースがここに表示されます。
 - [Classifiers] (分類子): クローラーに割り当てられたすべての分類子がここに表示されます。
 - [Tags] (タグ): タグが作成され、AWS リソースに割り当てられていれば、ここに表示されます。

クローラーによって設定されたデータカタログテーブルのパラメータ

これらのテーブルプロパティは、AWS Glue クローラーによって設定されます。ユーザーが `classification` および `compressionType` プロパティを使用することを想定しています。テーブルサイズの推定値など、その他のプロパティは内部計算に使用されるものであり、その精度やお客様のユースケースへの適用性は保証されません。これらのパラメータを変更すると、クローラーの動作が変わる可能性があるため、このワークフローはサポートされません。

プロパティキー	プロパティ値
UPDATED_BY_CRAWLER	更新を実行するクローラーの名前。
connectionName	データストアに接続するために使用されるクローラーのデータカタログ内にある接続名。
recordCount	ファイルサイズとヘッダーに基づいた、テーブル内の推定レコード数。
skip.header.line.count	スキップヘッダーにスキップされた行数。CSV に分類されたテーブルに設定します。

プロパティキー	プロパティ値
CrawlerSchemaSerializerVersion	内部使用目的
classification	クローラーによって推測されたデータの形式。AWS Glue クローラーでサポートされるデータ形式の詳細については、「 the section called “AWS Glue の組み込み分類子” 」を参照してください。
CrawlerSchemaDeserializerVersion	内部使用目的
sizeKey	クロールされたテーブル内のファイルの合計サイズ。
averageRecordSize	テーブル内の行の平均サイズ (バイト単位)。
compressionType	テーブル内のデータで使用される圧縮タイプ。AWS Glue クローラーでサポートされる圧縮タイプの詳細については、「 the section called “AWS Glue の組み込み分類子” 」を参照してください。
typeOfData	file、table、または view。
objectCount	テーブルの Amazon S3 パスの下にあるオブジェクトの数。

これらの追加のテーブルプロパティは、Snowflake データストアの AWS Glue クローラーによって設定されます。

プロパティキー	プロパティ値
aws:RawTableLastAltered	Snowflake テーブルの最後に変更されたタイムスタンプを記録します。
ViewOriginalText	SQL ステートメントを表示します。

プロパティキー	プロパティ値
ViewExpandedText	Base64 形式でエンコードされた SQL ステートメントを表示します。
ExternalTable:S3Location	Snowflake 外部テーブルの Amazon S3 の場所。
ExternalTable:FileFormat	Snowflake 外部テーブルの Amazon S3 ファイル形式。

これらの追加のテーブルプロパティは、Amazon Redshift、Microsoft SQL Server、MySQL、PostgreSQL、Oracle などの JDBC タイプのデータストアの AWS Glue クローラーによって設定されます。

プロパティキー	プロパティ値
aws:RawType	クローラーがデータカタログにデータを保存すると、データ型が Hive 互換型に変換されるため、ネイティブデータ型の情報が失われることがよくあります。クローラーは、ネイティブレベルのデータ型を提供する <code>aws:RawType</code> パラメータを出力します。
aws:RawColumnComment	コメントがデータベース内の列に関連付けられている場合、クローラーはカタログテーブル内の対応するコメントを出力します。コメント文字列は 255 バイトに切り捨てられます。 Microsoft SQL Server ではコメントはサポートされていません。
aws:RawTableComment	コメントがデータベース内のテーブルに関連付けられている場合、クローラーはカタログテーブル内の対応するコメントを出力します。コメント文字列は 255 バイトに切り捨てられます。 Microsoft SQL Server ではコメントはサポートされていません。

クローラーの動作のカスタマイズ

クローラーを実行すると、データストアの変更が検出される場合があります。これらの変更に伴って、以前のクローラーとは異なるスキーマやパーティションが生じることがあります。AWS Management Console または AWS Glue API を使用して、特定のタイプの変更をクローラーで処理する方法を設定できます。

トピック

- [新しいパーティションを追加するための増分クロール](#)
- [パーティションインデックスクローラー設定オプションの設定](#)
- [Amazon S3 イベント通知を使用した加速クロール](#)
- [既存のスキーマをクローラーで変更しないための方法](#)
- [各 Amazon S3 インクルードパスの単一のスキーマを作成する方法](#)
- [テーブルの場所とパーティションレベルの指定方法](#)
- [クローラーが作成できるテーブルの最大数を指定する方法](#)
- [Delta Lake データストアの設定オプションを指定する方法](#)
- [Lake Formation の認証情報を使用するようにクローラーを設定する方法](#)

Console

AWS Glue コンソールを使用してクローラーを定義する場合、クローラーの動作を設定するためのオプションをいくつか使用できます。AWS Glue コンソールを使用してクローラーを追加する方法の詳細については、「[クローラーの設定](#)」を参照してください。

以前にクロールしたデータストアに対してクローラーを実行すると、データストアでのスキーマの変更やオブジェクトの削除が検出される場合があります。クローラーは、スキーマの変更をログに記録します。クローラーのソースタイプに応じて、スキーマの変更ポリシーにかかわらず、新しいテーブルとパーティションが作成される可能性があります。

クローラーがスキーマの変更を検出したときの動作を指定するには、コンソールで以下のいずれかのアクションを選択できます。

- Data Catalog でテーブル定義を更新する – AWS Glue Data Catalog で、新しい列を追加し、欠落している列を削除して、既存の列の定義を変更します。クローラーで設定されていないすべてのメタデータを削除します。これはデフォルトの設定です。

- 新しい列の追加のみ – Amazon S3 データストアにマッピングされるテーブルの場合、検出した新しい列は追加されますが、既存の列のタイプは Data Catalog で削除または変更されません。データカタログの現在の列が正しく、クローラーで既存の列のタイプを削除または変更しない場合は、このオプションを選択します。Amazon S3 の基本的なテーブル属性 (分類、圧縮タイプ、CSV 区切り記号など) が変わった場合は、テーブルを廃止としてマークします。Data Catalog に存在する入力形式と出力形式そのままにします。SerDe パラメータは、クローラーで設定されたものである場合に限り、更新します。他のすべてのデータストアについては、既存の列定義を変更します。
- 変更を無視し、Data Catalog のテーブルを更新しない – 新しいテーブルとパーティションのみが作成されます。

これは、増分クローラーのデフォルト設定です。

新規のパーティションや変更されたパーティションがクローラーで検出される場合もあります。デフォルトでは、変更が行われると、新しいパーティションが追加され、既存のパーティションは更新されます。さらに、コンソールですべての新規および既存のパーティションを更新してテーブルのメタデータを反映するAWS Glueのようにクローラーの設定オプションを設定できます。このオプションを設定すると、パーティションは、分類、入力形式、出力形式、SerDe 情報、スキーマなどのメタデータプロパティを親テーブルから継承します。テーブルでの上記プロパティに対する変更は、そのパーティションに伝播されます。この設定オプションを既存のクローラーに設定すると、既存のパーティションは、次回クローラーが実行されるときに親テーブルのプロパティと一致するよう更新されます。

データストアで削除されたオブジェクトを検出したときのクローラーの動作を指定するには、以下のいずれかのアクションを選択します。

- Data Catalog からテーブルとパーティションを削除する
- 変更を無視し、Data Catalog のテーブルを更新しない

これは、増分クローラーのデフォルト設定です。

- Data Catalog でテーブルを廃止としてマークする – これがデフォルトの設定です。

AWS CLI

```
aws glue create-crawler \  
--name "your-crawler-name" \  
--role "your-iam-role-arn" \  

```

```
--database-name "your-database-name" \  
--targets 'S3Targets=[{Path="s3://your-bucket-name/path-to-data"}]' \  
--configuration '{"Version": 1.0, "CrawlerOutput": {"Partitions":  
{"AddOrUpdateBehavior": "InheritFromTable"}, "Tables": {"AddOrUpdateBehavior":  
"MergeNewColumns"}}}'
```

API

AWS Glue API を使用してクローラーを定義する場合は、いくつかのフィールドから選択してクローラーを設定できます。クローラー API の SchemaChangePolicy は、変更されたスキーマや削除されたオブジェクトを検出したときのクローラーの動作を決定します。クローラーは、実行時にスキーマの変更をログに記録します。

クローラー設定オプションを示すサンプル Python コード

```
import boto3  
import json  
  
# Initialize a boto3 client for AWS Glue  
glue_client = boto3.client('glue', region_name='us-east-1') # Replace 'us-east-1'  
with your desired AWS region  
  
# Define the crawler configuration  
crawler_configuration = {  
    "Version": 1.0,  
    "CrawlerOutput": {  
        "Partitions": {  
            "AddOrUpdateBehavior": "InheritFromTable"  
        },  
        "Tables": {  
            "AddOrUpdateBehavior": "MergeNewColumns"  
        }  
    }  
}  
  
configuration_json = json.dumps(crawler_configuration)  
# Create the crawler with the specified configuration  
response = glue_client.create_crawler(  
    Name='your-crawler-name', # Replace with your desired crawler name  
    Role='crawler-test-role', # Replace with the ARN of your IAM role for Glue  
    DatabaseName='default', # Replace with your target Glue database name  
    Targets={
```

```

    'S3Targets': [
      {
        'Path': "s3://your-bucket-name/path/", # Replace with your S3 path
to the data
      },
    ],
    # Include other target types like 'JdbcTargets' if needed
  },
  Configuration=configuration_json,
  # Include other parameters like Schedule, Classifiers, TablePrefix,
  SchemaChangePolicy, etc., as needed
)

print(response)a

```

クローラーを実行すると、スキーマの変更ポリシーにかかわらず、常に新しいテーブルとパーティションが作成されます。変更されたテーブルスキーマを検出したときのクローラーの動作を決定するには、SchemaChangePolicy 構造の UpdateBehavior フィールドで以下のいずれかのアクションを選択できます。

- UPDATE_IN_DATABASE – のテーブルを更新します。AWS Glue Data Catalog 新しい列を追加し、欠落している列を削除して、既存の列の定義を変更します。クローラーで設定されていないすべてのメタデータを削除します。
- LOG – 変更を無視し、Data Catalog でテーブルを更新しない。

これは、増分クローラのデフォルト設定です。

クローラー API の Configuration フィールドに指定されている JSON オブジェクトを使用して SchemaChangePolicy 構造を上書きすることもできます。この JSON オブジェクトに含まれているキーと値のペアを使用して、既存の列を更新しないで新規の列のみを追加するようにポリシーを設定できます。たとえば、次の JSON オブジェクトを文字列として指定します。

```

{
  "Version": 1.0,
  "CrawlerOutput": {
    "Tables": { "AddOrUpdateBehavior": "MergeNewColumns" }
  }
}

```

このオプションは、コンソールの [Add new columns onlyAWS Glue] (新しい列のみを追加) オプションに対応します。これにより、Amazon S3 データストアをクローलした結果のテーブルの SchemaChangePolicy 構造のみが上書きされます。Data Catalog (信頼できる情報源) にあるがままにメタデータを維持する場合は、このオプションを選択します。新しい列が検出されると追加されます。これには、ネストされたデータ型も含まれます。ただし、既存の列は削除されず、そのタイプは変更されません。Amazon S3 のテーブル属性が大幅に変わる場合は、テーブルを廃止としてマークし、互換性のない属性を解決する必要があるという警告をログに記録します。このオプションは増分クローラーには適用できません。

クローラーが、以前にクローलしたデータストアに対して実行される場合、新規または変更されたパーティションが検出される場合があります。デフォルトでは、変更が行われると、新しいパーティションが追加され、既存のパーティションは更新されます。さらに、クローラーの設定オプションを InheritFromTable に設定できます。このオプションは、コンソールの [Update all new and existing partitions with metadata from the tableAWS Glue] (すべての新規および既存のパーティションを更新してテーブルのメタデータを反映する) オプションに対応します。このオプションを設定すると、親テーブルのメタデータプロパティ (分類、入力形式、出力形式、SerDe 情報、スキーマなど) がパーティションに継承されます。親テーブルでのすべてのプロパティの変更は、そのパーティションに伝播されます。

この設定オプションを既存のクローラーに設定すると、既存のパーティションは、次回クローラーが実行されるときに親テーブルのプロパティと一致するよう更新されます。この動作は、クローラー API の Configuration フィールドで設定します。たとえば、次の JSON オブジェクトを文字列として指定します。

```
{
  "Version": 1.0,
  "CrawlerOutput": {
    "Partitions": { "AddOrUpdateBehavior": "InheritFromTable" }
  }
}
```

クローラー API の Configuration フィールドでは、複数の設定オプションを設定できます。たとえば、パーティションとテーブルの両方のクローラー出力を設定するには、次の JSON オブジェクトの文字列表現を指定できます。

```
{
  "Version": 1.0,
  "CrawlerOutput": {
    "Partitions": { "AddOrUpdateBehavior": "InheritFromTable" },

```

```
    "Tables": {"AddOrUpdateBehavior": "MergeNewColumns" }  
  }  
}
```

データストアで削除されたオブジェクトを検出したときのクローラーの動作を決定するには、以下のいずれかのアクションを選択できます。クローラー API の `SchemaChangePolicy` 構造の `DeleteBehavior` フィールドでは、削除されたオブジェクトを検出したときのクローラーの動作を設定します。

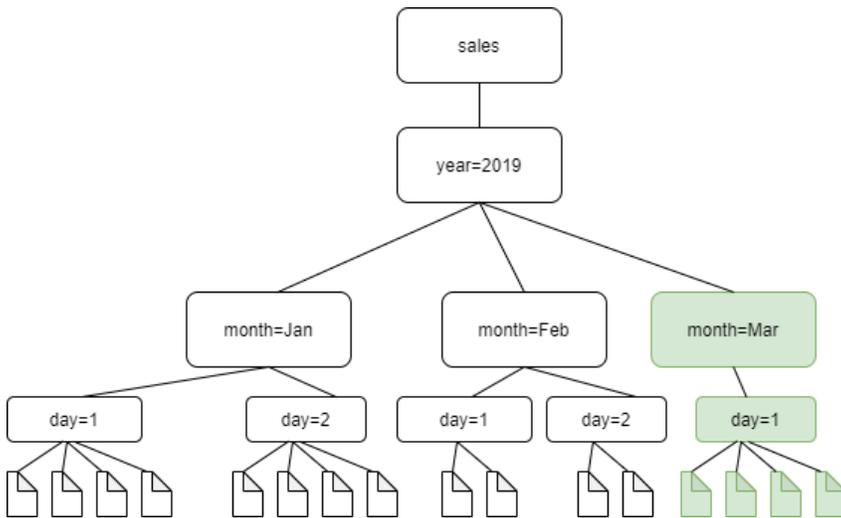
- `DELETE_FROM_DATABASE` – Data Catalog からテーブルとパーティションを削除します。
- `LOG` – 変更を無視します。Data Catalog を更新しません。代わりにログメッセージを書き込みます。
- `DEPRECATE_IN_DATABASE` – Data Catalog でテーブルを廃止としてマークします。これはデフォルトの設定です。

新しいパーティションを追加するための増分クロール

クローラーには新しいパーティションを追加するオプションがあり、安定したテーブルスキーマを持つ増分データセットのクロールが速くなります。典型的なユースケースは、スケジュールされたクローラーで、クロールごとに新しいパーティションが追加されます。このオプションをオンにすると、最初にターゲットデータセットで完全なクロールを実行し、クローラーが初期スキーマとパーティション構造を記録できるようになります。クロールの再実行時、スキーマに互換性がある場合にのみ、新しいパーティションが既存のテーブルに追加されます。最初のクロール実行後に、スキーマの変更は行われず、Data Catalog に新しいテーブルを追加することはありません。

Amazon S3 データソースを設定するときこのオプションを使用できます。CreateCrawler API で「Crawl_New_Folders」として RecrawlBehavior で RecrawlPolicy 設定することも、コンソールで [Subsequent crawler runs] を [Crawl new sub-folders only] と設定することもできます。

[the section called “クローラーは、どのようにパーティションを作成するタイミングを判断していますか?”](#) の例を続けます。次の図は、3月のファイルが追加されていることを示しています。



「Crawl_New_Folders」として RecrawlBehavior オプションを設定した場合、新しいフォルダ month=Mar のみがクローラされます。

注意と制限

このオプションをオンにすると、クローラーの編集時に Amazon S3 ターゲットデータストアを変更できなくなります。このオプションは、ある特定のクローラー設定に影響します。オンにすると、クローラーの更新動作と削除動作が LOG になります。これにより、以下のように処理されます。

- スキーマに互換性がないオブジェクトを検出した場合、クローラーはデータカタログにオブジェクトを追加せず、この詳細を CloudWatch Logs のログとして追加します。
- データカタログで削除されたオブジェクトは更新されません。

詳細については、「[the section called “クローラーの動作のカスタマイズ”](#)」を参照してください。

パーティションインデックスクローラー設定オプションの設定

データカタログは、特定のパーティションを効率的に検索できるようにパーティションインデックスをサポートしています。詳細については、「[AWS Glue でのパーティションインデックスの使用](#)」を参照してください。AWS Glue クローラーは、デフォルトで Amazon S3 および Delta Lake ターゲットのパーティションインデックスを作成します。

クローラーを定義すると、[出力とスケジュールの設定] ページの [詳細オプション] で、[パーティションインデックスを自動的に作成] するオプションがデフォルトで有効になります。

このオプションを無効にするには、コンソールの [パーティションインデックスを自動的に作成] チェックボックスの選択を解除できます。クローラー API を使用し、Configuration で

`CreatePartitionIndex` を設定してこのオプションを無効にすることもできます。デフォルト値は `True` です。

パーティションインデックスの使用に関する注意事項

- クローラーによって作成されたテーブルには、デフォルトでは `partition_filtering.enabled` 変数がありません。詳細については、「[AWS Glue パーティションインデックスとフィルタリング](#)」を参照してください。
- 暗号化されたパーティションのパーティションインデックスの作成はサポートされていません。

Amazon S3 イベント通知を使用した加速クローラ

Amazon S3 または Data Catalog ターゲットからオブジェクトを一覧表示する代わりに、Amazon S3 イベントを使用して変更を検索するようにクローラーを設定できます。この機能は、Amazon S3 または Data Catalog ターゲット全体を一覧表示するのではなく、Amazon S3 イベントを使用してイベントをトリガーしたサブフォルダからのすべてのファイルを一覧表示して 2 つのクローラ間の変更を識別することによって、再クローラ時間を短縮します。

最初のクローラでは、ターゲットからのすべての Amazon S3 オブジェクトを一覧表示します。最初のクローラの成功後、手動または設定されたスケジュールでリクローラを選択できます。クローラーは、すべてのオブジェクトをリストするのではなく、それらのイベントのオブジェクトのみをリストします。

Amazon S3 イベントベースのクローラーに移行する利点は以下のとおりです。

- ターゲットからのすべてのオブジェクトの一覧表示を要しない場合は、より速く再クローラできます。その代わりに、オブジェクトが追加または削除される特定のフォルダが一覧表示されます。
- オブジェクトが追加または削除される特定のフォルダの一覧表示を行うと、全体的なクローラコストが削減されます。

Amazon S3 イベントクローラは、クローラーのスケジュールに基づいて SQS キューから Amazon S3 イベントを使うことで実行します。キューにイベントがない場合、費用はかかりません。Amazon S3 イベントは、SQS キューに直接送信するように設定できます。また、複数のコンシューマーが同じイベント、SNS と SQS の組み合わせを必要とする場合にも設定できます。詳細については、「[the section called “Amazon S3 イベント通知のアカウントを設定します。”](#)」を参照してください。

イベントモードでクローラーを作成して設定した後の最初のクローールは、Amazon S3 または Data Catalog ターゲットの完全な一覧表示を行う一覧表示モードで実行されます。次のログは、最初に成功したクローール、「クローールは Amazon S3 イベントを使用して実行します」の後に、Amazon S3 イベントを使用してクローールのオペレーションを確認します。

Amazon S3 イベントクローールを作成し、クローールに影響を与える可能性のあるクローラーのプロパティを更新すると、クローールがリストモードで動作し、「クローールは S3 イベントモードで実行されていません」というログが追加されます。

Note

クローールごとに消費するメッセージの最大数は 10,000 メッセージです。

カタログターゲット

ターゲットが Data Catalog の場合、クローラーは変更内容 (テーブル内の追加パーティションなど) で Data Catalog 内の既存のテーブルを更新します。

トピック

- [Amazon S3 イベント通知のアカウントを設定します。](#)
- [Amazon S3 イベントクローラで暗号化を使用します。](#)

Amazon S3 イベント通知のアカウントを設定します。

このセクションでは、Amazon S3 イベント通知用にアカウントを設定する方法について説明します。また、スクリプトまたは AWS Glue コンソールを使用して行う手順について説明します。

前提条件

以下の設定タスクを実行します。括弧内の値は、スクリプトの構成可能な設定を参照している点に注意してください。

1. Amazon S3 バケットを作成します (s3_bucket_name)。
2. 識別バケットのパスであるクローラターゲット (folder_name 「test1」などの) を識別します。
3. クローラー名を準備します (crawler_name)
4. クローラー名と同じである可能性がある SNS トピック名 (sns_topic_name) を準備します。

5. クローラーが実行し、S3 バケットが存在する AWS リージョンを準備します。(region)
6. Amazon S3 イベントを取得するためにメールアドレスを使用する場合は、オプションでメールアドレスを準備します(subscribing_email)。

CloudFormation スタックを使用してリソースを作成することもできます。以下のステップを実行します。

1. 米国東部 (バージニア北部) で CloudFormation スタックを[起動](#)します。
2. [パラメータ] に、Amazon S3 バケットの名前 (アカウント番号を含む) を入力します。
3. I acknowledge that AWS CloudFormation might create IAM resources with custom names を選択します。
4. [Create stack] を選択します。

制限:

- ターゲットが Amazon S3 または Data Catalog であるかにかかわらず、クローラーがサポートするのは単一のターゲットのみです。
- プライベート VPC の SQS はサポートされていません。
- Amazon S3 サンプリングはサポートされていません。
- クローラーターゲットは、Amazon S3 ターゲットの場合はフォルダ、Data Catalog ターゲットの場合は 1 つ、または複数の AWS Glue Data Catalog テーブルにする必要があります。
- 「すべての」パスのワイルドカードをサポートしていません: s3://%
- Data Catalog ターゲットの場合、Amazon S3 イベントモードでは、すべてのカタログテーブルが同じ Amazon S3 バケットをポイントする必要があります。
- Data Catalog ターゲットの場合、カタログテーブルは Delta Lake 形式の Amazon S3 ロケーションをポイントしない必要があります (_symlink フォルダが含まれる、またはカタログテーブルの InputFormat を確認)。

Amazon S3 イベントベースのクローラーを使用するには、S3 ターゲットと同じプレフィックスからフィルタリングされたイベント付きの S3 バケットでイベント通知を有効にし、SQS に保存する必要があります。SQS とイベント通知は、[\[チュートリアル:通知のバケットを設定する\]](#) または、[the section called “SQS を生成し、ターゲットから Amazon S3 イベントを設定するスクリプト。”](#) の手順に従って、コンソールから設定できます。

SQS ポリシー

クローラーが使用するロールに添付する必要がある次の SQS ポリシーを追加します。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "VisualEditor0",
      "Effect": "Allow",
      "Action": [
        "sqs:DeleteMessage",
        "sqs:GetQueueUrl",
        "sqs:ListDeadLetterSourceQueues",
        "sqs:ReceiveMessage",
        "sqs:GetQueueAttributes",
        "sqs:ListQueueTags",
        "sqs:SetQueueAttributes",
        "sqs:PurgeQueue"
      ],
      "Resource": "arn:aws:sqs:{region}:{accountID}:cfn-sqs-queue"
    }
  ]
}
```

SQS を生成し、ターゲットから Amazon S3 イベントを設定するスクリプト。

前提条件が満たされていることを確認したら、以下の Python スクリプトを実行して SQS を作成します。構成可能な設定を、前提条件から用意された名前に置き換えます。

Note

スクリプトを実行した後、SQS コンソールにログインして、作成した SQS の ARN を見つけます。

Amazon SQS は、可視性タイムアウト、つまり Amazon SQS が他のカスタマーがそのメッセージの受信や処理できなくなる期間を設定します。可視性タイムアウトをクローラーのランタイム時間とほぼ等しく設定します。

```
#!/venv/bin/python
import boto3
```

```
import boto3

#-----Start : READ ME FIRST -----#
# 1. Purpose of this script is to create the SQS, SNS and enable S3 bucket
  notification.
#   The following are the operations performed by the scripts:
#   a. Enable S3 bucket notification to trigger 's3:ObjectCreated:' and
  's3:ObjectRemoved:' events.
#   b. Create SNS topic for fan out.
#   c. Create SQS queue for saving events which will be consumed by the crawler.
#   SQS Event Queue ARN will be used to create the crawler after running the
  script.
# 2. This script does not create the crawler.
# 3. SNS topic is created to support FAN out of S3 events. If S3 event is also used by
  another
#   purpose, SNS topic created by the script can be used.
# 1. Creation of bucket is an optional step.
#   To create a bucket set create_bucket variable to true.
# 2. The purpose of crawler_name is to easily locate the SQS/SNS.
#   crawler_name is used to create SQS and SNS with the same name as crawler.
# 3. 'folder_name' is the target of crawl inside the specified bucket 's3_bucket_name'
#
#-----End : READ ME FIRST -----#

#-----#
# Start : Configurable settings #
#-----#

#Create
region = 'us-west-2'
s3_bucket_name = 's3eventtestuswest2'
folder_name = "test"
crawler_name = "test33S3Event"
sns_topic_name = crawler_name
sqs_queue_name = sns_topic_name
create_bucket = False

#-----#
# End : Configurable settings #
#-----#

# Define aws clients
dev = boto3.session.Session(profile_name='myprofile')
```

```
boto3.setup_default_session(profile_name='myprofile')
s3 = boto3.resource('s3', region_name=region)
sns = boto3.client('sns', region_name=region)
sqs = boto3.client('sqs', region_name=region)
client = boto3.client("sts")
account_id = client.get_caller_identity()["Account"]
queue_arn = ""

def print_error(e):
    print(e.message + ' RequestId: ' + e.response['ResponseMetadata']['RequestId'])

def create_s3_bucket(bucket_name, client):
    bucket = client.Bucket(bucket_name)
    try:
        if not create_bucket:
            return True
        response = bucket.create(
            ACL='private',
            CreateBucketConfiguration={
                'LocationConstraint': region
            },
        )
        return True
    except botocore.exceptions.ClientError as e:
        print_error(e)
        if 'BucketAlreadyOwnedByYou' in e.message: # we own this bucket so continue
            print('We own the bucket already. Lets continue...')
            return True
    return False

def create_s3_bucket_folder(bucket_name, client, directory_name):
    s3.put_object(Bucket=bucket_name, Key=(directory_name + '/'))

def set_s3_notification_sns(bucket_name, client, topic_arn):
    bucket_notification = client.BucketNotification(bucket_name)
    try:
        response = bucket_notification.put(
            NotificationConfiguration={
                'TopicConfigurations': [
                    {
                        'Id' : crawler_name,
                        'TopicArn': topic_arn,
```

```
        'Events': [
            's3:ObjectCreated:*',
            's3:ObjectRemoved:*',
        ],
        'Filter' : {'Key': {'FilterRules': [{'Name': 'prefix',
'Value': folder_name}]}}
    },
]
)
return True
except boto3.exceptions.ClientError as e:
    print_error(e)
return False

def create_sns_topic(topic_name, client):
    try:
        response = client.create_topic(
            Name=topic_name
        )
        return response['TopicArn']
    except boto3.exceptions.ClientError as e:
        print_error(e)
    return None

def set_sns_topic_policy(topic_arn, client, bucket_name):
    try:
        response = client.set_topic_attributes(
            TopicArn=topic_arn,
            AttributeName='Policy',
            AttributeValue='''{
                "Version": "2008-10-17",
                "Id": "s3-publish-to-sns",
                "Statement": [{
                    "Effect": "Allow",
                    "Principal": { "AWS" : "*" },
                    "Action": [ "SNS:Publish" ],
                    "Resource": "%s",
                    "Condition": {
                        "StringEquals": {
                            "AWS:SourceAccount": "%s"
                        }
                    },
                }
            ]
        }''')
    
```

```
        "ArnLike": {
            "aws:SourceArn": "arn:aws:s3:*:*:*s"
        }
    }
}]]
}''' % (topic_arn, account_id, bucket_name)
)
return True
except botocore.exceptions.ClientError as e:
    print_error(e)

return False

def subscribe_to_sns_topic(topic_arn, client, protocol, endpoint):
    try:
        response = client.subscribe(
            TopicArn=topic_arn,
            Protocol=protocol,
            Endpoint=endpoint
        )
        return response['SubscriptionArn']
    except botocore.exceptions.ClientError as e:
        print_error(e)
    return None

def create_sqs_queue(queue_name, client):
    try:
        response = client.create_queue(
            QueueName=queue_name,
        )
        return response['QueueUrl']
    except botocore.exceptions.ClientError as e:
        print_error(e)
    return None

def get_sqs_queue_arn(queue_url, client):
    try:
        response = client.get_queue_attributes(
            QueueUrl=queue_url,
            AttributeNames=[
                'QueueArn',
```

```

    ]
    )
    return response['Attributes']['QueueArn']
except botocore.exceptions.ClientError as e:
    print_error(e)
return None

def set_sqs_policy(queue_url, queue_arn, client, topic_arn):
    try:
        response = client.set_queue_attributes(
            QueueUrl=queue_url,
            Attributes={
                'Policy': '''{
                    "Version": "2012-10-17",
                    "Id": "AllowSNSPublish",
                    "Statement": [
                        {
                            "Sid": "AllowSNSPublish01",
                            "Effect": "Allow",
                            "Principal": "*",
                            "Action": "SQS:SendMessage",
                            "Resource": "%s",
                            "Condition": {
                                "ArnEquals": {
                                    "aws:SourceArn": "%s"
                                }
                            }
                        }
                    ]
                }''' % (queue_arn, topic_arn)
            )
        return True
    except botocore.exceptions.ClientError as e:
        print_error(e)
    return False

if __name__ == "__main__":
    print('Creating S3 bucket %s.' % s3_bucket_name)
    if create_s3_bucket(s3_bucket_name, s3):
        print('\nCreating SNS topic %s.' % sns_topic_name)
        topic_arn = create_sns_topic(sns_topic_name, sns)
        if topic_arn:

```

```

print('SNS topic created successfully: %s' % topic_arn)

print('Creating SQS queue %s' % sqs_queue_name)
queue_url = create_sqs_queue(sqs_queue_name, sqs)
if queue_url is not None:
    print('Subscribing sqs queue with sns.')
    queue_arn = get_sqs_queue_arn(queue_url, sqs)
    if queue_arn is not None:
        if set_sqs_policy(queue_url, queue_arn, sqs, topic_arn):
            print('Successfully configured queue policy.')
            subscription_arn = subscribe_to_sns_topic(topic_arn, sns,
'sqs', queue_arn)
            if subscription_arn is not None:
                if 'pending confirmation' in subscription_arn:
                    print('Please confirm SNS subscription by visiting the
subscribe URL.')
                else:
                    print('Successfully subscribed SQS queue: ' +
queue_arn)
            else:
                print('Failed to subscribe SNS')
        else:
            print('Failed to set queue policy.')
    else:
        print("Failed to get queue arn for %s" % queue_url)
# ----- End subscriptions to SNS topic -----

print('\nSetting topic policy to allow s3 bucket %s to publish.' %
s3_bucket_name)
if set_sns_topic_policy(topic_arn, sns, s3_bucket_name):
    print('SNS topic policy added successfully.')
    if set_s3_notification_sns(s3_bucket_name, s3, topic_arn):
        print('Successfully configured event for S3 bucket %s' %
s3_bucket_name)
        print('Create S3 Event Crawler using SQS ARN %s' % queue_arn)
    else:
        print('Failed to configure S3 bucket notification.')
else:
    print('Failed to add SNS topic policy.')
else:
    print('Failed to create SNS topic.')

```

コンソールを使用した Amazon S3 イベント通知用のクローラーの設定 (Amazon S3 ターゲット)

Amazon S3 ターゲットのために AWS Glue コンソールを使用して Amazon S3 イベント通知用のクローラーを設定するには、以下の手順を実行します。

1. クローラーのプロパティを設定します。詳細については、「[AWS Glue コンソールでのクローラー設定オプションの設定](#)」を参照してください。
2. [データソースの設定] セクションに、[データは AWS Glue テーブルにマッピング済みですか?] という質問が表示されています。

デフォルトでは、[Not yet] (まだです) が選択されています。Amazon S3 のデータソースを使用しており、データがまだ AWS Glue テーブルにマップされていないため、これはデフォルトままにしておきます。

3. [Data sources] (データソース) セクションで、[Add a data source] (データソースを追加) を選択します。

Step 1
Set crawler properties

Step 2
Choose data sources and classifiers

Step 3
Configure security settings

Step 4
Set output and scheduling

Step 5
Review and create

Choose data sources and classifiers

Data source configuration

Is your data already mapped to Glue tables?

Not yet
Select one or more data sources to be crawled.

Yes
Select existing tables from your Glue Data Catalog.

Data sources (0) Edit Remove Add a data source
The list of data sources to be scanned by the crawler.

Type	Data source	Parameters
You don't have any data sources.		

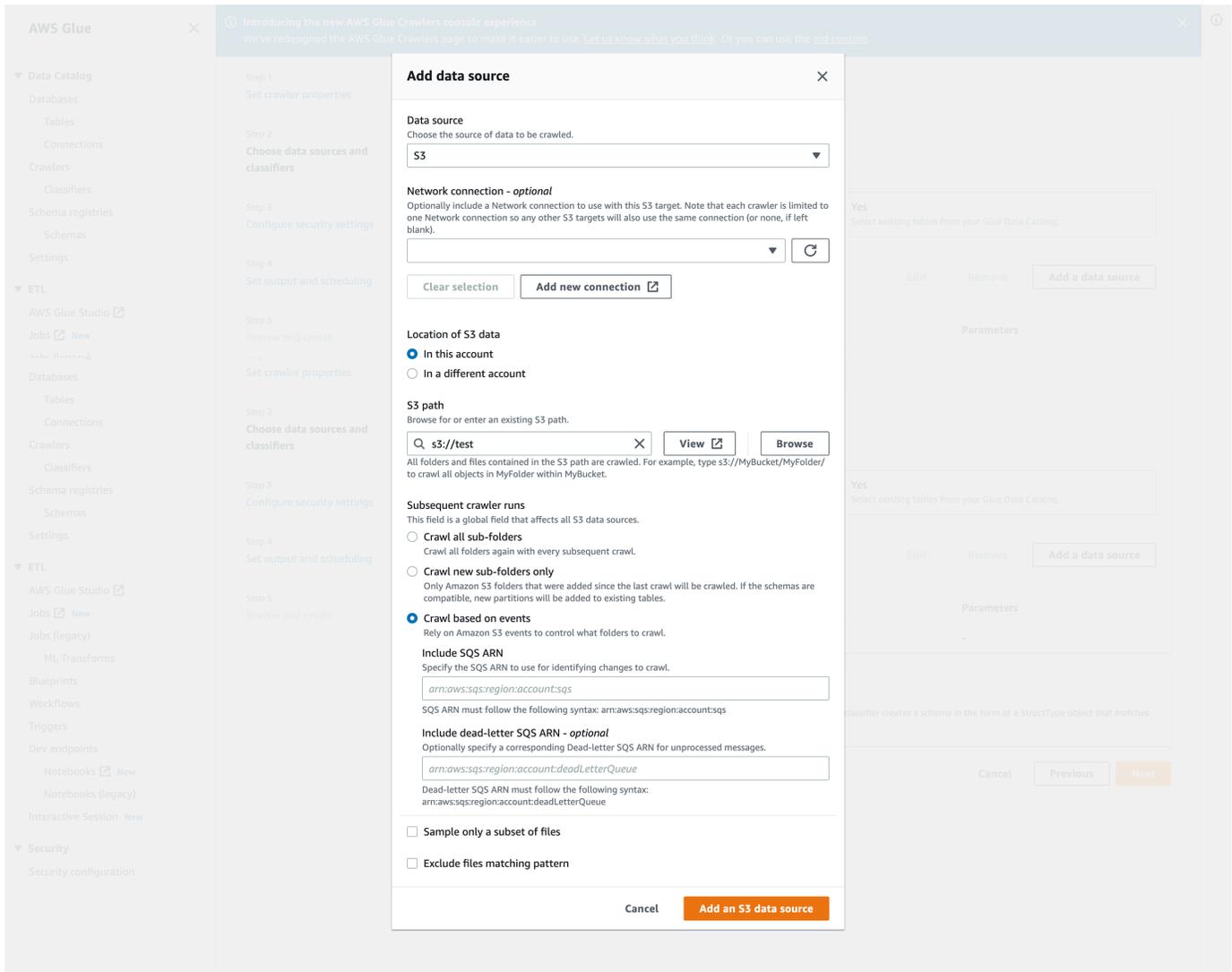
Add a data source

► **Custom classifiers - optional**
A classifier checks whether a given file is in a format the crawler can handle. If it is, the classifier creates a schema in the form of a StructType object that matches that data format.

Cancel Previous Next

4. [Add data source] (データソースの追加) ダイアログで、Amazon S3 データソースを以下のように設定します。
 - [Data source] (データソース): デフォルトで、Amazon S3 が選択されています。
 - [Network connection] (ネットワーク接続) (オプション): [Add new connection] (新しい接続を追加) を選択します。

- [Location of Amazon S3 data] (Amazon S3 データの場所): デフォルトで、[In this account] (このアカウント内) が選択されています。
- [Amazon S3 path] (Amazon S3 パス): フォルダとファイルがクローलされる Amazon S3 パスを指定します。
- [Subsequent crawler runs] (それ以降のクローラー実行): クローラーに関する Amazon S3 イベント通知を使用するには、[Crawl based on events] (イベントに基づくクロール) を選択します。
- [Include SQS ARN] (SQS ARN を含める): 有効な SQS ARN を含むデータストアパラメータを指定します。(例えば、arn:aws:sqs:region:account:sqs)
- [Include dead-letter SQS ARN] (配信不能 SQS ARN を含める) (オプション): 有効な Amazon 配信不能 SQS ARN を指定します。(例えば、arn:aws:sqs:region:account:deadLetterQueue)
- [Add an Amazon S3 data source] (Amazon S3 データソースを追加) を選択します。



AWS CLI を使用した Amazon S3 イベント通知のクローラーの設定

次に示すのは、SQS キューを作成し、Amazon S3 ターゲットバケットでイベント通知を設定するための Amazon S3 AWS CLI コールの例です。

```
S3 Event AWS CLI
aws sqs create-queue --queue-name MyQueue --attributes file://create-queue.json
create-queue.json
...
{
  "Policy": {
```

```

"Version": "2012-10-17",
"Id": "example-ID",
"Statement": [
  {
    "Sid": "example-statement-ID",
    "Effect": "Allow",
    "Principal": {
      "Service": "s3.amazonaws.com"
    },
    "Action": [
      "SQS:SendMessage"
    ],
    "Resource": "SQS-queue-ARN",
    "Condition": {
      "ArnLike": {
        "aws:SourceArn": "arn:aws:s3:*:*:awsexamplebucket1"
      },
      "StringEquals": {
        "aws:SourceAccount": "bucket-owner-account-id"
      }
    }
  }
]
}
...
aws s3api put-bucket-notification-configuration --bucket customer-data-pdx --
notification-configuration file://s3-event-config.json
s3-event-config.json
...
{
  "QueueConfigurations": [
    {
      "Id": "s3event-sqs-queue",
      "QueueArn": "arn:aws:sqs:{region}:{account}:queuename",
      "Events": [
        "s3:ObjectCreated:*",
        "s3:ObjectRemoved:*"
      ],
      "Filter": {
        "Key": {
          "FilterRules": [
            {
              "Name": "Prefix",

```

```
        "Value": "/json"
      }
    ]
  }
}
...
Create Crawler:
```

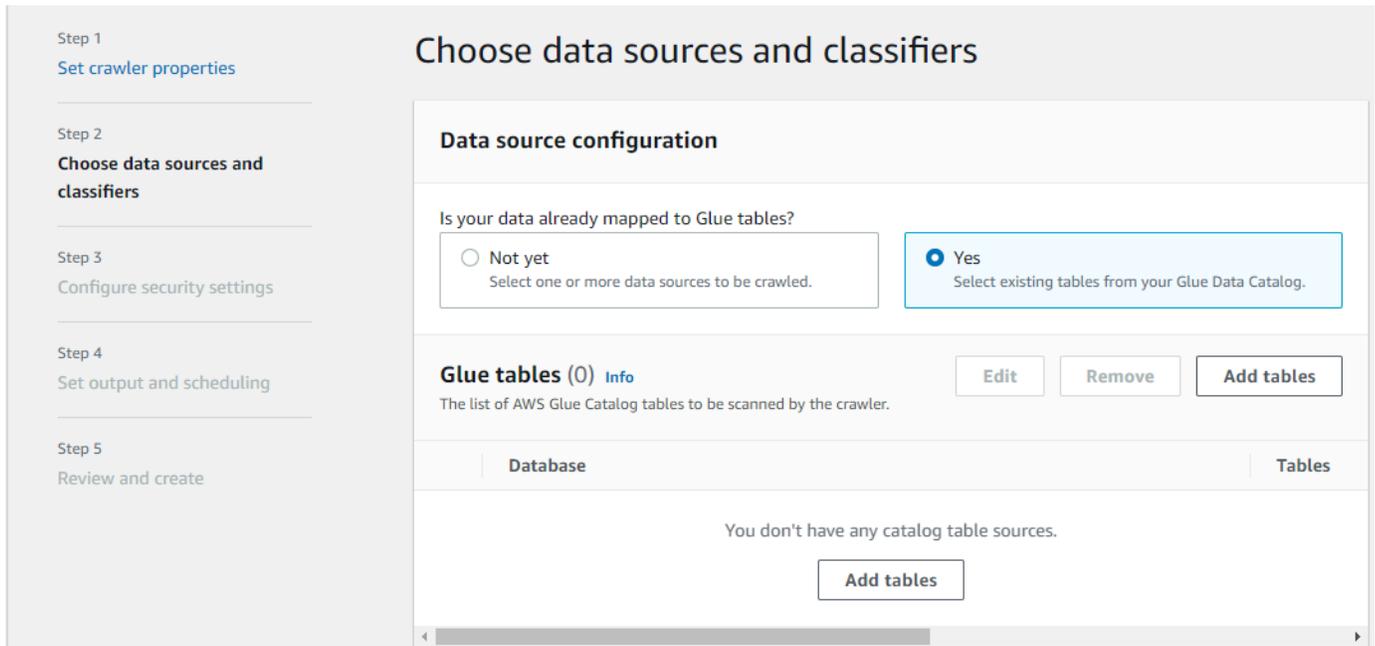
コンソールを使用した Amazon S3 イベント通知用のクローラーの設定 (Data Catalog ターゲット)

カタログターゲットがあるときは、AWS Glue コンソールを使用して Amazon S3 イベント通知用のクローラーを設定します。

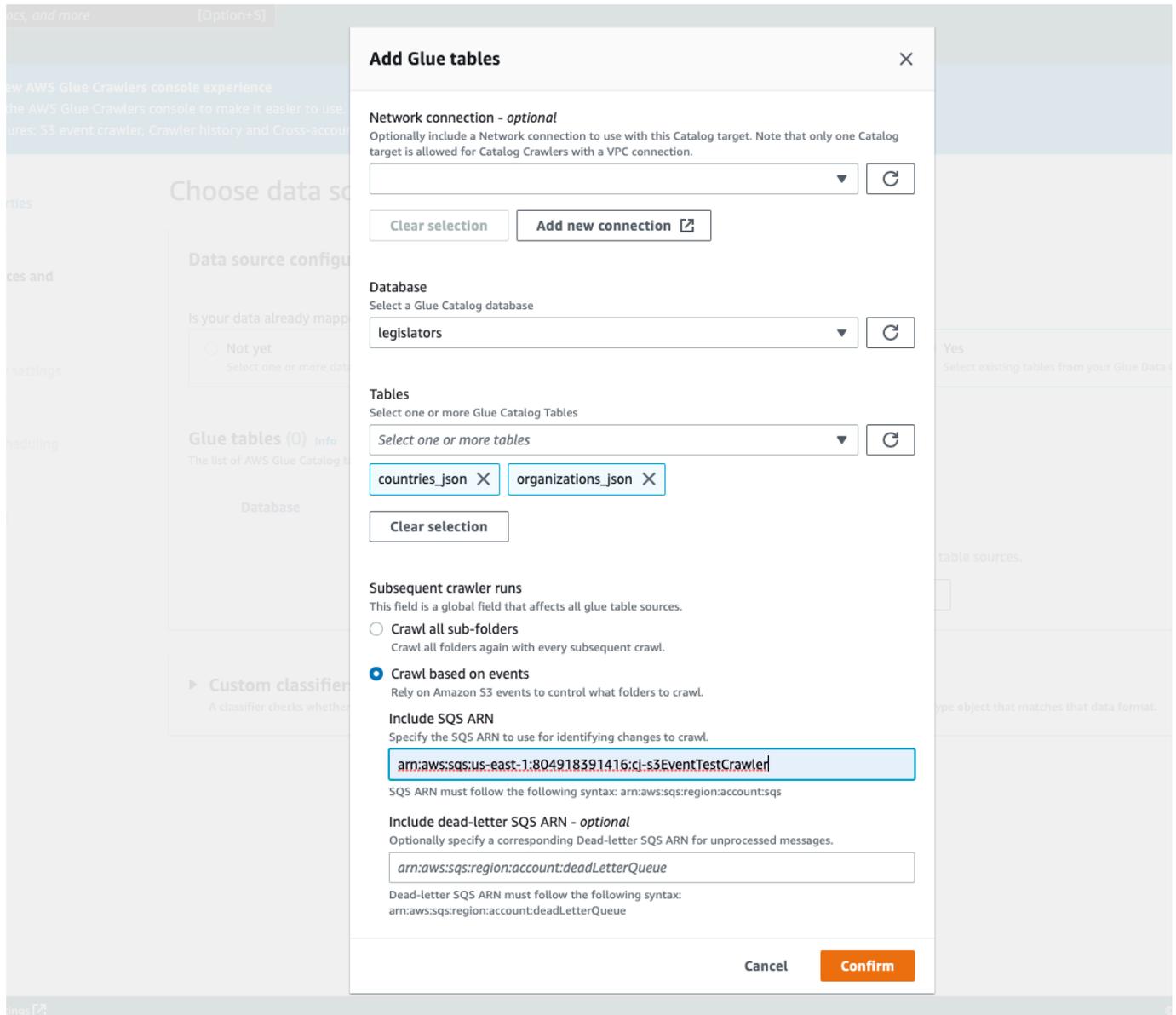
1. クローラーのプロパティを設定します。詳細については、「[AWS Glue コンソールでのクローラー設定オプションの設定](#)」を参照してください。
2. [データソースの設定] セクションに、[データは AWS Glue テーブルにマッピング済みですか?] という質問が表示されています。

[Yes] (はい) を選択して、Data Catalog からの既存のテーブルをデータソースとして選択します。

3. [Glue tables] (Glue テーブル) セクションで、[Add tables] (テーブルを追加する) を選択します。



4. [Add table] (テーブルを追加する) モーダルで、データベースとテーブルを設定します。
 - [Network connection] (ネットワーク接続) (オプション): [Add new connection] (新しい接続を追加) を選択します。
 - [Database] (データベース): Data Catalog 内のデータベースを選択します。
 - [Tables] (テーブル): Data Catalog 内のデータベースから 1 つ、または複数のテーブルを選択します。
 - [Subsequent crawler runs] (それ以降のクローラー実行): クローラーに関する Amazon S3 イベント通知を使用するには、[Crawl based on events] (イベントに基づくクローラ) を選択します。
 - [Include SQS ARN] (SQS ARN を含める): 有効な SQS ARN を含むデータストアパラメータを指定します。(例えば、arn:aws:sqs:region:account:sqs)
 - [Include dead-letter SQS ARN] (配信不能 SQS ARN を含める) (オプション): 有効な Amazon 配信不能 SQS ARN を指定します。(例えば、arn:aws:sqs:region:account:deadLetterQueue)
 - [確認] を選択します。



Amazon S3 イベントクローラで暗号化を使用します。

このセクションでは SQS のみ、または SQS と Amazon S3 の両方での暗号化の使用について説明します。

トピック

- [SQS でのみ暗号化を有効にします](#)
- [SQS と S3 の両方での暗号化の有効にします。](#)
- [よくある質問](#)

SQS でのみ暗号化を有効にします

Amazon SQS は、デフォルトで送信中の暗号化を提供します。オプションのサーバー側の暗号化 (SSE) をキューに追加するため、編集パネルの[カスタマーマスターキー \(CMK\)](#) にアタッチします。つまり、SQS は SQS サーバー上のすべての保管中の顧客データを暗号化します。

カスタマーマスターキー (CMK) を作成します

1. [キー管理服务 (KMS)] > [カスタマーマネージドキー > [キーを作成する] を選択します。
2. 手順に従って、独自のエイリアスと説明を追加します。
3. このキーを使用できるようにする IAM ロールを追加します。
4. キーポリシーで、別のステートメントを「ステートメント」リストに追加することで、[カスタムキーポリシー](#)は Amazon SNS に十分なキー使用許可を付与します。

```
"Statement": [  
  {  
    "Effect": "Allow",  
    "Principal": {  
      "Service": "sns.amazonaws.com"  
    },  
    "Action": [  
      "kms:GenerateDataKey",  
      "kms:Decrypt"  
    ],  
    "Resource": "*"   
  }  
]
```

キューのサーバー側の暗号化 (SSE) を有効にします。

1. [Amazon SQS] > [キュー] > [sqs_queue_name] > [暗号化] タブを選択します。
2. [編集] を選択し、[暗号化] まで下にスクロールしてドロップダウンします。
3. [有効] を選択し、SSE を追加します。
4. 名前 alias/aws/sqs のデフォルトキーではなく、前に作成した CMK をクリックします。

▼ Encryption - *Optional*

Amazon SQS provides in-transit encryption by default. To add at-rest encryption to your queue, enable server-side encryption. [Info](#)

Server-side encryption

- Disabled
 Enabled

Customer master key [Info](#)

alias/sqs-key ▼

これを追加すると、追加したキーで暗号化タブを更新します。

The screenshot shows the AWS Glue console interface. At the top, there are navigation tabs: SNS subscriptions, Lambda triggers, Dead-letter queue, Monitoring, Tagging, Access policy, and Encryption (which is highlighted). Below the tabs, the 'Encryption' section is displayed. It includes an 'Edit' button in the top right corner. The main content area shows 'CMK alias' set to 'alias/sqs-key' and 'Data key reuse period' set to '5 Minutes'.

Note

Amazon SQS は最長メッセージ保持期間を超えてキューに残っているメッセージを自動的に削除します。デフォルトのメッセージ保持期間は 4 日間です。イベントの欠落を回避するには、SQS MessageRetentionPeriod を 14 日以内に変更します。

SQS と S3 の両方での暗号化の有効にします。

SQS でサーバー側の暗号化 (SSE) を有効にします。

1. 「[the section called “SQS でのみ暗号化を有効にします”](#)」の手順を実行します。
2. CMK セットアップの最後のステップで、Amazon S3 に十分なキー使用許可を付与します。

「ステートメント」リストに以下を貼り付けます。

```
"Statement": [
  {
    "Effect": "Allow",
    "Principal": {
      "Service": "s3.amazonaws.com"
```

```
    },
    "Action": [
      "kms:GenerateDataKey",
      "kms:Decrypt"
    ],
    "Resource": "*"
  }
]
```

Amazon S3 バケットのサーバー側の暗号化 (SSE) を有効にします。

1. 「[the section called “SQS でのみ暗号化を有効にします”](#)」の手順を実行します。
2. 次のいずれかを行います。
 - S3 バケット全体で SSE を有効にするには、ターゲットバケットの [プロパティ] タブを操作します。

ここで SSE を有効にして、使用する暗号化の種類を選択できます。Amazon S3 は、Amazon S3 が作成、管理、使用する暗号化キーを提供したり、KMS からキーを選択することもできます。

Edit default encryption

Default encryption

Automatically encrypt new objects stored in this bucket. [Learn more](#) 

Server-side encryption

- Disable
 Enable

Encryption key type

To upload an object with a customer-provided encryption key (SSE-C), use the AWS CLI, AWS SDK, or Amazon S3 REST API.

- Amazon S3 key (SSE-S3)**
An encryption key that Amazon S3 creates, manages, and uses for you. [Learn more](#) 
- AWS Key Management Service key (SSE-KMS)**
An encryption key protected by AWS Key Management Service (AWS KMS). [Learn more](#) 

Cancel

Save changes

- 特定のフォルダで SSE を有効にするには、ターゲットフォルダの横にある [チェックボックス] を選択し、[アクション] をドロップダウンして [サーバー側の暗号化の編集] を選択します。

The screenshot shows the AWS S3 console interface. At the top, there's a header 'Objects (1)' and a sub-header explaining that objects are fundamental entities stored in Amazon S3. Below this, there are several action buttons: 'Refresh', 'Delete', 'Actions' (with a dropdown arrow), 'Create folder', and 'Upload'. A search bar is also present with the placeholder text 'Find objects by prefix'. Below the search bar, there's a table with columns for selection (checkbox) and name. The 's3events/' folder is selected. The 'Actions' dropdown menu is open, showing various options: 'Open', 'Calculate total size', 'Copy', 'Move', 'Initiate restore', 'Query with S3 Select', 'Download actions' (with sub-options 'Download' and 'Download as'), 'Edit actions' (with sub-options 'Rename object', 'Edit storage class', 'Edit server-side encryption', and 'Edit metadata'). The 'Edit server-side encryption' option is highlighted.

よくある質問

Amazon SNS トピックに公開したメッセージが、サーバー側暗号化 (SSE) が有効になっている登録済みの Amazon SQS キューに配信されないのはなぜですか。

Amazon SQS キューが使用していることを再確認してください。

1. [カスタマーマスターキー \(CMK\)](#) は、カスタマーが管理します。SQS によって提供されるデフォルトではありません。
2. (1) の CMK は [カスタムキーポリシー](#) を含み、Amazon SNS に十分なキー使用許可を付与します。

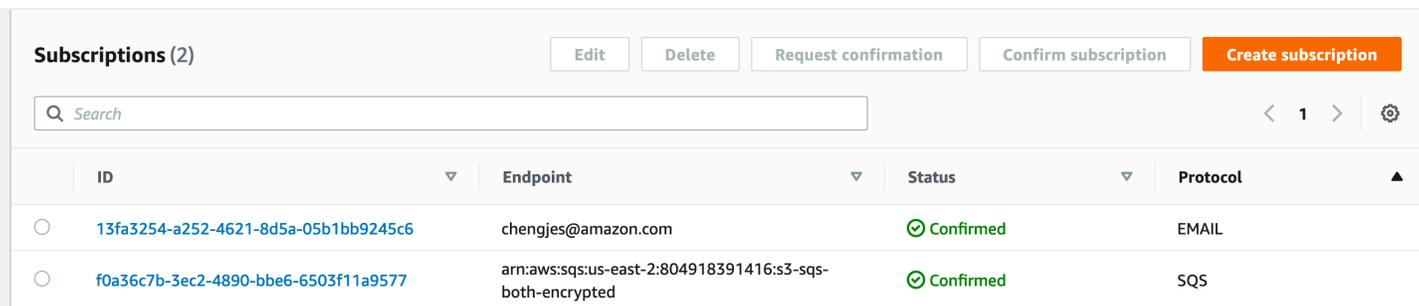
詳細情報は、ナレッジセンターの [この記事](#) を参照してください。

E メール通知の登録をしましたが、Amazon S3 バケットを編集しても E メールの更新が届きません。

メールの「登録確認」リンクをクリックして、メールアドレスを確認したことを確認してください。SNS トピックの下の [サブスクリプション] テーブルで確認ステータスを確認できます。

[Amazon SNS] > [トピック] > [sns_topic_name] > [サブスクリプション] テーブルを選択します。

前提条件のスクリプトに従えば、sns_topic_name はユーザーの sqs_queue_name と同等だと分かります。これは次のように表示されます。



The screenshot shows the 'Subscriptions (2)' page in the AWS console. It includes buttons for 'Edit', 'Delete', 'Request confirmation', 'Confirm subscription', and 'Create subscription'. A search bar is present. Below is a table with columns for ID, Endpoint, Status, and Protocol.

ID	Endpoint	Status	Protocol
13fa3254-a252-4621-8d5a-05b1bb9245c6	chengjes@amazon.com	Confirmed	EMAIL
f0a36c7b-3ec2-4890-bbe6-6503f11a9577	arn:aws:sqs:us-east-2:804918391416:s3-sqs-both-encrypted	Confirmed	SQS

SQS キューでサーバー側の暗号化を有効にした後、追加したフォルダの一部だけがテーブルに表示されます。なぜ parquet が一部欠けているのでしょうか。

SQS キューで SSE を有効にする前に Amazon S3 バケットを変更した場合、クローラによって変更が適用されない可能性があります。S3 バケットにすべての更新を確実にクローラするために、リストモード (「すべてのフォルダをクローラする」) でクローラを再度実行します。もう 1 つのオプションは、S3 イベントを有効にして新しいクローラを作成することで、やり直すことです。

既存のスキーマをクローラーで変更しないための方法

Amazon S3 テーブル定義の既存のフィールドに対する更新をクローラーで上書きしない場合は、コンソールでオプションとして [Add new columns only] (新しい列の追加のみ) を選択するか、設定オプションとして MergeNewColumns を設定します。これはテーブルとパーティションに適用されます (Partitions.AddOrUpdateBehavior を InheritFromTable で上書きしていない場合)。

クローラーの実行時にテーブルスキーマを一切変更しない場合は、スキーマ変更ポリシーを LOG に設定します。設定オプションにより、テーブルから継承するようにパーティションスキーマを設定することもできます。

コンソールでクローラーを設定する場合は、以下のアクションを選択できます。

- 変更を無視し、Data Catalog のテーブルを更新しない
- すべての新規および既存のパーティションを更新してテーブルのメタデータを反映する

API を使用してクローラーを設定する場合は、以下のパラメータを設定します。

- SchemaChangePolicy 構造の UpdateBehavior フィールドを LOG に設定します。
- クローラー API で次の JSON オブジェクトの文字列表現を使用して Configuration フィールドを設定します。

```
{
  "Version": 1.0,
  "CrawlerOutput": {
    "Partitions": { "AddOrUpdateBehavior": "InheritFromTable" }
  }
}
```

各 Amazon S3 インクルードパスの単一のスキーマを作成する方法

デフォルトでは、Amazon S3 に保存されたデータのテーブルをクローラーが定義するときに、データの互換性とスキーマの類似性の両方が考慮されます。考慮されるデータ互換性要因には、データが同じ形式 (JSON など) かどうか、同じ圧縮タイプ (GZIP など) かどうか、Amazon S3 パスの構造、他のデータ属性などです。スキーマの類似性は、別個の Amazon S3 オブジェクトのスキーマがどれくらい似ているかの尺度です。

可能な場合は、CombineCompatibleSchemas へのクローラーを共通テーブル定義に設定できます。このオプションを使用しても、クローラーはデータ互換性を考慮に入れますが、指定されたインクルードパスで Amazon S3 オブジェクトを評価するときに特定のスキーマの類似性を無視します。

コンソールでクローラーを設定する場合、スキーマを組み合わせるには、クローラーオプション [Create a single schema for each S3 path (各 S3 パスの単一のスキーマを作成する)] を選択します。

API を使用してクローラーを設定する場合は、以下の設定オプションを設定します。

- クローラー API で次の JSON オブジェクトの文字列表現を使用して Configuration フィールドを設定します。

```
{
  "Version": 1.0,
  "Grouping": {
    "TableGroupingPolicy": "CombineCompatibleSchemas" }
}
```

このオプションを説明するため、インクルードパス `s3://bucket/table1/` を使用してクローラーを定義するとします。クローラーが実行されると、次の特性を持つ 2 つの JSON ファイルが検索されます。

- ファイル 1 – `S3://bucket/table1/year=2017/data1.json`
- ファイルのコンテンツ – `{"A": 1, "B": 2}`
- スキーマ – `A:int, B:int`

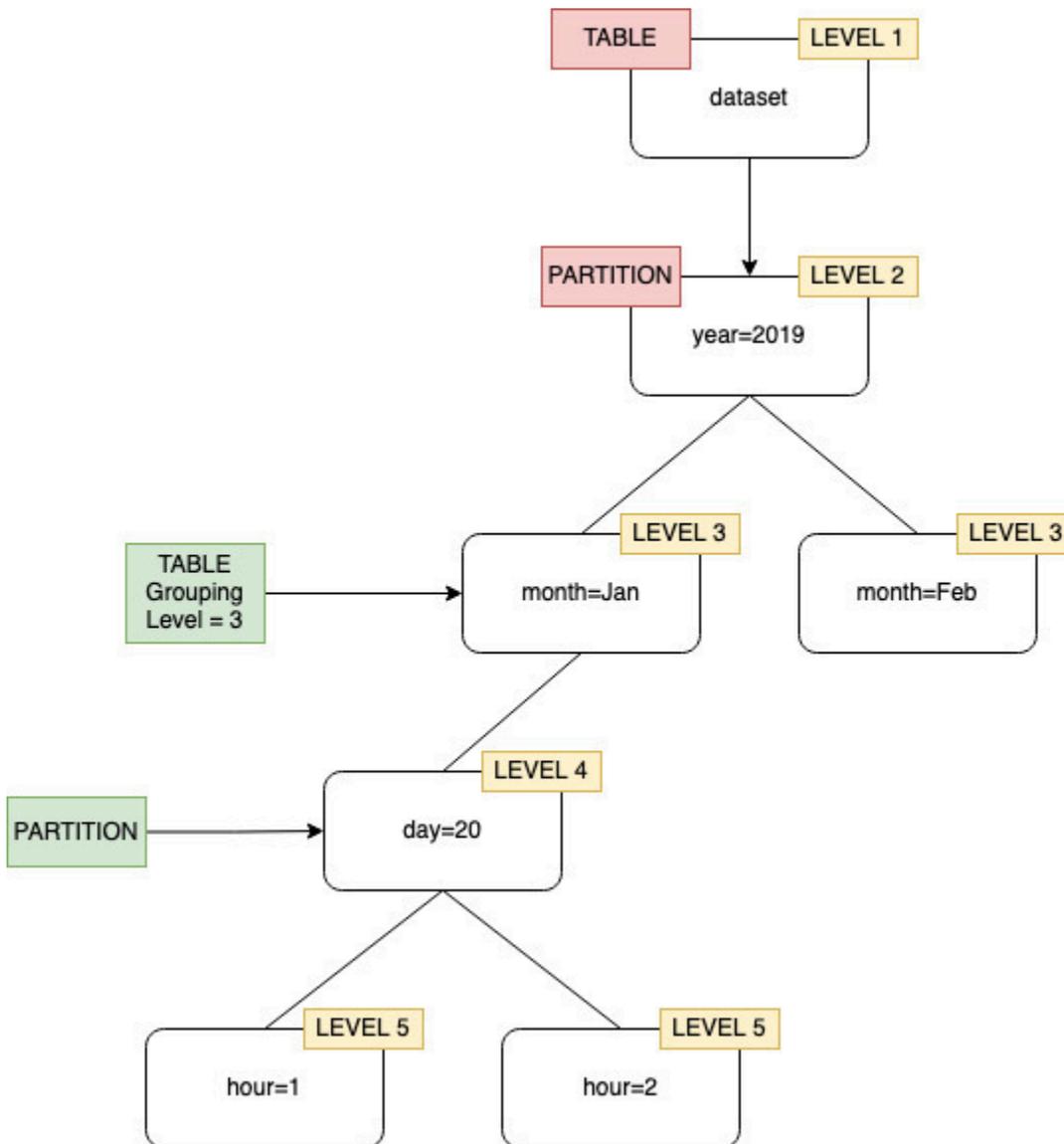
- ファイル 2 – `S3://bucket/table1/year=2018/data2.json`
- ファイルのコンテンツ – `{"C": 3, "D": 4}`
- スキーマ – `C: int, D: int`

デフォルトでは、スキーマの類似性が十分ではないため、クローラーは `year_2017` および `year_2018` という 2 つのテーブルを作成します。ただし、オプション [Create a single schema for each S3 path (各 S3 パスの単一のスキーマを作成する)] を選択し、データに互換性がある場合、クローラーは 1 つのテーブルを作成します。テーブルにはスキーマ `A:int,B:int,C:int,D:int` と `partitionKey year:string` があります。

テーブルの場所とパーティションレベルの指定方法

デフォルトでは、Amazon S3 に保存されたデータのテーブルをクローラーが定義するときに、クローラーはスキーマを結合してトップレベルのテーブル (year=2019) を作成しようとします。場合によっては、フォルダ month=Jan のテーブルをクローラーが作成することを期待することがありますが、兄弟フォルダ (month=Mar) が同じテーブルにマージされているので、代わりにクローラーはパーティションを作成します。

テーブルレベルのクローラーオプションを使用すると、クローラーにテーブルの配置場所やパーティションの作成方法を柔軟に指定できます。テーブルレベルを指定すると、その絶対レベルにAmazon S3 バケットからテーブルが作成されます。



コンソールでクローラーを設定するとき、テーブルレベル クローラーオプションの値を指定できます。値は、テーブルの場所 (データセット内の絶対レベル) を示す正の整数である必要があります。最上位レベルのフォルダのレベルは 1 です。例えば、mydataset/year/month/day/hour というパスで、レベルが 3 に設定されている場合、テーブルは mydataset/year/month という場所に作成されます。

Console

The screenshot shows the 'Set output and scheduling' configuration page in the AWS Glue console. The page is divided into several sections:

- Step 1: Set crawler properties** (selected)
- Step 2: Choose data sources and classifiers**
- Step 3: Configure security settings**
- Step 4: Set output and scheduling** (current step)
- Step 5: Review and create**

The main configuration area is titled 'Set output and scheduling' and contains the following options:

- Output configuration**
 - Target database:** A dropdown menu showing '0-test-catalog' with a refresh icon and buttons for 'Clear selection' and 'Add database'.
 - Table name prefix - optional:** A text input field with a placeholder 'Type a prefix added to table names'.
 - Maximum table threshold - optional:** A text input field with a placeholder 'Type a number greater than 0'. Below it is a note: 'This field sets the maximum number of tables the crawler is allowed to generate. In the event that this number is surpassed, the crawl will fail with an error. If not set, the crawler will automatically generate the number of tables depending on the data schema.'
- Advanced options**
 - S3 schema grouping:**
 - Create a single schema for each S3 path**
By default, when a crawler defines tables for data stored in S3, it considers both data compatibility and schema similarity. Select this check box to group compatible schemas into a single table definition across all S3 objects under the provided include path. Other criteria will still be considered to determine proper grouping.
 - Table level - optional:**
The value must be a positive integer that indicates table location (the absolute level in the dataset). The level for the top level folder is 1. For example, for the path mydataset/a/b, if the level is set to 3, the table is created at location mydataset/a/b.
A dropdown menu shows the value '3'.

API

API を使用してクローラー設定をすると、次の JSON オブジェクトの文字列表現をする Configuration フィールドを設定します。例:

```
configuration = jsonencode(
{
  "Version": 1.0,
  "Grouping": {
    TableLevelConfiguration = 2
  }
})
```

CloudFormation

この例では、CloudFormation テンプレート内のコンソールで使用できる [Table Level] (テーブルレベル) オプションを設定します。

```
"Configuration": "{
  \"Version\":1.0,
  \"Grouping\":{\"TableLevelConfiguration\":2}
}"
```

クローラーが作成できるテーブルの最大数を指定する方法

オプションで、クローラーが作成できるテーブルの最大数を指定する場合は、AWS Glue コンソールまたは CLI 経由で `TableThreshold` を指定します。クローリング中にクローラーによって検出されたテーブルがこの入力値より多い場合、クローリングが失敗し、データカタログにデータが書き込まれません。

このパラメータは、クローラーによって検出および作成されるテーブルが予想よりもはるかに多い場合に役立ちます。これには、以下のような複数の理由が考えられます。

- AWS Glue ジョブを使用して Amazon S3 ロケーションに入力する場合、フォルダと同じレベルに空のファイルができる可能性があります。このような場合、この Amazon S3 ロケーションでクローラーを実行すると、ファイルとフォルダが同じレベルに存在するため、クローラーは複数のテーブルを作成します。
- `"TableGroupingPolicy": "CombineCompatibleSchemas"` を設定しない場合、予想よりも多くのテーブルが作成される可能性があります。

`TableThreshold` を 1 以上の整数値として指定します。この値はクローラーごとに設定されます。つまり、クローリングごとに、この値が考慮されます。例えば、あるクローラーでは `TableThreshold` 値を 5 に設定したとします。各クローリングで、AWS Glue は検出されたテーブルの数をこのテーブルのしきい値 (5) と比較します。検出されたテーブルの数が 5 未満の場合、AWS Glue はデータカタログにテーブルを書き込み、そうでない場合、データカタログに書き込むことなくクローリングが失敗します。

コンソール

AWS コンソールを使用して `TableThreshold` を設定する場合は、次の画面で行います。

Set output and scheduling

Output configuration [Info](#)

Target database

Table name prefix - optional

Maximum table threshold - optional
This field sets the maximum number of tables the crawler is allowed to generate. In the event that this number is surpassed, the crawl will fail with an error. If not set, the crawler will automatically generate the number of tables depending on the data schema.

▶ Advanced options

CLI

AWS CLI を使用して TableThreshold を設定する場合は、次のようになります。

```
"{"Version":1.0,
"CrawlerOutput":
{"Tables":{"AddOrUpdateBehavior":"MergeNewColumns",
"TableThreshold":5}}";
```

エラーメッセージはログに記録され、テーブルパスを特定してデータをクリーンアップするのに役立ちます。テーブル数が指定したテーブルしきい値を超えているためにクローラーが失敗した場合は、次のようなログの例がアカウントで確認できます。

```
Table Threshold value = 28, Tables detected - 29
```

CloudWatch では、検出されたすべてのテーブルの場所を INFO メッセージとしてログに記録します。失敗の理由としてエラーがログに記録されます。

```
ERROR com.amazonaws.services.glue.customerLogs.CustomerLogService - CustomerLogService
received CustomerFacingException with message
The number of tables detected by crawler: 29 is greater than the table threshold value
provided: 28. Failing crawler without writing to Data Catalog.
com.amazonaws.services.glue.exceptions.CustomerFacingInternalException: The number of
tables detected by crawler: 29 is greater than the table threshold value provided:
28.
Failing crawler without writing to Data Catalog.
```

Delta Lake データストアの設定オプションを指定する方法

Delta Lake データストアのクローラーを設定する場合は、次の設定パラメータを指定します。

Connection

任意で、この Amazon S3 ターゲットで使用するネットワーク接続を選択または追加します。接続の詳細については、「[データへの接続](#)」を参照してください

クエリ用テーブルの作成

Delta Lake テーブルの作成方法を次の中から選択します。

- [Create Native tables] (ネイティブテーブルの作成): Delta トランザクションログの直接クエリをサポートするクエリエンジンとの統合を可能にします。
- [Create Symlink tables] (Symlink テーブルの作成): 指定された設定パラメータに基づいて、マニフェストファイルを分割キーで分割して Symlink マニフェストを作成します。

マニフェストの書き込みを有効にします (Delta Lake ソースで [Create Symlink tables] (Symlink テーブルの作成) を選択している場合のみ設定可能)。

Delta Lake のトランザクションログでテーブルメタデータまたはスキーマの変更を検出するかどうかを選択します。マニフェストファイルを再生成します。Delta Lake SET TBLPROPERTIES で自動マニフェスト更新を構成している場合は、このオプションを選択しないでください。

Delta Lake のテーブルパスを含める

Delta テーブルへの Amazon S3 パスを `s3://bucket/prefix/object` として 1 つ以上指定します。

Add data source ✕

Data source
Choose the source of data to be crawled.

Delta Lake ▼

Connection - optional
Select a connection to access the data sources below.

▼ ↺

Clear selection Add new connection [↗](#)

Include delta lake table paths
Browse for or enter an existing S3 path.

`s3://bucket/prefix/object` Remove

Add new delta table path

Enable write manifest
When enabled, if the crawler detects table metadata or schema changes in the Delta Lake transaction log, it regenerates the manifest file. You should not choose this option if you configured automatic manifest updates with Delta Lake SET TBLPROPERTIES.

Cancel Add a Delta Lake data source

Lake Formation の認証情報を使用するようにクローラーを設定する方法

AWS Lake Formation 認証情報を使用するようにクローラーを設定すると、同じ AWS アカウント または別の AWS アカウント 内の基盤となる Amazon S3 ロケーションで Amazon S3 データストアまたはデータカタログテーブルにアクセスできます。クローラーとデータカタログテーブルが同じア

ウントに存在する場合は、既存のデータカタログテーブルをクローラーのターゲットとして設定できます。現在、データカタログテーブルをクローラーのターゲットとして使用する場合、単一のカタログテーブルを含む単一のカタログターゲットのみが許可されています。

Note

データカタログテーブルをクローラーターゲットとして定義する場合、データカタログテーブルの基盤となるロケーションが Amazon S3 ロケーションであることを確認してください。Lake Formation の認証情報を使用するクローラーは、基盤となる Amazon S3 ロケーションを備えたデータカタログターゲットのみをサポートしています。

クローラーと登録された Amazon S3 ロケーションまたはデータカタログテーブルが同じアカウントに存在する場合は、セットアップが必要になります (アカウント内クローリング)

クローラーが Lake Formation の認証情報を使用してデータストアまたはデータカタログテーブルにアクセスできるようにするには、データロケーションを Lake Formation で登録する必要があります。また、クローラーの IAM ロールには、Amazon S3 バケットが登録されている送信先からデータを読み込む許可が付与されている必要があります。

AWS Management Console または AWS Command Line Interface (AWS CLI) を使用して、次の設定手順を完了できます。

AWS Management Console

- クローラーソースにアクセスするようにクローラーを設定する前に、データストアまたはデータカタログのデータロケーションを Lake Formation で登録します。Lake Formation コンソール (<https://console.aws.amazon.com/lakeformation/>) で、クローラーが定義されている AWS アカウントで Amazon S3 ロケーションをデータレイクのルートロケーションとして登録します。詳細については、「[Amazon S3 ロケーションの登録](#)」を参照してください。
- クローラーの実行に使用される IAM ロールにデータロケーション許可を付与し、クローラーが Lake Formation の送信先からデータを読み込めるようにします。詳細については、「[データロケーション許可の付与 \(同じアカウント\)](#)」を参照してください。
- クローラーロールにデータベースへのアクセス権限 (Create) を付与します。このデータベースは出力データベースとして指定されています。詳細については、「[Lake Formation コンソールと名前付きリソース方式を使用したデータベース許可の付与](#)」を参照してください。
- IAM コンソール (<https://console.aws.amazon.com/iam/>) で、クローラー用の IAM ロールを作成します。ロールに lakeformation:GetDataAccess ポリシーを追加します。

5. AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) で、クローラーを設定する際に、オプション [Use Lake Formation credentials for crawling Amazon S3 data source] (Amazon S3 データソースのクローリングに Lake Formation の認証情報を使用する) を選択してください。

 Note

accountId フィールドは、アカウント内クローリングのためのオプションです。

AWS CLI

```
aws glue --profile demo create-crawler --debug --cli-input-json '{
  "Name": "prod-test-crawler",
  "Role": "arn:aws:iam::111122223333:role/service-role/AWSGlueServiceRole-prod-
test-run-role",
  "DatabaseName": "prod-run-db",
  "Description": "",
  "Targets": {
    "S3Targets": [
      {
        "Path": "s3://crawl-testbucket"
      }
    ]
  },
  "SchemaChangePolicy": {
    "UpdateBehavior": "LOG",
    "DeleteBehavior": "LOG"
  },
  "RecrawlPolicy": {
    "RecrawlBehavior": "CRAWL_EVERYTHING"
  },
  "LineageConfiguration": {
    "CrawlerLineageSettings": "DISABLE"
  },
  "LakeFormationConfiguration": {
    "UseLakeFormationCredentials": true,
    "AccountId": "111122223333"
  },
  "Configuration": {
    "Version": 1.0,
    "CrawlerOutput": {
```

```
        "Partitions": { "AddOrUpdateBehavior": "InheritFromTable" },
        "Tables": {"AddOrUpdateBehavior": "MergeNewColumns" }
    },
    "Grouping": { "TableGroupingPolicy": "CombineCompatibleSchemas" }
},
"CrawlerSecurityConfiguration": "",
"Tags": {
    "KeyName": ""
}
}'
```

クローラーと登録された Amazon S3 ロケーションが異なるアカウントに存在する場合は、セットアップが必要になります (クロスアカウントクローリング)

クローラーが Lake Formation の認証情報を使用して別のアカウントのデータストアにアクセスできるようにするには、まず Amazon S3 のデータロケーションを Lake Formation で登録する必要があります。次に、以下の手順を実行して、クローラーのアカウントにデータロケーション許可を付与します。

AWS Management Console または AWS CLI を使用して、次のステップを実行できます。

AWS Management Console

- Amazon S3 ロケーションが登録されているアカウント (アカウント B) では、次の操作を行います。
 - Lake Formation で Amazon S3 パスを登録します。詳細については、「[Amazon S3 ロケーションの登録](#)」を参照してください。
 - クローラーを実行するアカウント (アカウント A) にデータロケーション許可を付与します。詳細については、「[Grant data location permissions](#)」(データロケーション許可を付与する) を参照してください。
 - 基盤となるロケーションをターゲットの Amazon S3 ロケーションとして、Lake Formation に空のデータベースを作成します。詳細については、「[データベースの作成](#)」を参照してください。
 - アカウント A (クローラーを実行するアカウント) に、前のステップで作成したデータベースへのアクセス権を付与します。詳細については、「[Granting database permissions](#)」(データベース許可の付与) を参照してください。
- クローラーが作成され実行されるアカウント (アカウント A) では、次の操作を行います。

- a. AWS RAM コンソールを使用して、外部アカウント (アカウント B) から共有されたデータベースを受け入れます。詳細については、「[AWS Resource Access Manager からのリソース共有招待の承諾](#)」を参照してください。
- b. クローラー用の IAM ロールを作成します。ロールに `lakeformation:GetDataAccess` ポリシーを追加します。
- c. Lake Formation コンソール (<https://console.aws.amazon.com/lakeformation/>) で、ターゲットの Amazon S3 ロケーションに対するデータロケーション許可をクローラーの実行に使用される IAM ロールに付与して、クローラーが Lake Formation の送信先からデータを読み込めるようにします。詳細については、「[Granting data location permissions](#)」(データロケーション許可の付与) を参照してください。
- d. 共有データベースでリソースリンクを作成します。詳細については、「[Create a resource link](#)」(リソースリンクを作成する) を参照してください。
- e. クローラーロールに共有データベースと (Describe) リソースリンクに対するアクセス許可 (Create) を付与します。リソースリンクはクローラーの出力で指定されます。
- f. AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) で、クローラーを設定する際に、オプション [Use Lake Formation credentials for crawling Amazon S3 data source] (Amazon S3 データソースのクローリングに Lake Formation の認証情報を使用する) を選択してください。

クロスアカウントクローリングの場合は、ターゲットの Amazon S3 ロケーションが Lake Formation で登録されている AWS アカウント ID を指定します。アカウント内クローリングの場合、`accountId` フィールドはオプションです。

Step 1
Set crawler properties

Step 2
Choose data sources and classifiers

Step 3
Configure security settings

Step 4
Set output and scheduling

Step 5
Review and create

Configure security settings

IAM role

Existing IAM role

↻
View ↗

Create new IAM role
Update chosen IAM role

Only IAM roles created by the AWS Glue console and have the prefix "AWSGlueServiceRole-" can be updated.

Lake Formation configuration - optional

Allow the crawler to use Lake Formation credentials for crawling the data source.

Use Lake Formation credentials for crawling S3 data source

Checking this box will allow the crawler to use Lake Formation credentials for crawling the data source. If the data source belongs to another account, you must provide the registered account ID. Otherwise, the crawler will crawl only those data sources associated to the account. Only applicable to S3 and Glue Catalog data sources.

Location of S3 data

In this account

In a different account

Account ID

111111111111

Must be a valid account ID, containing only numbers (0-9) and 12 characters long.

▶ **Security configuration - optional**

Enable at-rest encryption with a security configuration.

Cancel
Previous
Next

AWS CLI

```
aws glue --profile demo create-crawler --debug --cli-input-json '{
  "Name": "prod-test-crawler",
  "Role": "arn:aws:iam::111122223333:role/service-role/AWSGlueServiceRole-prod-
test-run-role",
  "DatabaseName": "prod-run-db",
  "Description": "",
  "Targets": {
    "S3Targets": [
      {
        "Path": "s3://crawl-testbucket"
      }
    ]
  },
  "SchemaChangePolicy": {
    "UpdateBehavior": "LOG",
    "DeleteBehavior": "LOG"
  },
  "RecrawlPolicy": {
    "RecrawlBehavior": "CRAWL_EVERYTHING"
  }
},
```

```
"LineageConfiguration": {
  "CrawlerLineageSettings": "DISABLE"
},
"LakeFormationConfiguration": {
  "UseLakeFormationCredentials": true,
  "AccountId": "111111111111"
},
"Configuration": {
  "Version": 1.0,
  "CrawlerOutput": {
    "Partitions": { "AddOrUpdateBehavior": "InheritFromTable" },
    "Tables": { "AddOrUpdateBehavior": "MergeNewColumns" }
  },
  "Grouping": { "TableGroupingPolicy": "CombineCompatibleSchemas" }
},
"CrawlerSecurityConfiguration": "",
"Tags": {
  "KeyName": ""
}
}'
```

Note

- Lake Formation 認証情報を使用するクローラーは、Amazon S3 およびデータカタログターゲットでのみサポートされます。
- Lake Formation 認証情報供給を使用するターゲットの場合、基盤となる Amazon S3 ロケーションは同じバケットに属している必要があります。例えば、すべてのターゲットロケーションが同じバケット (bucket1) にある限り、ユーザーは複数のターゲット (s3://bucket1/folder1、s3://bucket1/folder2) を使用できます。別のバケット (s3://bucket1/folder1、s3://bucket2/folder2) を指定することはできません。
- 現在、データカタログのターゲットクローラーでは、単一のカatalogテーブルを含む単一のカatalogターゲットのみが許可されています。

チュートリアル: AWS Glue クローラーの追加

この AWS Glue シナリオでは、主要な航空会社の到着データを分析して、各月の出発空港の人気度を計算するよう求められています。2016 年のフライトデータが CSV 形式で Amazon S3 に保存され

ています。データの変換と分析を行う前に、AWS Glue Data Catalog 内のそのメタデータをカタログ化します。

このチュートリアルでは、これらのフライトログからメタデータを推測し、データ Amazon S3 にテーブルを作成するクローラを追加します。

トピック

- [前提条件](#)
- [ステップ 1: クローラの追加](#)
- [ステップ 2: クローラを実行する](#)
- [ステップ 3: AWS Glue Data Catalog オブジェクトを表示する](#)

前提条件

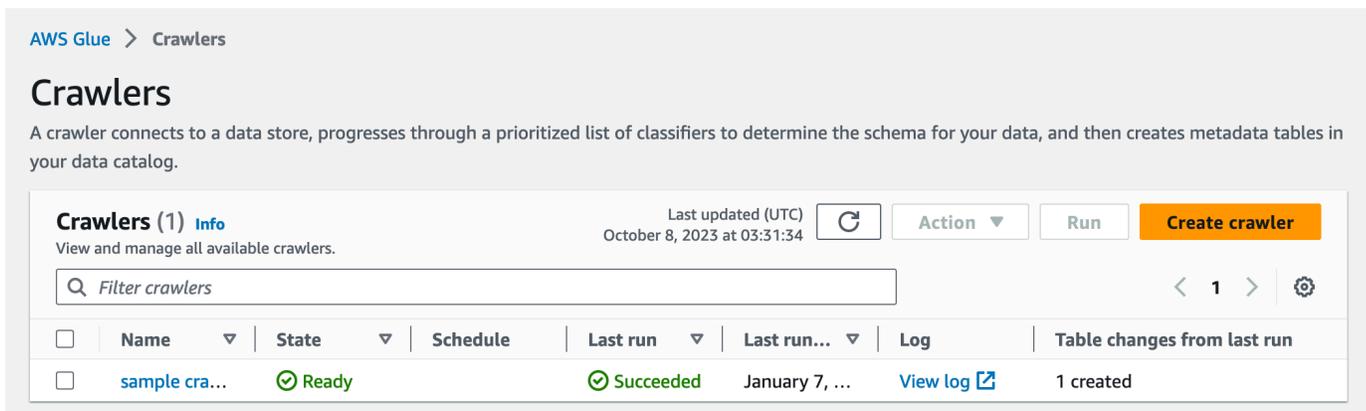
このチュートリアルでは、AWS アカウントを持ち、AWS Glue にアクセスできることを前提としています。

ステップ 1: クローラの追加

Amazon S3 に保存されている CSV ファイルからメタデータを抽出するクローラを設定および実行するには、以下の手順に従ってください。

Amazon S3 に保存されているファイルを読み込むクローラを作成するには

1. AWS Glue サービスコンソールの左側のメニューで、[Crawlers] (クローラー) を選択します。
2. クローラーページで、[Create crawler] を選択します。これにより、クローラの詳細に関する入力を求める一連のページが表示されます。



The screenshot shows the AWS Glue console interface for the 'Crawlers' section. At the top, there is a breadcrumb 'AWS Glue > Crawlers' and a title 'Crawlers'. Below the title is a brief description: 'A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.' The main content area features a 'Crawlers (1) Info' header with a refresh button and a 'Create crawler' button. Below this is a search bar labeled 'Filter crawlers' and a table listing the available crawlers. The table has columns for Name, State, Schedule, Last run, Last run..., Log, and Table changes from last run. One crawler named 'sample cra...' is listed with a state of 'Ready', a last run status of 'Succeeded', and a log link.

<input type="checkbox"/>	Name	State	Schedule	Last run	Last run...	Log	Table changes from last run
<input type="checkbox"/>	sample cra...	Ready		Succeeded	January 7, ...	View log	1 created

3. [クローラ名] フィールドに、「**Flights Data Crawler**」を入力して、[Next] (次へ) を選択します。

クローラは分類子を呼び出して、データのスキーマを推測します。このチュートリアルでは、デフォルトの CSV 用の組み込みの分類子を使用します。

4. [クローラソースタイプ] で、[Data stores] (データストア) を選択し、[Next] (次へ) を選択します。
5. それでは、クローラをデータに指定します。[Add a data store] (データストアの追加) ページで、[Amazon S3 data store] (Amazon S3 データストア) を選択します。このチュートリアルでは接続を使用しないため、[Connection] (接続) フィールドが表示されている場合は、そのフィールドを空白のままにしてください。

[Crawl data in] (クローラするデータの場所) で、[Specified path in another account] (別のアカウントで指定されたパス) を選択します。次に、[Include path] (インクルードパス) に、クローラがフライトデータ (`s3://crawler-public-us-east-1/flight/2016/csv`) に対するパスを入力します。パスを入力すると、このフィールドのタイトルが [Include path] (インクルードパス) に変わります。[Next] を選択します。

6. 単一のクローラで複数のデータストアをクローラできます。ただし、このチュートリアルでは 1 つのデータストアのみを使用しているため、[No] (なし) を選択して、次に [Next] (次へ) を選択します。
7. クローラーには、データストアにアクセスして AWS Glue Data Catalog でオブジェクトを作成するためのアクセス許可が必要です。これらのアクセス許可を設定するために、[Create an IAM role] (IAM ロールを作成する) を選択します。IAM ロール名は「AWSGlueServiceRole-」で始まります。フィールドにロール名の後ろの部分を入力します。「**CrawlerTutorial**」と入力し、[Next] (次へ) を選択します。

Note

IAM ロールを作成するには、AWS ユーザーは CreateRole、CreatePolicy、および AttachRolePolicy アクセス許可が必要です。

ウィザードが「AWSGlueServiceRole-CrawlerTutorial」という名前の IAM ロールを作成し、AWS 管理ポリシー AWSGlueServiceRole をこのロールにアタッチし、Amazon S3 ロケーション `s3://crawler-public-us-east-1/flight/2016/csv` への読み取りアクセスを許可するインラインポリシーを追加します。

8. クローラのスケジュールを作成します。[Frequency] (頻度) で、[Run on demand (オンデマンドで実行する)] (オンデマンドで実行する) を選択してから、[Next] (次へ) を選択します。

9. クローラは Data Catalog 内にテーブルを作成します。テーブルは、Data Catalog 内のデータベースに格納されます。まず、Add database (データベースを追加する) を選択してデータベースを作成します。ポップアップウィンドウで、データベース名に「**test-flights-db**」と入力し、[Create] (作成) を選択します。

次に、[Prefix added to tables] (テーブルに追加されたプレフィックス) に「**flights**」を入力します。残りのフィールドにはデフォルト値を使用し、[Next] (次へ) を選択します。

10. 選択した内容を確認するには、[Add crawler] (クローラの追加) ウィザードを使用します。誤りを見つけた場合は、Back (戻る) をクリックして前のページに戻り、修正します。

情報を確認したら、[Finish] (完了) をクリックしてクローラを作成します。

ステップ 2: クローラを実行する

クローラを作成した後、ウィザードがクローラビューページを表示します。オンデマンドのスケジュールでクローラを作成するため、クローラを実行するオプションを使用できます。

クローラを実行するには

1. このページの上部にあるバナーでは、クローラが作成されたことが表示され、それを実行するかどうかを尋ねられます。[Run it now?] (今すぐ実行しますか?) をクリックしてクローラを実行します。

バナーでは、クローラの「Attempting to run (実行しようとしています)」および「Running (実行中)」というメッセージの表示が変わります。クローラの実行を開始すると、バナーが消え、クローラの表示が更新され、クローラの [Starting] (起動) ステータスが表示されます。1 分後、[Refresh] (更新) アイコンをクリックして、テーブルに表示されるクローラのステータスを更新できます。

2. クローラが完了すると、クローラの変更を説明する新しいバナーが表示されます。test-flights-db リンクを選択して、Data Catalog オブジェクトを表示できます。

ステップ 3: AWS Glue Data Catalog オブジェクトを表示する

クローラがソースの場所からデータを読み込み、Data Catalog にテーブルを作成します。テーブルは、スキーマを含むデータを表すメタデータ定義です。Data Catalog のテーブルにはデータが含まれていません。代わりに、ジョブ定義のソースまたはターゲットとしてこれらのテーブルを使用します。

クローラが作成した Data Catalog オブジェクトを表示する方法

1. 左側のナビゲーションペインの [Data catalog] (データカタログ) で、[Databases] (データベース) を選択します。ここでは、クローラによって作成された flights-db データベースを表示できます。
2. 左側のナビゲーションの [Data catalog] (データカタログ) で、[Databases] (データベース) の下の Tables (テーブル) を選択します。ここでは、クローラが作成した flightscsv テーブルを表示できます。テーブル名を選択すると、テーブルの設定、パラメータ、およびプロパティを表示できます。このビューで下にスクロールすると、スキーマを表示できます。スキーマとは、テーブルの列とデータ型に関する情報です。
3. テーブルビューページで [View partitions] (パーティションの表示) を選択すると、データ用に作成されたパーティションを表示できます。最初の列はパーティションキーです。

メタデータの手動定義

AWS Glue データカタログは、データソースとデータセットに関するメタデータを保存する中央リポジトリです。クローラは、サポートされているデータソースのメタデータを自動的にクロールして入力できますが、データカタログでメタデータを手動で定義する必要があるシナリオがいくつかあります。

- サポートされていないデータ形式 – クローラでサポートされていないデータソースがある場合は、データカタログでそれらのデータソースのメタデータを手動で定義する必要があります。
- カスタムメタデータ要件 — AWS Glue クローラは、事前定義されたルールと規則に基づいてメタデータを推測します。AWS Glue クローラで推測されたメタデータでカバーされない特定のメタデータ要件がある場合は、ニーズに合わせてメタデータを手動で定義できます。
- データガバナンスと標準化 — データガバナンス、コンプライアンス、またはセキュリティ上の理由から、メタデータ定義をより細かくコントロールする必要がある場合があります。メタデータを手動で定義することで、そのメタデータが組織の標準とポリシーに準拠していることを保証できます。
- 将来のデータインジェストのためのプレースホルダー – すぐに使用またはアクセスできないデータソースがある場合は、プレースホルダーとして空のスキーマテーブルを作成できます。データソースが利用可能になったら、事前定義された構造を維持しながら、テーブルに実際のデータを入力できます。

メタデータを手動で定義するには、AWS Glue コンソール、Lake Formation コンソール、AWS Glue API、または AWS Command Line Interface (AWS CLI) を使用できます。データベース、テーブル、

パーティションを作成し、列名、データ型、説明、その他の属性などのメタデータプロパティを指定できます。

データベースの作成

データベースは、メタデータテーブルを AWS Glue で組織化するために使用されます。AWS Glue Data Catalog でテーブルを定義すると、データベースに追加します。1 つのテーブルは、1 つのデータベースにしか含まれません。

データベースには、数多くのさまざまなデータストアからのデータを定義するテーブルを含めることができます。このデータには、Amazon Simple Storage Service (Amazon S3) のオブジェクトと、Amazon Relational Database Service のリレーショナルテーブルを含めることができます。

Note

データベースを AWS Glue Data Catalog から削除すると、データベース内のすべてのテーブルも削除されます。

データベースのリストを表示するには、AWS Management Console にサインインして、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。[Databases] (データベース) を選択し、リスト内のデータベース名を選択して詳細を表示します。

コンソールの [DatabasesAWS Glue] (データベース) タブからデータベースの追加、編集、削除ができます。

- 新しいデータベースを作成するには、[Add database] (データベースを追加) を選択し名前と説明を指定します。Apache Hive など、その他のメタデータストアとの互換性を考慮して、名前は小文字に変換されます。

Note

Amazon Athena からデータベースにアクセスする場合は、英数字とアンダースコア文字のみを使用して名前を指定してください。詳細については、[Athena での名前](#)を参照してください。

- データベースの説明を編集するには、データベース名の横にあるチェックボックスを選択し、[Edit] (編集) を選択します。

- データベースを削除するには、データベース名の横にあるチェックボックスを選択し、[Remove] (削除) を選択します。
- データベースに含まれるテーブルのリストを表示するには、データベース名を選択します。そうすると、データベースプロパティにデータベース内のすべてのテーブルが表示されます。

クローラーで書き込みが行われたデータベースを変更するには、クローラー定義を変更する必要があります。詳細については、「[クローラーを使用したデータカタログへの入力](#)」を参照してください。

データベースリソースリンク

AWS Glue コンソールの更新が行われました。コンソールの現行バージョンでは、データベースリソースリンクはサポートされません。

Data Catalog にはデータベースへのリソースリンクを含めることもできます。データベースリソースリンクは、ローカルデータベースまたは共有データベースへのリンクです。現在、AWS Lake Formation のみでリソースリンクを作成できます。データベースへのリソースリンクを作成した後、データベース名を使用する任意の場所で、リソースリンク名を使用できます。所有しているデータベースまたは共有しているデータベースとともに、データベースリソースリンクは `glue:GetDatabases()` から返され、AWS Glue コンソールの [Databases] (データベース) ページにエントリとして表示されます。

Data Catalog には、テーブルリソースリンクを含めることもできます。

リソースリンクの詳細については、AWS Lake Formation デベロッパーガイドの「[Creating Resource Links](#)」を参照してください。

テーブルの作成

データストア内のデータのインベントリを作成するには、クローラーの実行が推奨されますが、メタデータテーブルを手動で AWS Glue Data Catalog に追加できます。このアプローチにより、メタデータ定義をより詳細にコントロールし、特定の要件に応じてカスタマイズできます。

以下の方法で、データカタログにテーブルを手動で追加することもできます。

- AWS Glue コンソールを使用して AWS Glue Data Catalog に手動でテーブルを作成します。詳細については、「[AWS Glue コンソールでのテーブルの使用](#)」を参照してください。

- [AWS Glue API](#) の CreateTable オペレーションを使用し、AWS Glue Data Catalog にテーブルを作成します。詳細については、「[CreateTable アクション \(Python: create_table\)](#)」を参照してください。
- AWS CloudFormation テンプレートを使用します。詳細については、「[AWS Glue の場合は AWS CloudFormation](#)」を参照してください。

コンソールまたは API を使用して手動でテーブルを定義する場合、テーブルスキーマおよびデータソースでデータのタイプとフォーマットを示す分類フィールドの値を指定します。クローラーでテーブルを作成する場合、データ形式とスキーマは、組み込み分類子またはカスタム分類子のいずれかによって決定されます。AWS Glue コンソールを使用してテーブルを作成する方法の詳細については、[AWS Glue コンソールでのテーブルの使用](#) を参照してください。

トピック

- [テーブルパーティション](#)
- [テーブルリソースリンク](#)
- [クローラーを使用して手動で作成されたデータカタログテーブルを更新する](#)
- [Data Catalog テーブルのプロパティ](#)
- [AWS Glue コンソールでのテーブルの使用](#)
- [AWS Glue でのパーティションインデックスの使用](#)

テーブルパーティション

Amazon Simple Storage Service (Amazon S3) フォルダの AWS Glue テーブル定義によって、パーティションテーブルを記述できます。たとえば、クエリのパフォーマンスを向上させるために、パーティションテーブルでは月の名前をキーとして毎月のデータを別のファイルに分割する場合があります。AWS Glue では、テーブル定義にテーブルのパーティションキーが含まれています。AWS Glue は、Amazon S3 フォルダのデータを評価してテーブルをカタログ化するとき、個々のテーブルまたはパーティション化されたテーブルを追加するかどうかを決定します。

テーブル内のすべてのパーティションをロードする代わりに、テーブル上にパーティションインデックスを作成して、パーティションのサブセットを取得できます。パーティションインデックスの使用方法については、「[AWS Glue でのパーティションインデックスの使用](#)」を参照してください。

AWS Glue によって Amazon S3 フォルダのパーティションテーブルを作成するには、次の条件がすべて満たされている必要があります。

- ファイルのスキーマは、AWS Glue によって決定されるものと似ている。
- ファイルのデータ形式が同じである。
- ファイルの圧縮形式が同じである。

例えば、iOS と Android 両方のアプリケーションの販売データを保存する my-app-bucket という名前の Amazon S3 バケットを所有しているとします。データは、年、月、日で分割されます。iOS および Android の販売に関するデータファイルは、同じスキーマ、データ形式、および圧縮形式です。AWS Glue Data Catalog では、AWS Glue クローラーが、年、月、日のパーティションキーを使用して 1 つのテーブル定義を作成します。

次の my-app-bucket の Amazon S3 リストでは、パーティションのいくつかが表示されています。= シンボルは、パーティションキー値の割り当てに使用されます。

```
my-app-bucket/Sales/year=2010/month=feb/day=1/iOS.csv
my-app-bucket/Sales/year=2010/month=feb/day=1/Android.csv
my-app-bucket/Sales/year=2010/month=feb/day=2/iOS.csv
my-app-bucket/Sales/year=2010/month=feb/day=2/Android.csv
...
my-app-bucket/Sales/year=2017/month=feb/day=4/iOS.csv
my-app-bucket/Sales/year=2017/month=feb/day=4/Android.csv
```

テーブルリソースリンク

AWS Glue コンソールの更新が行われました。コンソールの現行バージョンでは、テーブルリソースリンクはサポートされません。

Data Catalog には、テーブルへのリソースリンクを含めることができます。テーブルリソースリンクは、ローカルまたは共有テーブルへのリンクです。現在、AWS Lake Formation のみでリソースリンクを作成できます。テーブルへのリソースリンクを作成すると、テーブル名を使用する任意の場所で、リソースリンク名を使用できます。自分が所有している、または自分と共有しているテーブルに加えて、テーブルリソースリンクは `glue:GetTables()` から返され、AWS Glue コンソールの [Tables] (テーブル) ページにエントリとして表示されます。

Data Catalog には、データベースリソースリンクを含めることもできます。

リソースリンクの詳細については、AWS Lake Formation デベロッパーガイドの「[Creating Resource Links](#)」を参照してください。

クローラーを使用して手動で作成されたデータカタログテーブルを更新する

手動で AWS Glue Data Catalog テーブルを作成し、AWS Glue クローラーを使用して更新された状態を維持することもできます。スケジュールで実行されているクローラーは、新しいパーティションを追加して、スキーマの変更によりテーブルを更新できます。これは Apache Hive メタストアから移行されたテーブルにも適用されます。

これを行うには、クローラーを定義するときに、1 つ以上のデータストアをクローラーのソースとして指定する代わりに、既存のデータカタログテーブルを 1 つ以上指定します。クローラーは、カタログテーブルで指定されたデータストアをクローラーします。この場合、新しいテーブルは作成されず、手動で作成されたテーブルが更新されます。

以下に、カタログテーブルを手動で作成し、クローラーソースとしてカタログテーブルを指定するその他の理由を示します。

- カタログテーブル名命名アルゴリズムに依存せず、カタログテーブル名を選択することをお勧めします。
- パーティションの検出を中断する可能性がある形式のファイルが誤ってデータソースパスに保存されないように、新しいテーブルが作成されないようにする必要があります。

詳細については、「[ステップ 2: データソースと分類子を選択する](#)」を参照してください。

Data Catalog テーブルのプロパティ

AWS CLI で知られているテーブルのプロパティ、つまりパラメータは無効なキーと値の文字列です。AWS Glue の外部での Data Catalog の使用をサポートする場合は、テーブルに独自のプロパティを設定します。この操作は、Data Catalog を使用する他のサービスでも行うことができます。AWS Glue がジョブまたはクローラーの実行時に一部のテーブルプロパティを設定します。特に説明がない限り、これらのプロパティは内部使用のためのものであり、現在の形式での存続や、これらのプロパティを手動で変更した場合の製品の動作についてはサポートされません。

AWS Glue クローラーで設定できるテーブルプロパティの詳細については、「[the section called “クローラーによって設定されたデータカタログテーブルのパラメータ”](#)」を参照してください。

AWS Glue コンソールでのテーブルの使用

AWS Glue Data Catalog のテーブルは、データストア内のデータを表すメタデータ定義です。クローラの実行時にテーブルを作成するか、または、AWS Glue コンソールで手動でテーブルを作成できます。コンソールの [TablesAWS Glue] (テーブル) リストに、テーブルのメタデータの値が表示されます。ETL (抽出、変換、ロード) ジョブを作成するときに、テーブル定義を使用してソースとターゲットを指定します。

Note

AWS マネジメントコンソールの最新の変更に伴い、既存の IAM ロール を [SearchTables](#) 許可を持つように変更することが必要になる場合があります。新しいロールを作成する場合、SearchTables API 許可は既にデフォルトとして追加されています。

開始するには、AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。[Tables] (テーブル) タブをクリックし、[Add tables] (テーブルの追加) ボタンを使用して、クローラで、または属性を手動で入力してテーブルを作成します。

コンソールでテーブルを追加する

クローラを使用してテーブルを追加するには、[Add tables] (テーブルの追加)、[Add tables using a crawler] (クローラを使用してテーブルを追加) の順に選択します。次に、[Add crawler] (クローラの追加) ウィザードの手順に従います。クローラが実行されると、テーブルが AWS Glue Data Catalog に追加されます。詳細については、[クローラーを使用したデータカタログへの入力](#) を参照してください。

Data Catalog の Amazon Simple Storage Service (Amazon S3) テーブル定義の作成に必要な属性が分かっている場合は、テーブルウィザードで作成できます。[Add tables] (テーブルの追加)、[Add table manually] (手動でのテーブルを追加) の順に選択し、[Add table] (テーブルの追加) ウィザードの手順に従います。

コンソールで手動でテーブルを追加するときは、以下の点を考慮します。

- Amazon Athena からテーブルにアクセスする場合は、英数字とアンダースコア文字のみを使用して名前を指定してください。詳細については、「[テーブル、データベース、および列の名前](#)」を参照してください。
- ソースデータの場所は Amazon S3 パスにする必要があります。

- データのデータ形式は、ウィザードに表示されているいずれかの形式と一致する必要があります。対応する分類、SerDe、およびその他のテーブルのプロパティは、選択された形式に基づいて自動的に入力されます。次の形式でテーブルを定義できます。

Avro

Apache Avro JSON バイナリ形式。

CSV

文字で区切られた値。また、区切り文字として、カンマ、パイプ、セミコロン、タブ、または Ctrl-A を指定します。

JSON

JavaScript Object Notation。

XML

Extensible Markup Language 形式。データの行を定義する XML タグを指定します。列は行のタグ内で定義されます。

Parquet

Apache Parquet 列指向ストレージ。

ORC

Optimized Row Columnar (ORC) ファイル形式。Hive データを効率的に保存するために設計された形式。

- テーブルのパーティションキーを定義できます。
- 現在、コンソールで作成した分割されたテーブルは、ETL ジョブで使用することはできません。

テーブル属性

以下に重要なテーブル属性を示します。

名前

名前は、テーブルの作成時に決定され、変更することはできません。多くの AWS Glue オペレーションでテーブル名を参照します。

データベース

テーブルが存在するコンテナオブジェクト。このオブジェクトには、AWS Glue Data Catalog 内に存在するテーブルの組織名が含まれ、データストアの組織名とは異なる場合があります。データベースを削除すると、データベースに含まれるすべてのテーブルも Data Catalog から削除されます。

説明

テーブルの説明。テーブルの内容を理解しやすくするために説明を記入できます。

テーブル形式

標準 AWS Glue テーブル、または Apache Iceberg 形式のテーブルの作成を指定します。

圧縮を有効にする

[圧縮を有効にする] を選択して、テーブル内の小さな Amazon S3 オブジェクトを大きなオブジェクトに圧縮します。

IAM ロール

圧縮を実行するために、サービスはユーザーに代わって IAM ロールを引き受けます。ドロップダウンを使用して IAM ロールを選択できます。圧縮を有効にするために必要な許可がロールに付与されているようにしてください。

IAM ロールに必要な許可の詳細については、「[テーブル最適化の前提条件](#)」を参照してください。

ロケーション

このテーブル定義が表すデータストア内のデータの場所へのポインタ。

分類

テーブルの作成時に指定された分類の値。通常、これはクローラが実行されてソースデータの形式を指定するときに書き込まれます。

最終更新日

Data Catalog でこのテーブルが更新された日付と時刻 (UTC)。

追加された日付

Data Catalog にこのテーブルが追加された日付と時刻 (UTC)。

廃止済み

AWS Glue により、Data Catalog のテーブルは元のデータストアに存在しなくなったことが分かると、そのテーブルは廃止されたとしてデータカタログにマークされます。廃止されたテーブルを参照するジョブを実行する場合、ジョブは失敗する可能性があります。廃止されたテーブルを参照するジョブを編集し、ソースおよびターゲットとして削除します。廃止されたテーブルが不要になったら削除することをお勧めします。

Connection

AWS Glue でデータストアへの接続が必要な場合は、接続の名前がテーブルに関連付けられます。

テーブルの詳細の表示と編集

既存のテーブルの詳細を表示するには、リスト内のテーブル名を選択し、[Action, View details] (アクション、詳細を表示) を選択します。

テーブルの詳細にはテーブルのプロパティとスキーマが含まれます。このビューには、テーブルに定義された順序の列名、データ型、およびパーティションのキー列を含む、テーブルのスキーマが表示されます。列が複合型の場合は、以下の例に示すように、[View properties] (プロパティの表示) を選択して、そのフィールドの構造の詳細を表示します。

```
{
  "StorageDescriptor":
    {
      "cols": {
        "FieldSchema": [
          {
            "name": "primary-1",
            "type": "CHAR",
            "comment": ""
          },
          {
            "name": "second ",
            "type": "STRING",
            "comment": ""
          }
        ]
      },
      "location": "s3://aws-logs-111122223333-us-east-1",
      "inputFormat": "",
    }
}
```

```
"outputFormat": "org.apache.hadoop.hive.q1.io.HiveIgnoreKeyTextOutputFormat",
"compressed": "false",
"numBuckets": "0",
"SerDeInfo": {
  "name": "",
  "serializationLib": "org.apache.hadoop.hive.serde2.OpenCSVSerde",
  "parameters": {
    "separatorChar": "|"
  }
},
"bucketCols": [],
"sortCols": [],
"parameters": {},
"SkewedInfo": {},
"storedAsSubDirectories": "false"
},
"parameters": {
  "classification": "csv"
}
}
```

StorageDescriptor などのテーブルのプロパティの詳細については、「[StorageDescriptor 構造](#)」を参照してください。

テーブルのスキーマを変更するには、[Edit schema] (スキーマの編集) を選択し、列の追加および削除、列名の変更、データ型の変更をします。

スキーマを含め、テーブルの異なるバージョンを比較するには、[Compare versions] (バージョンの比較) を選択し、テーブルの 2 つのバージョンのスキーマを並べて比較します。詳細については、「[テーブルスキーマのバージョンの比較](#)」を参照してください。

Amazon S3 パーティションを構成するファイルを表示するには、[View partition] (パーティションの表示) を選択します。Amazon S3 のテーブルでは、[Key] (キー) 列に、ソースデータストアでテーブルを分割するために使用されるパーティションキーが表示されます。パーティションは、日付、場所、または部門などのキー列の値に基づいて、テーブルを関連する部分に分割する方法です。パーティションの詳細については、インターネットで「hive パーティション」を検索してください。

Note

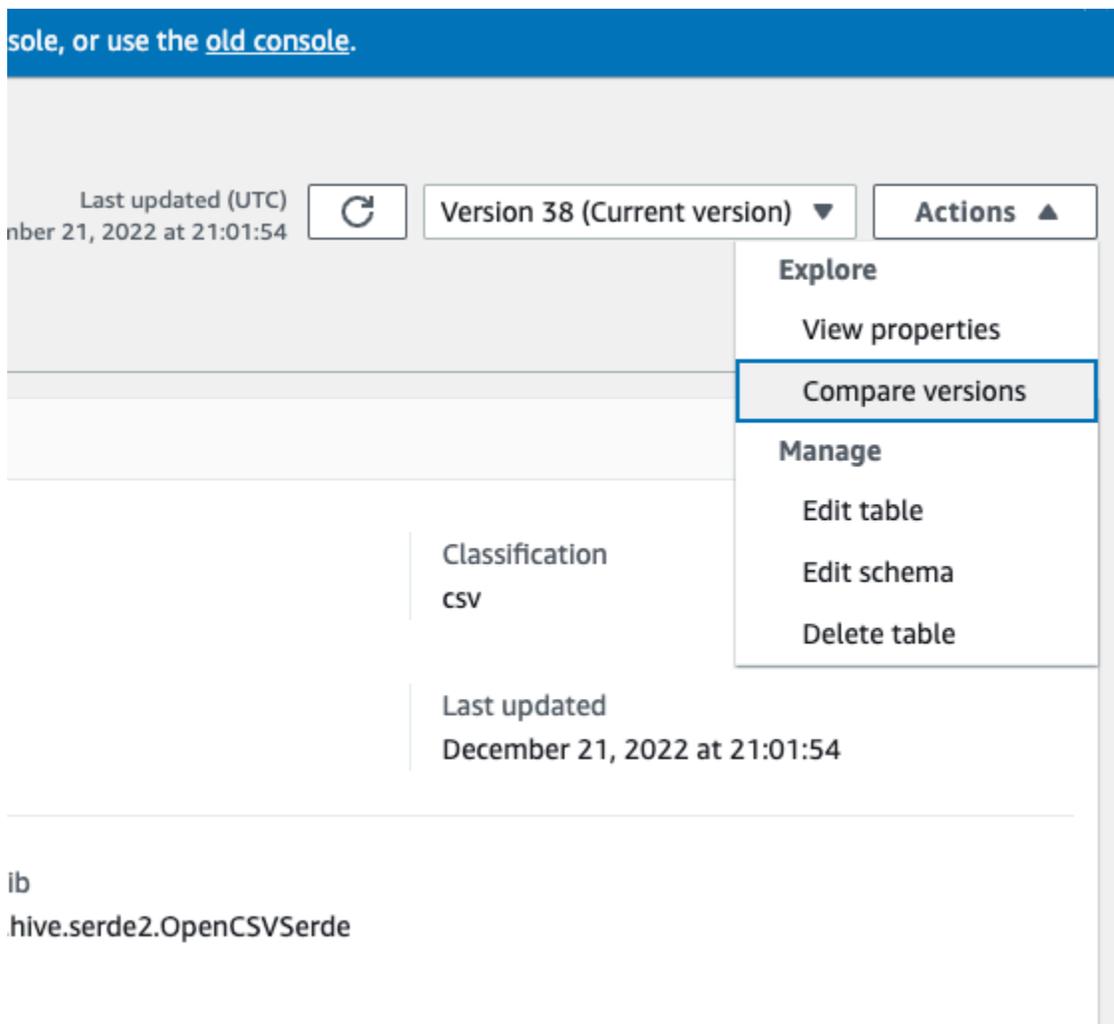
テーブルの詳細を表示するための詳細な手順については、コンソールの [Explore table] (テーブルの確認) チュートリアルを参照してください。

テーブルスキーマのバージョンの比較

テーブルスキーマの2つのバージョンを比較する場合、展開・折りたたみによってネストされた行の変更点を比較、2つのバージョンのスキーマを並べて比較、テーブルプロパティを並べて表示、などが行えます。

バージョンの比較方法

1. AWS Glue コンソールから [テーブル] > [アクション] の順に選択し、[バージョンの比較] を選択します。



2. バージョンのドロップダウンメニューから、比較するバージョンを選択します。スキーマを比較すると、[スキーマ] タブがオレンジ色で強調表示されます。
3. 2つのバージョン間でテーブルを比較すると、テーブルスキーマが画面の左側と右側に表示されます。列名、データ型、キー、コメントの各フィールドを並べて比較することで、変更点を視覚的に判断できます。変更があると、色付きのアイコンでその変更の種類が表示されます。

- 削除済み — 赤いアイコンによる表示は、以前のバージョンのテーブルスキーマから列が削除されたことを示します。
- 編集済みまたは移動済み — 青いアイコンによる表示は、新しいバージョンのテーブルスキーマで列が変更または移動されたことを示します。
- 追加 — 緑色のアイコンによる表示は、新しいバージョンのテーブルスキーマに列が追加されたことを示します。
- ネストされた変更 — 黄色のアイコンによる表示は、ネストされた列に変更が含まれていることを示します。列を選択して展開すると、削除、編集、移動、追加のいずれかが実行された列が表示されます。

Compare versions: cloudtrail_data

Legend: Deleted Edited/Moved Added Nested Changes Deleted

Version 0 (Last updated (UTC) January 17, 2023 at 19:08:58) | Version 2 (Current version) (Last updated (UTC) January 17, 2023 at 19:16:04)

Schema Properties

Table fields (33)

Field name	Data type	Key	Comment
eventversion	string	-	-
useridentity	struct	-	-
eventtime	string	-	-
eventsource	string	-	-
eventname	string	-	-
awsregion	string	-	-
sourceipaddress	string	-	-
useragent	string	-	-
requestparameters	struct	-	-
bucketName	string	-	-
Host	string	-	-
acl	string	-	-
lookupAttributes	array	-	-
startTime	string	-	-
endTime	string	-	-
maxResults	int	-	-
nextToken	string	-	-
filter	struct	-	-
aggregateField	string	-	-
responseelements	string	-	-
additionaleventdata	struct	-	-
requestid	string	-	-
eventid	string	-	-
readonly	boolean	-	-
resources	array	-	-
eventtype	string	-	-
managementevent	boolean	-	-
recipientaccountid	string	-	-
shareeventid	string	-	-
eventcategory	string	-	-
sessioncredentialfromconsole	string	-	-
errorcode	string	-	-
errormessage	string	-	-
new_col	string	-	-
eventid	string	(0)	-

4. 検索バーのフィルタフィールドを使用すると、ここに入力した文字に基づいてフィールドを表示できます。いずれかのテーブルバージョンで列名を入力すると、フィルタリングされたフィールドが両方のテーブルバージョンに表示され、変更が行われた箇所がわかります。
5. プロパティを比較するには、[プロパティ] タブを選択します。
6. バージョンの比較を停止するには、[比較の停止] を選択してテーブルのリストに戻ります。

AWS Glue でのパーティションインデックスの使用

時間の経過とともに、数十万のパーティションがテーブルに追加されます。[GetPartitions API](#) は、テーブル内のパーティションをフェッチするために使用されます。この API は、リクエストで指定された式と一致するパーティションを返します。

「国」、「カテゴリ」、「年」、「月」、「creationDate」のキーで分割された例として、sales_data テーブルを参考にします。2020 年の 2020 年 8 月 15 日以降に「書籍」のカテゴリで販売されたアイテムの売り上げデータを取得する場合、Data Catalog に「Category = 'Books' and creationDate > '2020-08-15'」という式で GetPartitions リクエストを行います。

テーブルにパーティションインデックスが存在しない場合、AWS Glue では、テーブルのすべてのパーティションがロードされ、GetPartitions リクエストでユーザーが指定したクエリ式を使用してロードされたパーティションがフィルタリングされます。インデックスのないテーブルでパーティション数が増えると、クエリの実行に時間がかかります。インデックスを使用すると、GetPartitions クエリでは、テーブル内のすべてのパーティションをロードするのではなく、パーティションのサブセットを取得しようとしています。

トピック

- [パーティションインデックスについて](#)
- [パーティションインデックスを持つテーブルの作成](#)
- [パーティションインデックスの既存テーブルへの追加](#)
- [テーブル上のパーティションインデックスの説明](#)
- [パーティションインデックスの使用に関する制限事項](#)
- [GetPartitions 呼び出しを最適化するためのインデックスの使用](#)
- [エンジンとの統合](#)

パーティションインデックスについて

パーティションインデックスを作成する際、指定したテーブルに既に存在するパーティションキーのリストを指定します。パーティションインデックスは、テーブルで定義されているパーティションキーのサブリストです。パーティションインデックスは、テーブルで定義されたパーティションキーを任意の順序にして作成できます。上記の sales_data テーブルでは、可能なインデックスには (国、カテゴリ、creationDate)、(国、カテゴリ、年)、(国、カテゴリ)、(国)、(カテゴリ、国、年、月) などがあります。

Data Catalog は、インデックスの作成時に指定された順序でパーティション値を連結します。インデックスは、テーブルに追加されたパーティションと整合して構築されます。インデックスは文字列 (string、char、varchar)、数値 (int、bigint、long、tinyint、smallint)、日付 (yyyy-MM-dd) の列タイプに対して作成できます。

サポートされているデータ型

- Date – YYYY-MM-DD などの ISO 形式の日付。たとえば、2020-08-15 の日付。この形式では、ハイフン (-) を使用して年、月、日を区切ります。インデックス用の日付の許容範囲は、0000-01-01 から 9999-12-31 までです。
- String – 一重引用符または二重引用符で囲まれた文字列リテラルです。
- Char – char(10) など、1 から 255 までの指定された長さを持つ固定長の文字データです。
- Varchar – varchar(10) など、1 から 65535 までの指定された長さを持つ可変長の文字データです。
- Numeric — int、bigint、long、tinyint、smallint

「Numeric」、「String」、「Date」データ型のインデックスは、=、>、>=、<、<=、演算子間をサポートします。インデックス作成ソリューションは現在、AND 論理演算子のみをサポートしています。「LIKE」、「IN」、「OR」、および「NOT」の演算子のある部分式は、インデックスを使用してフィルタリングする場合、この式では無視されます。無視された部分式のフィルタリングは、インデックスフィルタリングを適用した後にフェッチされたパーティションに対して行われます。

テーブルに追加されたパーティションごとに、対応するインデックス項目が作成されます。「n」個のパーティションを持つテーブルの場合、1つのパーティションインデックスは「n」個のパーティションインデックス項目になります。同じテーブルの「m」個のパーティションインデックスは、「m*n」個のパーティションインデックス項目になります。各パーティションインデックス項目は、現在のデータカタログストレージの AWS Glue 料金ポリシーに従って課金されます。ストレージオブジェクトの料金の詳細については、「[AWS Glue の料金](#)」を参照してください。

パーティションインデックスを持つテーブルの作成

テーブルの作成中にパーティションインデックスを作成できます。CreateTable リクエストは、[PartitionIndex オブジェクト](#)のリストを入力として指定します。1つのテーブルには、最大3つのパーティションインデックスを作成できます。各パーティションインデックスには、テーブルに対して定義された名前と partitionKeys のリストが必要です。テーブル上に作成されたインデックスは、[GetPartitionIndexes API](#) を使用してフェッチできます。

パーティションインデックスの既存テーブルへの追加

既存のテーブルにパーティションインデックスを追加するには、CreatePartitionIndex オペレーションを使用します。CreatePartitionIndex オペレーションごとに、1つの PartitionIndex を作成できます。インデックスを追加しても、テーブルの可用性には影響しません。これは、インデックスの作成中もテーブルが使用可能になるためです。

追加されたパーティションのインデックスステータスが CREATING に設定され、インデックスデータの作成が開始されます。インデックスの作成プロセスが正常に終了すると、indexStatus は ACTIVE に更新され、プロセスが正常に終了しなかった場合、インデックスステータスは FAILED に更新されます。インデックスが作成できない理由は、複数あり得ます。GetPartitionIndexes オペレーションを使用して、障害の詳細を取得できます。考えられる障害は次のとおりです。

- ENCRYPTED_PARTITION_ERROR – 暗号化されたパーティションを持つテーブルでのインデックスの作成はサポートされていません。
- INVALID_PARTITION_TYPE_DATA_ERROR – partitionKey の値が対応する partitionKey データ型に対して有効な値でない場合に発生します。例えば、「int」データ型の partitionKey の値が「foo」の場合などです。
- MISSING_PARTITION_VALUE_ERROR – indexedKey の partitionValue がない場合に発生します。これは、テーブルのパーティション化に整合性がない場合に発生します。
- UNSUPPORTED_PARTITION_CHARACTER_ERROR – インデックス付きパーティションキーの値に「\u0000」、「\u0001」、または「\u0002」という文字が含まれている場合に発生します。
- INTERNAL_ERROR – インデックスの作成中に内部エラーが発生しました。

テーブル上のパーティションインデックスの説明

テーブル上に作成されたパーティションインデックスを取得するには、GetPartitionIndexes オペレーションを使用します。レスポンスとして、テーブル上のすべてのインデックスと、各インデックスの現在のステータス (IndexStatus) が返ります。

パーティションインデックスの `IndexStatus` は、次のいずれかになります。

- `CREATING` – インデックスは現在作成されていますが、まだ使用することはできません。
- `ACTIVE` – インデックスは、すぐに使用できます。リクエストにインデックスを使用して最適化されたクエリを実行できます。
- `DELETING` – インデックスは現在削除されているため、使用できなくなっています。アクティブ状態のインデックスは、`DeletePartitionIndex` リクエストを使用してステータスを `ACTIVE` から `DELETING` に移行することによって削除できます。
- `FAILED` – 既存のテーブルでインデックスの作成に失敗しました。失敗したインデックスは、直近の 10 個まで各テーブルに保存されます。

既存のテーブルで作成されたインデックスには、次のような状態遷移があり得ます。

- `CREATING` → `ACTIVE` → `DELETING`
- `CREATING` → `FAILED`

パーティションインデックスの使用に関する制限事項

パーティションインデックスを作成したら、テーブルとパーティションの機能に対する次の変更点に注意してください。

新規パーティションの作成 (インデックス追加後)

テーブルにパーティションインデックスが作成されると、テーブルに追加されたすべての新しいパーティションは、インデックス付きキーのデータ型チェックのために検証されます。インデックス付きキーのパーティション値は、データ型形式について検証されます。データ型のチェックが失敗すると、パーティションの作成操作は失敗します。sales_data テーブルに対して、キー (カテゴリ、年) のインデックスが作成され、カテゴリの型が `string`、年の型が `int` の場合、`YEAR` の値が「foo」であれば、新しいパーティションの作成は失敗します。

インデックスを有効にすると、`U+0000`、`U+00001`、`U+0002` という文字を含むインデックス付きキー値のパーティションの追加が失敗するようになります。

テーブルの更新

テーブル上にパーティションインデックスを作成すると、既存のパーティションキーのパーティションキー名を変更することはできません。また、インデックスに登録されているキーのタイプまたは順序を変更することはできません。

GetPartitions 呼び出しを最適化するためのインデックスの使用

インデックス付きテーブルで GetPartitions を呼び出すとき、式を含めることができます。また、適用可能な場合は、Data Catalogで、可能であればインデックスが使用されます。インデックスの最初のキーは、フィルタリングに使用されるインデックスの式に含めて渡す必要があります。フィルタリングのインデックスの最適化は、ベストエフォートとして適用されます。Data Catalogでは、可能な限りインデックスの最適化を使用しようとはしますが、インデックスが見つからない場合、またはサポートされていない演算子の場合は、すべてのパーティションをロードする既存の実装に戻します。

上記の sales_data テーブルには、[Country, Category, Year] (国、カテゴリ、年) のインデックスを追加できます。式に「Country」が渡されない場合、登録されたインデックスは、インデックスを使用したパーティションのフィルタリングができません。最大3つのインデックスを追加して、さまざまなクエリパターンをサポートできます。

いくつかの式を例に取り、そこでインデックスが機能する様子を示します。

表現	インデックスの使用方法
Country = 'US'	インデックスは、パーティションをフィルタリングするために使用されます。
Country = 'US' and Category = 'Shoes'	インデックスは、パーティションをフィルタリングするために使用されます。
Category = 'Shoes'	式に「country」が指定されていないため、インデックスは使用されません。レスポンスを返すため、すべてのパーティションがロードされます。
Country = 'US' and Category = 'Shoes' and Year > '2018'	インデックスは、パーティションをフィルタリングするために使用されます。
Country = 'US' and Category = 'Shoes' and Year > '2018' and month = 2	インデックスを使用して、「country = "US" and category = "shoes" and year > 2018」に該当するすべてのパーティションがフェッチされます。その後、月に関する式によるフィルタリングが実行されます。

表現	インデックスの使用方法
Country = 'US' AND Category = 'Shoes' OR Year > '2018'	OR 演算子が式に存在するため、インデックスは使用されません。
Country = 'US' AND Category = 'Shoes' AND (Year = 2017 OR Year = '2018')	インデックスを使用して、「country = "US" and category = "shoes"」に該当するのすべてのパーティションがフェッチされ、その後、年に関する式によるフィルタリングが実行されます。
Country in ('US', 'UK') AND Category = 'Shoes'	IN 演算子は現在サポートされていないため、インデックスはフィルタリングに使用されません。
Country = 'US' AND Category in ('Shoes', 'Books')	インデックスを使用して、「country = "US"」に該当するのすべてのパーティションが取得され、その後、カテゴリに関する式によるフィルタリングが実行されます。
Country = 'US' AND Category in ('Shoes', 'Books') AND (creationDate > '2023-9-01')	「creationDate > '2023-9-01」に設定し、インデックスを使用して「country = "US"」に該当するすべてのパーティションをフェッチしたら、カテゴリの式にフィルタリングが実行されます。

エンジンとの統合

Redshift Spectrum、Amazon EMR、および AWS Glue ETL Spark DataFrames は、インデックスが AWS Glue で ACTIVE 状態になった後にパーティションを取得するためにインデックスを利用できます。[Athena](#) と [AWS Glue ETL DynamicFrame](#) では、クエリを向上させるためにインデックスを利用するには追加のステップに従う必要があります。

パーティションのフィルタリングの有効化

Athena でパーティションフィルタリングを有効にするには、テーブルのプロパティを次のように更新する必要があります。

1. AWS Glue コンソールで、[データカタログ] の [テーブル] を選択します。

2. テーブルを選択します。
3. [アクション] で [テーブルの編集] を選択します。
4. [テーブルのプロパティ] で、次の内容を追加します。
 - キー - partition_filtering.enabled
 - 値 - true
5. [Apply] を選択します。

または、Athena で [ALTER TABLE SET PROPERTIES](#) クエリを実行してこのパラメータを設定することもできます。

```
ALTER TABLE partition_index.table_with_index
SET TBLPROPERTIES ('partition_filtering.enabled' = 'true')
```

他の AWS のサービスとの統合

AWS Glue クローラー を使用して AWS Glue Data Catalog にデータを入力することもできますが、いくつかの AWS サービスではカタログに自動的に統合し、データを入力することができます。以下のセクションでは、データカタログに入力できる AWS サービスでサポートされる特定のユースケースについて詳しく説明します。

トピック

- [AWS Lake Formation](#)
- [Amazon Athena](#)

AWS Lake Formation

AWS Lake Formation は、AWS で簡単にセキュアなデータレイクを構築できるサービスです。Lake Formation は AWS Glue 上に構築され、Lake Formation と AWS Glue は同じ AWS Glue Data Catalog を共有します。Amazon S3 データロケーションを Lake Formation に登録し、Lake Formation コンソールを使用して、AWS Glue データカタログにデータベースとテーブルを作成し、データアクセスポリシーを定義し、データレイク全体のデータアクセスを一元的に監査できます。Lake Formation の細粒度のアクセスコントロールを使用して、既存のデータカタログリソースと Amazon S3 データロケーションを管理できます。

Lake Formation に登録されたデータを使用すると、IAM プリンシパル、AWS アカウント、AWS 組織、組織単位間でデータカタログリソースを安全に共有できます。

Lake Formation を使用したデータカタログリソースの作成の詳細については、AWS Lake Formation デベロッパーガイドの「[データカタログのテーブルとデータベースの作成](#)」を参照してください。

Amazon Athena

Amazon Athena は、AWS アカウントの Amazon S3 データに関するテーブルメタデータの保存と取得にデータカタログを使用します。テーブルメタデータは、Athena クエリエンジンがクエリするデータの検索、読み込み、および処理方法を把握できるようにします。

Athena CREATE TABLE ステートメントを直接使用して、AWS Glue Data Catalog にデータを入力できます。クローラーを実行することなく、データカタログでスキーマとパーティションメタデータを手動で定義してデータを入力できます。

1. Athena コンソールで、テーブルメタデータをデータカタログに保存するデータベースを作成します。
2. CREATE EXTERNAL TABLE ステートメントを使用して、データソースのスキーマを定義します。
3. データがパーティション化されている場合は、PARTITIONED BY 句を使用してパーティションキーを定義します。
4. LOCATION 句を使用して、実際のデータファイルを保存する Amazon S3 パスを指定します。
5. CREATE TABLE ステートメントを実行します。

このクエリは、実際にデータをクローリングすることなく、定義したスキーマとパーティションに基づいてデータカタログにテーブルメタデータを作成します。

Athena でテーブルにクエリを実行すると、データカタログのメタデータを使用して Amazon S3 のデータファイルにアクセスしてクエリを実行できます。

詳細については、Amazon Athena ユーザーガイドの「[データベースとテーブルの作成](#)」を参照してください。

データカタログの設定

データカタログの設定には、アカウントのデータカタログに対する暗号化およびアクセス許可のオプションが含まれています。

Data catalog settings

Last updated (UTC)
January 1, 1970 at 00:00:00



Choose encryption and permission options for your accounts data catalog.

Encryption options

Metadata encryption

Enable at-rest encryption for metadata stored in the data catalog.

Encrypt connection passwords

When enabled, the password you provide when you create a connection is encrypted with the given KMS key.

Permissions

Add a policy to define fine-grained access control of the data catalog.

1	

JSON Ln 1, Col 1 Errors: 0 Warnings: 0

Cancel

Save

Data Catalog のきめ細かなアクセスコントロールを変更するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. 暗号化オプションを選択します。
 - メタデータの暗号化 – データカタログのメタデータを暗号化するには、このチェックボックスをオンにします。メタデータは、指定した AWS Key Management Service (AWS KMS) キーを使用して、保存時に暗号化されます。詳細については、[Data Catalog の暗号化](#) を参照してください。
 - 接続パスワードの暗号化 – 接続を作成または更新する際、AWS Glue 接続オブジェクトのパスワードを暗号化するには、このチェックボックスをオンにします。パスワードは、指定した AWS KMS キーを使用して暗号化されます。パスワードが返ったら、暗号化されます。このオプションは、Data Catalog のすべての AWS Glue 接続のグローバル設定です。このチェックボックスをオフにしても、以前暗号化したパスワードは、作成または更新する際に使用したキーを使用して暗号化された状態に維持されます。AWS Glue 接続の詳細については、「[データへの接続](#)」を参照してください。

このオプションを有効にする場合は、AWS KMS キーを選択するか、[キーの ARN を入力] を選択して、キーの Amazon リソースネーム (ARN) を指定します。
arn:aws:kms:*region*:*account-id*:key/*key-id* の形式で ARN を入力します。
arn:aws:kms:*region*:*account-id*:alias/*alias-name* など、キーのエイリアスで ARN を指定することもできます。

Important

このオプションが選択されている場合は、接続を作成または更新するユーザーまたはロールに、指定された KMS キーの kms:Encrypt アクセス許可が必要です。

詳細については、[接続パスワードの暗号化](#) を参照してください。

3. [Settings] (設定) をクリックして、[Permissions] (アクセス権限) エディタで、アカウントの Data Catalog のきめ細かなアクセスコントロールを変更するポリシーステートメントを追加します。一度に Data Catalog にアタッチできるポリシーは 1 つだけです。JSON リソースポリシーをこのコンソールに貼り付けることができます。詳細については、[AWS Glue 内のリソースベースのポリシー](#) を参照してください。
4. [Save] (保存) をクリックして、行った変更を含めて Data Catalog を更新します。

また、AWS Glue API オペレーションを使用して、リソースポリシーを格納、取得、削除できます。詳細については、[AWS Glue のセキュリティ API](#) を参照してください。

トランザクションテーブルへのデータ入力と管理

[Apache Iceberg](#)、[Apache Hudi](#)、および Linux Foundation [Delta Lake](#) は、Apache Spark で大規模なデータ分析とデータレイクワークロードを処理するために設計されたオープンソースのテーブル形式です。

次の方法を使用して、AWS Glue Data Catalog で Iceberg、Hudi、および Delta Lake の各テーブルにデータを入力できます。

- AWS Glue クローラー – AWS Glue クローラー は、Iceberg、Hudi、Delta Lake のテーブルメタデータを自動的に検出してデータカタログに入力できます。詳細については、「[クローラーを使用したデータカタログへの入力](#)」を参照してください。
- AWS Glue ETL ジョブ – ETL ジョブを作成して Iceberg、Hudi、Delta Lake の各テーブルにデータを書き込み、そのメタデータをデータカタログに入力できます。詳細については、「[AWS Glue ETL ジョブでのデータレイクフレームワークの使用](#)」を参照してください。
- AWS Glue コンソール、AWS Lake Formation コンソール、AWS CLI または API – AWS Glue コンソール、Lake Formation コンソール、または API を使用して、データカタログに Iceberg テーブル定義を作成し、管理できます。

トピック

- [Apache Iceberg テーブルの作成](#)
- [Iceberg テーブルの最適化](#)

Apache Iceberg テーブルの作成

Amazon S3 に格納されたデータにより、AWS Glue Data Catalog で Apache Parquet データ形式を使用する Apache Iceberg テーブルを作成できます。Data Catalog のテーブルは、データストア内のデータを表すメタデータ定義です。デフォルトでは、AWS Glue は Iceberg v2 テーブルを作成します。v1 テーブルと v2 テーブルの違いについては、Apache Iceberg ドキュメントの「[形式バージョンの変更](#)」を参照してください。

[Apache Iceberg](#) は、非常に大規模な分析データセット用のオープンテーブル形式です。Iceberg では、スキーマの変更 (スキーマ進化とも呼ばれます) を簡単に行うことができます。つまり、

基になるデータを中断することなく、データテーブルの列を追加、名前変更、または削除できます。Iceberg はデータのバージョニングもサポートしているため、データの変更を経時的に追跡できます。これにより、タイムトラベル機能が有効になるため、過去のバージョンのデータにアクセスしてクエリを実行し、更新と削除の間に行われたデータの変更を分析できます。

AWS Glue または Lake Formation コンソールあるいは AWS Glue API の CreateTable オペレーションを使用して、データカタログに Iceberg テーブルを作成することができます。詳細については、「[CreateTable アクション \(Python: create_table\)](#)」を参照してください。

Data Catalog に Iceberg テーブルを作成する場合、読み取りと書き込みを実行できるように、Amazon S3 でテーブル形式とメタデータファイルのパスを指定する必要があります。

Amazon S3 のデータロケーションを AWS Lake Formation に登録すると、Lake Formation を使用してきめ細かなアクセス制御の許可により Iceberg テーブルを保護できます。Amazon S3 のソースデータと、Lake Formation に登録されていないメタデータへのアクセスは、Amazon S3 および AWS Glue アクションの IAM アクセス許可ポリシーによって決定されます。詳細については、「[アクセス許可の管理](#)」を参照してください。

Note

Data Catalog は、パーティションの作成と Iceberg テーブルプロパティの追加をサポートしていません。

前提条件

Data Catalog に Iceberg テーブルを作成し、Lake Formation のデータアクセス許可を設定するには、次の要件を満たす必要があります。

1. Lake Formation にデータが登録されていない状態で Iceberg テーブルを作成するために必要なアクセス許可。

Data Catalog にテーブルを作成するために必要なアクセス許可に加えて、テーブル作成者には次のアクセス許可が必要です。

- リソース `arn:aws:s3:::{bucketName}` での `s3:PutObject`
- リソース `arn:aws:s3:::{bucketName}` での `s3:GetObject`
- リソース `arn:aws:s3:::{bucketName}` での `s3:DeleteObject`

2. Lake Formation にデータが登録されている状態で Iceberg テーブルを作成するために必要なアクセス許可:

Lake Formation を使用してデータレイク内のデータを管理および保護するには、テーブルのデータを含む Amazon S3 ロケーションを Lake Formation に登録します。これは、Lake Formation が Athena、Redshift Spectrum、Amazon EMR などの AWS 分析サービスに認証情報を提供して、データにアクセスできるようにするためです。Amazon S3 ロケーションの登録の詳細については、「[データレイクへの Amazon S3 ロケーションの追加](#)」を参照してください。

Lake Formation に登録されている、基盤となるデータを読み書きするプリンシパルには、次のアクセス許可が必要です。

- lakeformation:GetDataAccess
- DATA_LOCATION_ACCESS

ロケーションに対するデータロケーション許可を持つプリンシパルは、すべての子ロケーションに対するロケーション許可も持っています。

データロケーション許可の詳細については、「[基盤となるデータのアクセスコントロール](#)」を参照してください。

圧縮を有効にするには、Data Catalog 内のテーブルを更新するアクセス許可を持つ IAM ロールを、サービスが引き受ける必要があります。詳細については、「[テーブル最適化の前提条件](#)」を参照してください。

Iceberg テーブルの作成

このページに記載されているように、AWS Glue または Lake Formation コンソールあるいは AWS Command Line Interface を使用して Iceberg v1 および v2 テーブルを作成できます。AWS Glue クローラーを使用して Iceberg テーブルを作成することもできます。詳細については、「AWS Glue デベロッパーガイド」の「[Data Catalog とクローラー](#)」を参照してください。

Iceberg テーブルを作成するには

Console

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。

2. Data Catalog で [テーブル] を選択し、[テーブルの作成] ボタンを使用して次の属性を指定します。
 - テーブル名 - テーブルの名前を入力します。Athena を使用してテーブルにアクセスする場合は、「Amazon Athena ユーザーガイド」の[命名に関するヒント](#)を使用します。
 - データベース - 既存のデータベースを選択するか、新しいデータベースを作成します。
 - 説明 - テーブルの説明。テーブルの内容を理解しやすくするために説明を記入できます。
 - テーブル形式 - [テーブル形式] として、[Apache Iceberg] を選択します。
 - 圧縮を有効にする - [圧縮を有効にする] を選択すると、テーブル内の小さな Amazon S3 オブジェクトが圧縮されてより大きなオブジェクトにまとめられます。
 - IAM ロール:- 圧縮を実行する場合、サービスはユーザーに代わって IAM ロールを引き受けません。IAM ロールは、ドロップダウンを使用して選択できます。圧縮を有効にするために必要なアクセス許可がロールにあることを確認します。

必要なアクセス許可の詳細については、「[テーブル最適化の前提条件](#)」を参照してください。

- ロケーション - メタデータテーブルを保存する Amazon S3 内のフォルダへのパスを指定します。Iceberg が読み取りと書き込みを実行するには、メタデータファイルと Data Catalog 内のロケーションが必要です。
- スキーマ - [列の追加] を選択して、列と、列のデータ型を追加します。空のテーブルを作成して、後でスキーマを更新することもできます。Data Catalog は Hive データ型をサポートしています。詳細については、「[Hive データ型](#)」を参照してください。

Iceberg では、テーブルを作成した後でスキーマとパーティションを進化させることができます。[\[Athena クエリ\]](#) を使用してテーブルスキーマを更新し、[\[Spark クエリ\]](#) を使用してパーティションを更新できます。

AWS CLI

```
aws glue create-table \  
  --database-name iceberg-db \  
  --region us-west-2 \  
  --open-table-format-input '{  
    "IcebergInput": {  
      "MetadataOperation": "CREATE",  
      "Version": "2"  
    }  
  }' \  
  --table-name <table-name>
```

```
--table-input '{"Name":"test-iceberg-input-demo",
  "TableType": "EXTERNAL_TABLE",
  "StorageDescriptor":{
    "Columns":[
      {"Name":"col1", "Type":"int"},
      {"Name":"col2", "Type":"int"},
      {"Name":"col3", "Type":"string"}
    ],
    "Location":"s3://DOC_EXAMPLE_BUCKET_ICEBERG/"
  }
}'
```

Iceberg テーブルの最適化

Apache Iceberg などのオープンテーブル形式を使用する Amazon S3 データレイクは、データを Amazon S3 オブジェクトとして保存します。データレイクテーブルに数千の小さな Amazon S3 オブジェクトがある場合、Iceberg テーブルのメタデータのオーバーヘッドが増加し、読み取りパフォーマンスに悪影響が及びます。AWS 分析サービス (Amazon Athena、Amazon EMR、AWS Glue ETL ジョブなど) による読み取りパフォーマンスを向上させるために、AWS Glue Data Catalog は、Data Catalog 内の Iceberg テーブル用にマネージド圧縮 (小さな Amazon S3 オブジェクトを圧縮してより大きなオブジェクトにまとめるプロセス) を提供しています。Lake Formation コンソール、AWS Glue コンソール、AWS CLI、または AWS API を使用して、Data Catalog 内の個々の Iceberg テーブルの圧縮を有効または無効にすることができます。

テーブル最適化は、テーブルパーティションを継続的にモニタリングして、ファイル数とファイルサイズがしきい値を超えたときに圧縮プロセスを開始します。write.target-file-size-bytes プロパティで指定されたファイルサイズが 128MB から 512MB の範囲内にある場合、Iceberg テーブルは圧縮の対象となります。データカタログでは、テーブルに write.target-file-size-bytes プロパティの 75% 未満のファイルが 5 つ以上ある場合、圧縮プロセスが開始されます。

たとえば、ファイルサイズのしきい値を write.target-file-size-bytes プロパティで 512MB に設定し (128MB から 512MB の指定された範囲内)、テーブルに 10 個のファイルが含まれているとします。10 個のファイルのうち 6 個がそれぞれ 384MB (0.75*512) 未満の場合、データカタログは圧縮をトリガーします。

データカタログは、同時クエリに支障を来たすことなく圧縮を実行します。データカタログは、Parquet 形式のテーブルのデータ圧縮のみをサポートします。

サポートされているデータ型、圧縮形式、制限事項については、「[マネージドデータ圧縮でサポートされる形式と制限事項](#)」を参照してください。

トピック

- [テーブル最適化の前提条件](#)
- [圧縮を有効にする](#)
- [圧縮を無効にする](#)
- [圧縮の詳細の表示](#)
- [Amazon CloudWatch メトリクスの表示](#)
- [オプティマイザーの削除](#)
- [マネージドデータ圧縮でサポートされる形式と制限事項](#)

テーブル最適化の前提条件

テーブルオプティマイザーは、テーブルの圧縮を有効にする際に指定する AWS Identity and Access Management (IAM) ロールの許可を引き受けます。IAM ロールには、データカタログ内のデータを読み取ったり、メタデータを更新したりするための許可が必要です。IAM ロールを作成し、次のインラインポリシーをアタッチできます。

- Lake Formation に登録されていないデータの場所に対する Amazon S3 の読み取り/書き込み許可を付与する次のインラインポリシーを追加します。このポリシーには、データカタログ内のテーブルを更新するための許可、および AWS Glue が Amazon CloudWatch ログにログを追加したり、メトリクスを公開したりするのを許可するための許可も含まれています。Lake Formation に登録されていない Amazon S3 のソースデータへのアクセスは、Amazon S3 および AWS Glue アクションの IAM アクセス許可ポリシーによって決定されます。

次のインラインポリシーでは、`bucket-name` を Amazon S3 バケット名、`aws-account-id` および `region` をデータカタログの有効な AWS アカウント番号およびリージョン、`database_name` をデータベース名、`table_name` をテーブル名に置き換えます。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "s3:PutObject",
```

```

        "s3:GetObject",
        "s3:DeleteObject"
    ],
    "Resource": [
        "arn:aws:s3:::<bucket-name>/*"
    ]
},
{
    "Effect": "Allow",
    "Action": [
        "s3:ListBucket"
    ],
    "Resource": [
        "arn:aws:s3:::<bucket-name>"
    ]
},
{
    "Effect": "Allow",
    "Action": [
        "glue:UpdateTable",
        "glue:GetTable"
    ],
    "Resource": [
        "arn:aws:glue:<region>:<aws-account-id>:table/<database-name>/<table-
name>",
        "arn:aws:glue:<region>:<aws-account-id>:database/<database-name>",
        "arn:aws:glue:<region>:<aws-account-id>:catalog"
    ]
},
{
    "Effect": "Allow",
    "Action": [
        "logs:CreateLogGroup",
        "logs:CreateLogStream",
        "logs:PutLogEvents"
    ],
    "Resource": "arn:aws:logs:<region>:<aws-account-id>:log-group:/aws-glue/
iceberg-compaction/logs:*"
}
]
}

```

- 次のポリシーを使用して、Lake Formation に登録されたデータの圧縮を有効にします。

テーブルに対する IAM_ALLOWED_PRINCIPALS グループアクセス許可が圧縮ロールに付与されていない場合、ロールにはテーブルに対する Lake Formation の ALTER、DESCRIBE、INSERT、および DELETE アクセス許可が必要です。

Lake Formation に Amazon S3 バケットを登録する方法の詳細については、「[データレイクへの Amazon S3 ロケーションの追加](#)」を参照してください。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "lakeformation:GetDataAccess"
      ],
      "Resource": "*"
    },
    {
      "Effect": "Allow",
      "Action": [
        "glue:UpdateTable",
        "glue:GetTable"
      ],
      "Resource": [
        "arn:aws:glue:<region>:<aws-account-id>:table/<databaseName>/<tableName>",
        "arn:aws:glue:<region>:<aws-account-id>:database/<database-name>",
        "arn:aws:glue:<region>:<aws-account-id>:catalog"
      ]
    },
    {
      "Effect": "Allow",
      "Action": [
        "logs:CreateLogGroup",
        "logs:CreateLogStream",
        "logs:PutLogEvents"
      ],
      "Resource": "arn:aws:logs:<region>:<aws-account-id>:log-group:/aws-glue/iceberg-compaction/logs:*"
    }
  ]
}
```

```
}

```

- (オプション) [サーバー側の暗号化](#)を使用して暗号化された Amazon S3 バケットにデータがある Iceberg テーブルを圧縮するには、圧縮ロールに Amazon S3 オブジェクトを復号して、暗号化されたバケットにオブジェクトを書き込むための新しいデータキーを生成するアクセス許可が必要です。次のポリシーを目的の AWS KMS キーに追加します。バケットレベルの暗号化のみがサポートされています。

```
{
  "Effect": "Allow",
  "Principal": {
    "AWS": "arn:aws:iam::<aws-account-id>:role/<compaction-role-name>"
  },
  "Action": [
    "kms:Decrypt",
    "kms:GenerateDataKey"
  ],
  "Resource": "*"
}
```

- (オプション) Lake Formation で登録されたデータの場所については、場所を登録するために使用されるロールには、Amazon S3 オブジェクトを復号したり、暗号化されたバケットにオブジェクトを書き込むための新しいデータキーを生成したりする許可が必要です。詳細については、「[暗号化された Amazon S3 の場所の登録](#)」を参照してください。
- (オプション) AWS KMS キーが別の AWS アカウントに保存されている場合は、圧縮ロールに次のアクセス許可を含める必要があります。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "kms:Decrypt",
        "kms:GenerateDataKey"
      ],
      "Resource": ["arn:aws:kms:<REGION>:<KEY_OWNER_ACCOUNT_ID>:key/<KEY_ID>"]
    }
  ]
}
```

- 圧縮を実行するために使用するロールには、そのロールに対する iam:PassRole アクセス許可が必要です。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "iam:PassRole"
      ],
      "Resource": [
        "arn:aws:iam::<account-id>:role/<compaction-role-name>"
      ]
    }
  ]
}
```

- 圧縮プロセスを実行する IAM ロールを引き受けるには、次の信頼ポリシーを AWS Glue サービスのロールに追加します。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "",
      "Effect": "Allow",
      "Principal": {
        "Service": "glue.amazonaws.com"
      },
      "Action": "sts:AssumeRole"
    }
  ]
}
```

圧縮を有効にする

Lake Formation コンソール、AWS Glue コンソール、AWS CLI、または AWS API を使用して、Data Catalog 内の Apache Iceberg テーブルの圧縮を有効にすることができます。新しいテーブルの場合は、テーブル形式として Apache Iceberg を選択し、テーブルの作成時に圧縮を有効にすることができます。圧縮は、新しいテーブルのためにデフォルトで無効になっています。

Console

圧縮を有効にするには

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソールを開き、データレイク管理者、テーブル作成者、またはテーブルに対する および 許可を付与されたユーザーとしてサインインします。
2. ナビゲーションペインの [データカタログ] で [テーブル] を選択します。
3. [テーブル] ページで、圧縮を有効にするオープンテーブル形式のテーブルを選択し、[アクション] メニューで [圧縮を有効にする] を選択します。
4. テーブルを選択して [テーブルの詳細] ページを開いて、圧縮を有効にすることもできます。ページの下部セクションにある [テーブル最適化] タブを選択し、[圧縮を有効にする] を選択します。
5. 次に、[テーブル最適化の前提条件](#) セクションに示されている許可を持つ既存の IAM ロールをドロップダウンから選択します。

[新しい IAM ロールを作成] オプションを選択すると、サービスは圧縮を実行するために必要な許可を持つカスタムロールを作成します。

既存の IAM ロールを更新するには、以下のステップに従います。

- a. IAM ロールの許可ポリシーを更新するには、IAM コンソールで、圧縮の実行に使用されている IAM ロールにアクセスします。
- b. [Add permissions] (アクセス許可の追加) セクションで、[Create policy] (ポリシーの作成) を選択します。新しく開いたブラウザウィンドウで、ロールで使用する新しいポリシーを作成します。
- c. [Create policy] (ポリシーの作成) ページで、[JSON] タブを選択します。前提条件に示している JSON コードをポリシーエディタフィールドにコピーします。

AWS CLI

次の例は、圧縮を有効にする方法を示しています。アカウント ID を有効な AWS アカウント ID に置き換えます。データベース名とテーブル名を実際の Iceberg テーブル名とデータベース名に置き換えます。roleArn を、圧縮の実行に必要なアクセス許可を持つ IAM ロールの ARN (AWS リソースネーム) と IAM ロールの名前に置き換えます。

```
aws glue create-table-optimizer \  
  --catalog-id 123456789012 \  
  --database-name iceberg_db \  
  --table-name iceberg_table \  
  --table-optimizer-configuration  
  '{"roleArn":"arn:aws:iam::123456789012:role/compaction_role", "enabled":'true'}' \  
  --type compaction
```

AWS API

テーブルの圧縮を有効にするには、CreateTableOptimizer オペレーションを呼び出します。

圧縮を有効にすると、[テーブル最適化] タブに以下の圧縮の詳細が表示されます (約 15~20 分後)。

開始時刻

Lake Formation 内で圧縮処理が開始した時刻。値は UTC 時間のタイムスタンプです。

終了時刻

Data Catalog で圧縮処理が終了した時刻。値は UTC 時間のタイムスタンプです。

ステータス

圧縮実行のステータス。値は成功または失敗です。

圧縮ファイル数

圧縮したファイルの総数。

圧縮バイト数

圧縮したバイトの総数。

圧縮を無効にする

AWS Glue コンソールまたは AWS CLI を使用して、特定の Apache Iceberg テーブルの自動圧縮を無効にできます。

Console

1. [データカタログ]、[テーブル] の順に選択します。テーブルリストから、圧縮を無効にするオープンテーブル形式のテーブルを選択します。
2. Iceberg テーブルを選択し、[アクション] で [圧縮を無効にする] を選択できます。

[テーブルの詳細] ページの下部のセクションで [圧縮を無効にする] を選択して、テーブルの圧縮を無効にすることもできます。

3. 確認メッセージで [圧縮を無効にする] を選択します。圧縮は後で再度有効にすることができません。

確認すると、圧縮が無効になり、テーブルの圧縮ステータスが Off に戻ります。

AWS CLI

次の例で、アカウント ID を、有効な AWS アカウント ID に置き換えます。データベース名とテーブル名を、実際の Iceberg テーブル名とデータベース名に置き換えます。roleArn を、圧縮の実行に必要なアクセス許可を持つ IAM ロールの AWS リソースネーム (ARN) と IAM ロールの実際の名前に置き換えます。

```
aws glue update-table-optimizer \  
  --catalog-id 123456789012 \  
  --database-name iceberg_db \  
  --table-name iceberg_table \  
  --table-optimizer-configuration  
'{"roleArn":"arn:aws:iam::123456789012:role/compaction_role", "enabled":'false'}'\  
  --type compaction
```

AWS API

UpdateTableOptimizer オペレーションを呼び出して、特定のテーブルの圧縮を無効にします。

圧縮の詳細の表示

Apache Iceberg の圧縮ステータスは、Lake Formation コンソール、AWS CLI、または AWS API オペレーションを使用して表示できます。

Console

Iceberg テーブルの圧縮ステータスを表示するには (コンソール)

- Data Catalog の下部で [テーブル] を選択すると、Lake Formation コンソールで Iceberg テーブルの圧縮ステータスを表示できます。[圧縮ステータス] フィールドには、圧縮実行のステータスが表示されます。テーブルの設定を使用して、テーブルの形式と圧縮ステータスを表示できます。
- 特定のテーブルの圧縮実行履歴を表示するには、[AWS Glue Data Catalog] の [テーブル] を選択でテーブルを選択し、テーブルの詳細を表示します。[テーブル最適化] タブには、テーブルの圧縮履歴が表示されます。

AWS CLI

AWS CLI を使用して圧縮の詳細を表示できます。

次の例では、アカウントID を有効な AWS アカウント ID、データベース名、テーブル名を実際の Iceberg テーブル名に置き換えます。

- テーブルの前の圧縮実行の詳細を取得するには

```
aws get-table-optimizer \  
  --catalog-id 123456789012 \  
  --database-name iceberg_db \  
  --table-name iceberg_table \  
  --type compaction
```

- 次の例を使用して、特定のテーブルのオプティマイザーの履歴を取得します。

```
aws list-table-optimizer-runs \  
  --catalog-id 123456789012 \  
  --database-name iceberg_db \  
  --table-name iceberg_table \  
  --type compaction
```

- 次の例は、複数のオプティマイザーの圧縮実行と設定の詳細を取得する方法を示しています。最大 20 個のオプティマイザを指定できます。

```
aws glue batch-get-table-optimizer \  
--entries '[{"catalogId":"123456789012", "databaseName":"iceberg_db",  
"tableName":"iceberg_table", "type":"compaction"}]'
```

AWS API

- GetTableOptimizer オペレーションを使用して、前回実行したオプティマイザーの詳細を取得します。
- 特定のテーブル上の特定のオプティマイザーの履歴を取得するには、ListTableOptimizerRuns オペレーションを使用します。1 回の API 呼び出しで 20 個のオプティマイザーを指定できます。
- アカウント内の複数のオプティマイザーの設定の詳細を取得するには、BatchGetTableOptimizer オペレーションを使用します。このオペレーションではアカウント間呼び出しをサポートしていません。

Amazon CloudWatch メトリクスの表示

圧縮が正常に実行されると、サービスは圧縮ジョブのパフォーマンスに関する Amazon CloudWatch メトリクスを作成します。[CloudWatch メトリクス] に移動し、[メトリクス]、[すべてのメトリクス] の順に選択できます。メトリクスは、特定の名前空間 (AWS Glue など)、テーブル名、またはデータベース名でフィルタリングできます。

詳細については、「Amazon CloudWatch ユーザーガイド」の「[使用可能なメトリクスを表示する](#)」を参照してください。

- 圧縮されたバイト数
- 圧縮ファイル数
- ジョブに割り当てられた DPU 数
- ジョブの期間 (時間数)

オプティマイザーの削除

AWS CLI または AWS API オペレーションを使用して、テーブルのオプティマイザーと関連するメタデータを削除できます。

テーブルの圧縮履歴を削除するには、次の AWS CLI コマンドを実行します。

```
aws glue delete-table-optimizer \  
  --catalog-id 123456789012 \  
  --database-name iceberg_db \  
  --table-name iceberg_table \  
  --type compaction
```

テーブルのオプティマイザーを削除するには、DeleteTableOptimizer オペレーションを使用します。

マネージドデータ圧縮でサポートされる形式と制限事項

Amazon Athena、Amazon EMR、AWS Glue ETL ジョブなどの AWS 分析サービスによる読み取りパフォーマンスを向上させるために、は Data Catalog の Iceberg テーブルに対してマネージド圧縮 (小さな Amazon S3 オブジェクトを大きなオブジェクトに圧縮するプロセス) AWS Glue Data Catalog を提供します。

データ圧縮は、暗号化されたテーブルからのデータの読み取りなど、データの読み書きのためのさまざまなデータ型と圧縮形式をサポートしています。

データ圧縮は次をサポートします。

- ファイルタイプ - Parquet
- データ型 - ブール、整数、長整数、浮動小数点、倍精度浮動小数点数、文字列、10 進数、日付、時刻、タイムスタンプ、文字列、UUID、バイナリ
- 圧縮 - zstd、gzip、snappy、非圧縮
- 暗号化 - データ圧縮では、デフォルトの Amazon S3 暗号化 (SSE-S3) とサーバー側 KMS 暗号化 (SSE-KMS) のみがサポートされます。
- ビンパック圧縮
- スキーマ進化
- ターゲットファイルサイズ (iceberg 設定では write.target-file-size-bytes property) が 128MB ~ 512 MB の範囲内のテーブル。

- リージョン
 - アジアパシフィック (東京)
 - アジアパシフィック (ソウル)
 - アジアパシフィック (ムンバイ)
 - アジアパシフィック (シンガポール)
 - 欧州 (アイルランド)
 - 欧州 (ロンドン)
 - 欧州 (フランクフルト)
 - 米国東部 (バージニア北部)
 - 米国東部 (オハイオ)
 - 米国西部 (北カリフォルニア)
 - 南米 (サンパウロ)
- 基礎となるデータを保存する Amazon S3 バケットが別のアカウントにある場合、データカタログが存在するアカウントから圧縮を実行できます。これを実行するには、圧縮ルールが Amazon S3 バケットにアクセスできる必要があります。

データ圧縮は現在、次をサポートしていません。

- ファイルタイプ - Avro、ORC
- データ型 - 固定小数点
- 圧縮 - brotli、lz4
- パーティションの仕様が進化する中でのファイルの圧縮。
- 通常の並べ替えまたは Z オーダーの並べ替え
- ファイルのマージまたは削除 - 圧縮プロセスでは、削除ファイルが関連付けられているデータファイルはスキップされます。
- クロスアカウントテーブルでの圧縮 - クロスアカウントテーブルでは圧縮を実行できません。
- クロスリージョンテーブルでの圧縮 - クロスリージョンテーブルでは圧縮を実行できません。
- リソースのリンクでの圧縮の有効化
- Amazon S3 バケットの VPC エンドポイント
- [DynamoDB ロックマネージャー](#) - データ圧縮を使用する場合、他のデータロードジョブは `org.apache.iceberg.aws.dynamodb.lock-impl` としてを使用しないでください `DynamoDbLockManager`。

データカタログの管理

AWS Glue Data Catalog は、Amazon S3 データセットの構造メタデータと運用メタデータを保存する中央メタデータリポジトリです。データカタログを効果的に管理することは、データ品質、パフォーマンス、セキュリティ、ガバナンスを維持する上で重要です。

これらのデータカタログ管理プラクティスを理解して適用することで、データランドスケープの進化に合わせて、メタデータが正確でパフォーマンスが高く、安全で、十分に管理されるようにすることができます。

このセクションでは、データカタログ管理の以下の側面について説明します。

- テーブルスキーマとパーティションの更新: データの進化に伴い、データカタログで定義されているテーブルスキーマまたはパーティション構造を更新する必要がある場合があります。AWS Glue ETL を使用してプログラムでこれらの更新を行う方法の詳細については、「[AWS Glue ETL ジョブを使用してデータカタログのスキーマを更新し、新規パーティションを追加する](#)」を参照してください。
- 列統計の管理: 正確な列統計は、クエリプランの最適化とパフォーマンスの向上に役立ちます。列統計を生成、更新、管理する方法については、「[列統計を使用したクエリのパフォーマンスの最適化](#)」を参照してください。
- データカタログの暗号化: 機密メタデータを保護するために、AWS Key Management Service (AWS KMS) を使用してデータカタログを暗号化できます。このセクションでは、データカタログの暗号化を有効にして管理する方法について説明します。
- AWS Lake Formation を使用したデータカタログの保護: Lake Formation には、データレイクのセキュリティとアクセスコントロールに対する包括的なアプローチが用意されています。Lake Formation を使用して、データカタログおよび基盤となるデータへのアクセスを保護および管理できます。

トピック

- [AWS Glue ETL ジョブを使用してデータカタログのスキーマを更新し、新規パーティションを追加する](#)
- [列統計を使用したクエリのパフォーマンスの最適化](#)
- [Data Catalog の暗号化](#)
- [Lake Formation を使用したデータカタログの保護](#)

AWS Glue ETL ジョブを使用してデータカタログのスキーマを更新し、新規パーティションを追加する

抽出、変換、ロード (ETL) ジョブによって、ターゲットデータストアに新しいテーブルパーティションが作成される場合があります。データセットスキーマは、時間の経過とともに AWS Glue Data Catalog スキーマから進化し、拡散する可能性があります。AWS Glue ETL ジョブには、ETL スクリプト内で使用できるいくつかの機能が用意され、Data Catalog でスキーマおよびパーティションを更新できるようになりました。これらの機能を使用すると、クローラを再実行することなく、Data Catalog での ETL 作業の結果を確認できます。

新しいパーティション

AWS Glue Data Catalog で新しいパーティションを表示するには、次のいずれかを実行します。

- ジョブが終了したら、クローラを再実行し、クローラの完了時にコンソールで新しいパーティションを表示します。
- ジョブが終了すると、クローラを再実行することなく、コンソールで新しいパーティションがすぐに表示されます。この機能は、次の例に示すように、ETL スクリプトに数行のコードを追加して有効にすることができます。このコードは `enableUpdateCatalog` 引数を使用して、新しいパーティションの作成時のジョブ実行中に Data Catalog を更新する必要があることを示します。

方式 1

オプション引数に `enableUpdateCatalog` と `partitionKeys` を渡します。

Python

```
additionalOptions = {"enableUpdateCatalog": True}
additionalOptions["partitionKeys"] = ["region", "year", "month", "day"]

sink = glueContext.write_dynamic_frame_from_catalog(frame=last_transform,
    database=<target_db_name>,

    table_name=<target_table_name>, transformation_ctx="write_sink",

    additional_options=additionalOptions)
```

Scala

```
val options = JsonOptions(Map(
  "path" -> <S3_output_path>,
  "partitionKeys" -> Seq("region", "year", "month", "day"),
  "enableUpdateCatalog" -> true))
val sink = glueContext.getCatalogSink(
  database = <target_db_name>,
  tableName = <target_table_name>,
  additionalOptions = options)sink.writeDynamicFrame(df)
```

方式 2

enableUpdateCatalog に partitionKeys と getSink() を渡して、DataSink オブジェクトで setCatalogInfo() を呼び出します。

Python

```
sink = glueContext.getSink(
  connection_type="s3",
  path="<S3_output_path>",
  enableUpdateCatalog=True,
  partitionKeys=["region", "year", "month", "day"])
sink.setFormat("json")
sink.setCatalogInfo(catalogDatabase=<target_db_name>,
  catalogTableName=<target_table_name>)
sink.writeFrame(last_transform)
```

Scala

```
val options = JsonOptions(
  Map("path" -> <S3_output_path>,
    "partitionKeys" -> Seq("region", "year", "month", "day"),
    "enableUpdateCatalog" -> true))
val sink = glueContext.getSink("s3", options).withFormat("json")
sink.setCatalogInfo(<target_db_name>, <target_table_name>)
sink.writeDynamicFrame(df)
```

クローラを再実行することなく、AWS Glue ETLジョブ自体を使用して、Data Catalog での新しいカタログテーブルの作成、変更されたスキーマによる既存のテーブルの更新、および新しいテーブルパーティションの追加が可能になりました。

テーブルスキーマの更新

Data Catalog テーブルのスキーマを上書きする場合は、次のいずれかを実行します。

- ジョブが終了したら、クローラを再実行し、テーブル定義も更新するようにクローラが設定されていることを確認します。クローラが終了したら、新しいパーティションをスキーマの更新とともにコンソールに表示します。詳細については、「[API を使用したクローラの設定](#)」を参照してください。
- ジョブが終了すると、クローラを再実行することなく、コンソールで変更済みスキーマがすぐに表示されます。この機能は、次の例に示すように、ETL スクリプトに数行のコードを追加して有効にすることができます。このコードでは `enableUpdateCatalog` を `true` に設定し、`updateBehavior` を `UPDATE_IN_DATABASE` に設定しています。これは、ジョブ実行中にスキーマを上書きし、Data Catalog に新しいパーティションを追加することを示します。

Python

```
additionalOptions = {
    "enableUpdateCatalog": True,
    "updateBehavior": "UPDATE_IN_DATABASE"}
additionalOptions["partitionKeys"] = ["partition_key0", "partition_key1"]

sink = glueContext.write_dynamic_frame_from_catalog(frame=last_transform,
    database=<dst_db_name>,
    table_name=<dst_tbl_name>, transformation_ctx="write_sink",
    additional_options=additionalOptions)
job.commit()
```

Scala

```
val options = JsonOptions(Map(
    "path" -> outputPath,
    "partitionKeys" -> Seq("partition_0", "partition_1"),
    "enableUpdateCatalog" -> true))
val sink = glueContext.getCatalogSink(database = nameSpace, tableName = tableName,
    additionalOptions = options)
sink.writeDynamicFrame(df)
```

テーブルスキーマが上書きされないようにして、新しいパーティションを追加する場合は、`updateBehavior` 値を `LOG` に設定することもできます。`updateBehavior` のデフォルト値は

UPDATE_IN_DATABASE です。そのため、明示的に定義しない場合、テーブルスキーマは上書きされます。

enableUpdateCatalog が true に設定されていない場合、updateBehavior で選択したオプションに関係なく、ETL ジョブは Data Catalog 内のテーブルを更新しません。

新しいテーブルの作成

同じオプションを使用して、Data Catalog で新しいテーブルを作成することもできます。setCatalogInfo を使用して、データベースと新しいテーブル名を指定できます。

Python

```
sink = glueContext.getSink(connection_type="s3", path="s3://path/to/data",
    enableUpdateCatalog=True, updateBehavior="UPDATE_IN_DATABASE",
    partitionKeys=["partition_key0", "partition_key1"])
sink.setFormat("<format>")
sink.setCatalogInfo(database=<dst_db_name>, catalogTableName=<dst_tbl_name>)
sink.writeFrame(last_transform)
```

Scala

```
val options = JsonOptions(Map(
  "path" -> outputPath,
  "partitionKeys" -> Seq("<partition_1>", "<partition_2>"),
  "enableUpdateCatalog" -> true,
  "updateBehavior" -> "UPDATE_IN_DATABASE"))
val sink = glueContext.getSink(connectionType = "s3", connectionOptions =
  options).withFormat("<format>")
sink.setCatalogInfo(database = "<dst_db_name>", catalogTableName =
  "<dst_tbl_name>")
sink.writeDynamicFrame(df)
```

制限事項

次の制限事項に注意してください。

- Amazon Simple Storage Service (Amazon S3) のターゲットのみがサポートされています。
- enableUpdateCatalog 機能は、管理対象テーブルでは、サポートされていません。
- json、csv、avro、および parquet の形式のみがサポートされます。

- parquet 分類でテーブルを作成または更新するには、AWS Glue に最適化された DynamicFrames 用 parquet ライターを使用する必要があります。これは、次のいずれかの方法で実現できます。
- parquet 分類を使用してカタログ内の既存のテーブルを更新する場合は、更新前にテーブルでは "useGlueParquetWriter" テーブルプロパティを true に設定する必要があります。このプロパティは、AWS Glue API/SDK、コンソール、または Athena DDL ステートメントから設定できます。

The screenshot shows the AWS Glue console interface for editing a table. The left sidebar contains navigation options like 'Data Catalog tables', 'Data Catalog', and 'Data Integration and ETL'. The main content area is titled 'Edit table' and has three sections: 'Table details', 'Serde parameters', and 'Table properties'. The 'Table properties' section is a table with columns for 'Key', 'Value', and 'Remove'. The table contains several properties such as 'skip.header.line.count', 'has_encrypted_data', 'columnsOrdered', 'areColumnsQuoted', 'delimiter', 'classification', and 'typeOfData'. At the bottom of this table, there is an 'Add' button and two input fields for 'Enter a unique key' and 'Enter a value'. A red box highlights the 'Add' button and the input fields.

Key	Value	Remove
skip.header.line.count	1	Remove
has_encrypted_data	false	Remove
columnsOrdered	true	Remove
areColumnsQuoted	false	Remove
delimiter	,	Remove
classification	csv	Remove
typeOfData	file	Remove
Enter a unique key	Enter a value	Remove

カタログテーブルプロパティを設定したら、次のコードスニペットを使用して新しいデータでカタログテーブルを更新します。

```
glueContext.write_dynamic_frame.from_catalog(
    frame=frameToWrite,
    database="dbName",
    table_name="tableName",
    additional_options={
        "enableUpdateCatalog": True,
        "updateBehavior": "UPDATE_IN_DATABASE"
```

```
}  
)
```

- テーブルがカタログ内にまだ存在しない場合は、`connection_type="s3"` を用いた `getSink()` メソッドをスクリプトで使用して、データを Amazon S3 に書き込むとともに、テーブルとそのパーティションをカタログに追加します。ワークフローに適切な `partitionKeys` と `compression` を指定します。

```
s3sink = glueContext.getSink(  
    path="s3://bucket/folder/",  
    connection_type="s3",  
    updateBehavior="UPDATE_IN_DATABASE",  
    partitionKeys=[],  
    compression="snappy",  
    enableUpdateCatalog=True  
)  
  
s3sink.setCatalogInfo(  
    catalogDatabase="dbName", catalogTableName="tableName"  
)  
  
s3sink.setFormat("parquet", useGlueParquetWriter=true)  
s3sink.writeFrame(frameToWrite)
```

- `glueparquet` フォーマットの値は、AWS Glue parquet ライターを有効にする従来のメソッドです。
- `updateBehavior` を LOG に設定した場合、DynamicFrame スキーマが Data Catalog テーブルのスキーマに定義されている列のサブセットと同等であるか、またはそのサブセットを含んでいる場合にのみ、新しいパーティションが追加されます。
- スキーマの更新は、パーティション化されていないテーブル (「`partitionKeys`」オプションを使用していない) ではサポートされていません。
- `partitionKeys` は、ETL スクリプトで渡されるパラメータと Data Catalog テーブルスキーマの `partitionKeys` で同等で、同じ順序でなければなりません。
- この機能は、現在、更新スキーマがネストされているテーブルの更新/作成をサポートしていません (例えば、構造体の内部の配列)。

詳細については、[the section called “AWS Glue for Spark”](#) を参照してください。

ETL ジョブでの MongoDB 接続の操作

MongoDB の接続を作成し、その接続を AWS Glue ジョブで使用します。詳細については、「AWS Glue プログラミングガイド」の「[the section called “MongoDB 接続”](#)」を参照してください。接続の url、username、および password は MongoDB 接続に保存されます。ETL ジョブスクリプトで `glueContext.getCatalogSource` という `additionalOptions` パラメータを使用して、その他のオプションを指定できます。その他のオプションには次のようなものがあります。

- `database`: (必須) 読み込み元の MongoDB データベース。
- `collection`: (必須) 読み込み元の MongoDB コレクション。

ETL ジョブスクリプト内に `database` および `collection` 情報を書き込むと、複数のジョブに同じ接続を使用できます。

1. AWS Glue Data Catalog 接続を MongoDB データソース用に作成します。接続パラメータの説明については、「["connectionType": "mongodb"](#)」を参照してください。接続は、コンソール、API、または CLI を使用して作成できます。
2. AWS Glue Data Catalog でデータベースを作成して MongoDB データのテーブル定義を格納します。詳細については、「[データベースの作成](#)」を参照してください。
3. MongoDB に接続するための接続内の情報を使用して、MongoDB 内のデータをクローラするクローラを作成します。クローラによって、ジョブで使用する MongoDB データベース内のテーブルを記述するテーブルが AWS Glue Data Catalog に作成されます。詳細については、「[クローラを使用したデータカタログへの入力](#)」を参照してください。
4. カスタムスクリプトを使用してジョブを作成します。ジョブは、コンソール、API、または CLI を使用して作成できます。詳細については、「[AWS Glue でジョブを追加する](#)」を参照してください。
5. ジョブ用のデータターゲットを選択します。データターゲットを表すテーブルは、Data Catalog で定義することも、ジョブを実行するときに作成することもできます。ジョブを作成するときに、ターゲットの場所を選択します。ターゲットに接続が必要な場合、接続もジョブ内で参照されます。ジョブに複数のデータターゲットが必要な場合、後でスクリプトを編集して追加できます。
6. ジョブと生成されたスクリプトの引数を指定することで、ジョブ処理環境をカスタマイズします。

Data Catalog で定義されたテーブル構造に基づいて、MongoDB データベースから DynamicFrame を作成する例を次に示します。このコードは `additionalOptions` を使用して、追加のデータソースの情報を指定しています。

Scala

```
val resultFrame: DynamicFrame = glueContext.getCatalogSource(  
    database = catalogDB,  
    tableName = catalogTable,  
    additionalOptions = JsonOptions(Map("database" -> DATABASE_NAME,  
        "collection" -> COLLECTION_NAME))  
).getDynamicFrame()
```

Python

```
glue_context.create_dynamic_frame_from_catalog(  
    database = catalogDB,  
    table_name = catalogTable,  
    additional_options = {"database": "database_name",  
        "collection": "collection_name"})
```

7. オンデマンドで、またはトリガーを使用してジョブを実行します。

列統計を使用したクエリのパフォーマンスの最適化

追加のデータパイプラインを設定することなく、Parquet、ORC、JSON、ION、CSV、XML などのデータ形式で AWS Glue Data Catalog テーブルの列レベルの統計を計算できます。列統計は、列内の値に関するインサイトを得ることで、データプロファイルを理解するのに役立ちます。データカタログは、最小値、最大値、null 値の合計、個別の値の合計、値の平均長、true 値の合計出現数などの列の値の統計の生成をサポートしています。

Amazon Redshift および Amazon Athena などの AWS 分析サービスは、これらの列統計を使用してクエリ実行プランを生成し、クエリのパフォーマンスを改善する最適なプランを選択できます。

AWS Glue コンソールまたは AWS CLI を使用して列統計生成タスクを実行するように設定できます。プロセスを開始すると、AWS Glue はバックグラウンドで Spark ジョブを開始し、データカタログ内の AWS Glue テーブルメタデータを更新します。列統計は、AWS Glue コンソールもしくは AWS CLI を使用して、または [GetColumnStatisticsForTable](#) API オペレーションを呼び出すことによって表示できます。

Note

Lake Formation の許可を使用してテーブルに対するアクセスを制御している場合、列統計タスクによって引き受けられるロールには、統計を生成するための完全なテーブルアクセスが必要です。

トピック

- [列統計を生成するための前提条件](#)
- [列統計の生成](#)
- [列統計の表示](#)
- [列統計の生成](#)
- [列統計の削除](#)
- [列統計タスクの実行の表示](#)
- [列統計タスクの実行の停止](#)
- [考慮事項と制約事項](#)

列統計を生成するための前提条件

列統計を生成または更新するために、統計生成タスクはユーザーに代わって AWS Identity and Access Management (IAM) ロールを引き受けます。ロールに付与された許可に基づいて、列統計生成タスクは Amazon S3 データストアからデータを読み取ることができます。

Note

Lake Formation によって管理されるテーブルの統計を生成するには、統計の生成に使用される IAM ロールに完全なテーブルアクセスが必要です。

ロールベースのアクセスコントロールを使用するには、以下のポリシーにリストされている許可を持つ IAM ロールを作成し、そのロールを列統計生成タスクに追加する必要があります。

列統計を生成するための IAM ロールを作成するには

1. IAM ロールを作成するには、「[AWS Glue の IAM ロールを作成する](#)」を参照してください。

2. 既存のロールを更新するには、IAM コンソールで、列統計の生成プロセスで使用されている IAM ロールにアクセスします。
3. [許可を追加] セクションで、[ポリシーをアタッチ] を選択します。新しく開いたブラウザウィンドウで、AWSGlueServiceRole AWS マネージドポリシーを選択します。
4. Amazon S3 のデータの場所からデータを読み取るための許可も含める必要があります。

[Add permissions] (アクセス許可の追加) セクションで、[Create policy] (ポリシーの作成) を選択します。新しく開いたブラウザウィンドウで、ロールで使用する新しいポリシーを作成します。

5. [ポリシーを作成] ページで、[JSON] タブを選択します。次の JSON コードをポリシーエディタフィールドにコピーします。

Note

次のポリシーでは、アカウント ID を有効な AWS アカウント に置き換え、region をテーブルのリージョンに置き換え、bucket-name を Amazon S3 バケット名に置き換えます。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "S3BucketAccess",
      "Effect": "Allow",
      "Action": [
        "s3:ListBucket",
        "s3:GetObject"
      ],
      "Resource": [
        "arn:aws:s3:::<bucket-name>/*",
        "arn:aws:s3:::<bucket-name>"
      ]
    }
  ]
}
```

6. (オプション) Lake Formation の許可を使用してデータに対するアクセスを提供している場合、IAM ロールには lakeformation:GetDataAccess 許可が必要です。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "LakeFormationDataAccess",
      "Effect": "Allow",
      "Action": "lakeformation:GetDataAccess",
      "Resource": [
        "*"
      ]
    }
  ]
}
```

Amazon S3 のデータの場所が Lake Formation に登録されており、テーブルに対する IAM_ALLOWED_PRINCIPALS グループ許可が、列統計生成タスクによって引き受けられる IAM ロールに付与されていない場合、そのロールには、テーブルに対する Lake Formation ALTER および DESCRIBE 許可が必要です。Amazon S3 バケットの登録に使用されるロールには、テーブルに対する Lake Formation INSERT および DELETE 許可が必要です。

Amazon S3 のデータの場所が Lake Formation に登録されておらず、テーブルに対する IAM_ALLOWED_PRINCIPALS グループ許可が IAM ロールに付与されていない場合、そのロールには、テーブルに対する Lake Formation ALTER、DESCRIBE、INSERT、および DELETE 許可が必要です。

7. (オプション) 暗号化された Amazon CloudWatch Logs を書き込む列統計生成タスクには、キーポリシーで次の許可が必要です。

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Sid": "CWLogsKmsPermissions",
    "Effect": "Allow",
    "Action": [
      "logs:CreateLogGroup",
      "logs:CreateLogStream",
      "logs:PutLogEvents",
      "logs:AssociateKmsKey"
    ],
    "Resource": [
```

```

        "arn:aws:logs:<region>:111122223333:log-group:/aws-glue:*"
    ]
},
{
    "Sid": "KmsPermissions",
    "Effect": "Allow",
    "Action": [
        "kms:GenerateDataKey",
        "kms:Decrypt",
        "kms:Encrypt"
    ],
    "Resource": [
        "arn:aws:kms:<region>:111122223333:key/"arn of key used for ETL cloudwatch
        encryption"
    ],
    "Condition": {
        "StringEquals": {
            "kms:ViaService": ["glue.<region>.amazonaws.com"]
        }
    }
}
]
}

```

8. 列統計の実行に使用するロールには、そのロールに対する `iam:PassRole` のアクセス許可が必要です。

```

{
    "Version": "2012-10-17",
    "Statement": [{
        "Effect": "Allow",
        "Action": [
            "iam:PassRole"
        ],
        "Resource": [
            "arn:aws:iam::111122223333:role/<columnstats-role-name>"
        ]
    }]
}

```

9. 列統計を生成するための IAM ロールを作成する場合、そのロールには、サービスがそのロールを引き受けることを可能にする次の信頼ポリシーも必要です。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "TrustPolicy",
      "Effect": "Allow",
      "Principal": {
        "Service": "glue.amazonaws.com"
      },
      "Action": "sts:AssumeRole",
    }
  ]
}
```

列統計の生成

AWS Glue コンソールまたは AWS CLI を使用して、データカタログでの統計生成を管理するには、次のステップに従います。

Console

コンソールを使用して列統計を生成するには

1. AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) にサインインします。
2. データカタログテーブルを選択します。
3. リストからテーブルを選択します。
4. [アクション] メニューで [統計を生成] を選択します。

[テーブル] ページの下のセクションにある [列統計] タブで [統計を生成] ボタンを選択することもできます。

5. [統計を生成] ページで、次のオプションを指定します。

Generate statistics

Generate column statistics for the table to improve query performance and potentially save costs. [View pricing](#)

Choose columns

Table (All columns)

Generate statistics for all columns.

Selected columns

Choose the columns to generate statistics.

Row sampling options

We recommend to use all rows to compute accurate column statistics. You can use sampling when the dataset is potentially large and approximate results are acceptable.

All rows

Generate column statistics on entire data.

Sample rows

Generate approximate statistics using sample rows.

IAM role

To generate statistics, the IAM role assumed by the job should have necessary permissions. [Learn more](#)

Choose an existing IAM role

12495-pentestRole



[View](#)

[Create new IAM role](#)

▶ Security configuration - optional

Enable at-rest encryption with a security configuration.

Cancel

Generate statistics

- テーブル (すべての列) – テーブル内のすべての列の統計を生成するには、このオプションを選択します。
- 選択された列 – 特定の列の統計を生成するには、このオプションを選択します。ドロップダウンリストから列を選択できます。
- すべての行 – 正確な統計を生成するには、テーブルからすべての行を選択します。
- サンプル行 – テーブルから特定の割合の行のみを選択して統計を生成します。デフォルトは、すべての行です。上矢印と下矢印を使用してパーセント (%) の値を増減します。

Note

正確な統計を計算するには、テーブル内のすべての行を含めることをお勧めします。近似値が許容される場合にのみ、サンプル行を使用して列統計を生成します。

6. (オプション) 次に、ログの保管中の暗号化を有効にするセキュリティ設定を選択します。
7. [統計を生成] を選択してプロセスを実行します。

AWS CLI

次の例では、DatabaseName、TableName、および ColumnNameList の値を、実際のデータベース、テーブル、および列の名前に置き換えます。アカウント ID を有効な AWS アカウントに置き換え、統計の生成に使用している IAM ロールの名前にロール名を置き換えます。

```
aws glue start-column-statistics-task-run --input-cli-json file://input.json
{
  "DatabaseName": "<test-db>",
  "TableName": "<test-table>",
  "ColumnNameList": [
    "<column1>",
    "<column2>",
  ],
  "Role": "arn:aws:iam::<123456789012>:role/<Stats-Role>",
  "SampleSize": 10.0
}
```

[StartColumnStatisticsTaskRun](#) オペレーションを呼び出して列統計を生成することもできます。

列統計の表示

統計が正常に生成されると、データカタログは、Amazon Athena と Amazon Redshift のコストベースのオプティマイザーがクエリを実行する際に最適な選択を行うことができるように、この情報を保存します。統計は列のタイプによって異なります。

AWS Management Console

テーブルの列統計を表示するには

- 列統計タスクを実行すると、[テーブルの詳細] ページの [列統計] タブにテーブルの統計が表示されます。

AWS Glue > Tables > pentest_orders_xml

pentest_orders_xml Last updated (UTC)
October 25, 2023 at 19:14:47 Version 15 (Current version) Actions

Table overview | Data quality New

Table details | Advanced properties

Name pentest_orders_xml	Description -	Database pentest_db	Classification XML
Location s3://kietduon-column-statistics-table/orders.xml	Connection -	Deprecated -	Last updated October 25, 2023 at 19:14:47
Input format -	Output format -	Serde serialization lib -	

Schema | Partitions | Indexes | **Column statistics - new**

Column statistics (9) Last updated (UTC)
November 6, 2023 at 21:50:40 Stop View all runs Generate statistics

Get an overview of the data profile. We estimate the approximate number of distinct values in a data set with 5% average relative error.

Find columns

Column name	Last updated (UTC)	Average length	Distinct values	Max length	Null values	Max value	Min value	True values	False values
o_clerk	October 25, 2023 at 19:14:	15.00	919	15	-	-	-	-	-
o_comment	October 25, 2023 at 19:14:	88.38	3156	124559	-	-	-	-	-
o_custkey	October 25, 2023 at 19:14:	-	919	-	-	1499	1	-	-
o_order-priority	October 25, 2023 at 19:14:	8.45	5	15	-	-	-	-	-
o_orderdate	October 25, 2023 at 19:14:	10.00	1790	10	-	-	-	-	-
o_orderkey	October 25, 2023 at 19:14:	-	3098	-	-	12451	1	-	-
o_orderstatus	October 25, 2023 at 19:14:	1.00	3	1	-	-	-	-	-
o_ship-priority	October 25, 2023 at 19:14:	-	1	-	-	-	-	-	-
o_totalprice	October 25, 2023 at 19:14:	-	3062	-	-	422359.65	974.04	-	-

次の統計を使用できます。

- [列名]: 統計を生成するために使用される列名
- [最終更新日時]: 統計が生成された日時
- [平均長]: 列内の値の平均長
- [個別の値]: 列内の個別の値の合計数。列内の個別の値の数を 5% の相対誤差で推定します。
- [最大値]: 列内の最大値。
- [最小値]: 列内の最小値。
- [最大長]: 列内の最大値の長さ。
- [Null 値]: 列内の null 値の合計数。
- [True 値]: 列内の true 値の合計数。
- [False 値]: 列内の false 値の合計数。

AWS CLI

次の例は、AWS CLI を使用して列統計を取得する方法を示しています。

```
aws glue get-column-statistics-for-table \
```

列統計を使用したクエリのパフォーマンスの最適化

```
--database-name <test_db> \  
--table-name <test_tble> \  
--column-names <col1>
```

[GetColumnStatisticsForTable](#) API オペレーションを使用して列統計を表示することもできます。

列統計の生成

統計を最新状態に保つことで、クエリプランナーが最適なプランを選択できるようになるため、クエリのパフォーマンスが向上します。列の統計を更新するには、AWS Glue コンソールから [統計を生成] タスクを明示的に実行する必要があります。データカタログは統計を自動的に更新しません。

コンソールで AWS Glue の統計生成機能を使用していない場合は、[UpdateColumnStatisticsForTable](#) API オペレーションまたは AWS CLI を使用して列統計を手動で更新できます。次の例は、AWS CLI を使用して列統計を更新する方法を示しています。

```
aws glue update-column-statistics-for-table --cli-input-json:
```

```
{  
  "CatalogId": "111122223333",  
  "DatabaseName": "test_db",  
  "TableName": "test_table",  
  "ColumnStatisticsList": [  
    {  
      "ColumnName": "col1",  
      "ColumnType": "Boolean",  
      "AnalyzedTime": "1970-01-01T00:00:00",  
      "StatisticsData": {  
        "Type": "BOOLEAN",  
        "BooleanColumnStatisticsData": {  
          "NumberOfTrues": 5,  
          "NumberOfFalses": 5,  
          "NumberOfNulls": 0  
        }  
      }  
    }  
  ]  
}
```

列統計の削除

[DeleteColumnStatisticsForTable](#) API オペレーションまたは AWS CLI を使用して、列統計を削除できます。次の例は、AWS Command Line Interface (AWS CLI) を使用して列統計を削除する方法を示しています。

```
aws glue delete-column-statistics-for-table \  
  --database-name test_db \  
  --table-name test_table \  
  --column-name col1
```

列統計タスクの実行の表示

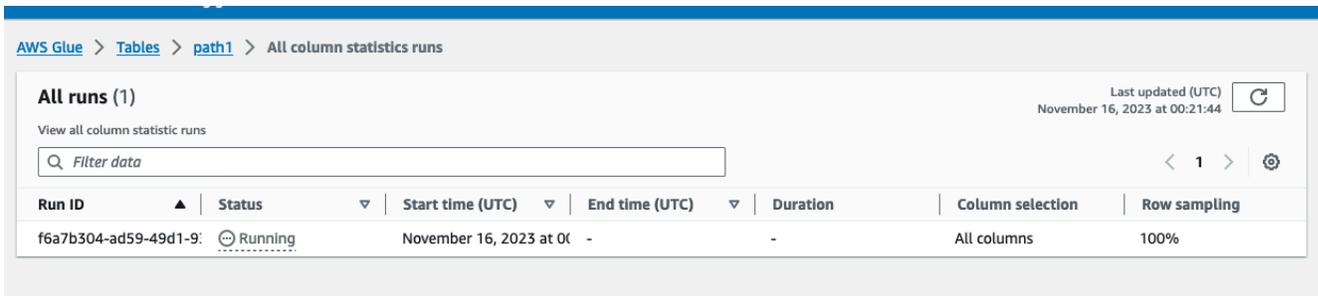
列統計タスクを実行した後、AWS Glue コンソール、AWS CLI、または [GetColumnStatisticsTaskRuns](#) オペレーションを使用して、テーブルのタスク実行の詳細を調査できます。

Console

列統計タスクの実行の詳細を表示するには

1. AWS Glue コンソールで、[データカタログ] の下の [テーブル] を選択します。
2. 列統計を含むテーブルを選択します。
3. [テーブルの詳細] ページで、[列統計] を選択します。
4. [実行を表示] を選択します。

指定されたテーブルに関連付けられたすべての実行に関する情報を確認できます。



The screenshot shows the AWS Glue console interface for viewing column statistics runs. The breadcrumb navigation is "AWS Glue > Tables > path1 > All column statistics runs". The main heading is "All runs (1)" with a refresh icon and the text "Last updated (UTC) November 16, 2023 at 00:21:44". Below the heading is a search bar labeled "Filter data". A table displays the run details:

Run ID	Status	Start time (UTC)	End time (UTC)	Duration	Column selection	Row sampling
f6a7b304-ad59-49d1-9	Running	November 16, 2023 at 00:21:44	-	-	All columns	100%

AWS CLI

次の例では、DatabaseName と TableName の値を、実際のデータベースとテーブルの名前に置き換えます。

```
aws glue get-column-statistics-task-runs --input-cli-json file://input.json
{
  "DatabaseName": "<test-db>",
  "TableName": "<test-table>"
}
```

列統計タスクの実行の停止

テーブルについての列統計タスクの実行は、AWS Glue コンソール、AWS CLI、または [StopColumnStatisticsTaskRun](#) オペレーションを使用して停止できます。

Console

列統計タスクの実行を停止するには

1. AWS Glue コンソールで、[データカタログ] の下の [テーブル] を選択します。
2. 列統計タスクの実行が進行中のテーブルを選択します。
3. [テーブルの詳細] ページで、[列統計] を選択します。
4. [Stop] (停止) を選択します。

実行が完了する前にタスクを停止すると、テーブルの列統計は生成されません。

AWS CLI

次の例では、DatabaseName と TableName の値を、実際のデータベースとテーブルの名前に置き換えます。

```
aws glue stop-column-statistics-task-run --input-cli-json file://input.json
{
  "DatabaseName": "<test-db>",
  "TableName": "<test-table>"
}
```

考慮事項と制約事項

列統計の生成には、次の考慮事項と制限事項が適用されます。

考慮事項

- サンプルングを使用して統計を生成すると、実行時間は短縮されますが、不正確な統計が生成される可能性があります。
- 各列統計の実行では、データセット全体の処理が必要です。
- データカタログは、異なるバージョンの統計を保存しません。
- テーブルごとに一度に 1 つの統計生成タスクのみを実行できます。
- データカタログに登録されているカスタマー AWS KMS キーを使用してテーブルが暗号化されている場合、AWS Glue は同じキーを使用して統計を暗号化します。

列統計タスクは、次の場合に統計の生成をサポートします。

- 完全なテーブル許可が IAM ロールに付与されている場合 (IAM または Lake Formation)。
- Lake Formation ハイブリッドアクセスモードを使用するテーブルに対する許可が IAM ロールに付与されている場合。

列統計タスクは、次についての統計の生成をサポートしていません。

- Lake Formation のセルベースのアクセスコントロールを備えたテーブル。
- トランザクションデータレイク - Linux Foundation Delta Lake、Apache Iceberg、Apache Hudi。
- フェデレーテッドデータベース内のテーブル - Hive メタストア、Amazon Redshift データ共有
- ネストされた列、配列、および構造体のデータ型。
- 別のアカウントから共有されているテーブル。

Data Catalog の暗号化

AWS Key Management Service (AWS KMS) によって管理される暗号化キーを使用して、保管中の AWS Glue Data Catalog に保存されているメタデータを保護できます。データカタログ設定を使用して、新しいデータカタログのデータカタログ暗号化を有効にできます。必要に応じて、既存のデータカタログの暗号化を有効または無効にできます。有効にすると、AWS Glue はカタログに書き込まれたすべての新しいメタデータを暗号化しますが、既存のメタデータは暗号化されません。

データカタログの暗号化の詳細については、「[Data Catalog の暗号化](#)」を参照してください。

Lake Formation を使用したデータカタログの保護

AWS Lake Formation は、AWS で簡単にセキュアなデータレイクを構築できるサービスです。きめ細かなアクセスコントロール許可を定義することで、データレイクを作成および安全に管理するための一元的な場所を提供します。Lake Formation は、データカタログを使用して、テーブル定義、スキーマ情報、データアクセスコントロール設定など、データレイクに関するメタデータを保存および取得します。

メタデータテーブルまたはデータベースの Amazon S3 データロケーションを Lake Formation に登録し、それを使用してデータカタログリソースに対するメタデータレベルのアクセス許可を定義できます。また、Lake Formation を使用して、統合された分析エンジンに代わって、Amazon S3 に保存されている基盤となるデータへのストレージアクセス許可も管理できます。

詳細については、「[AWS Lake Formation とは](#)」を参照してください。

データカタログにアクセスする

AWS Glue Data Catalog を使用して、データを検出して理解できます。データカタログは、スキーマの定義、データ型、場所などのメタデータを一貫して維持する方法を提供します。次の方法を使用してデータカタログにアクセスできます。

- AWS Glue コンソール – ウェブベースのユーザーインターフェイスである AWS Glue コンソールからデータカタログにアクセスして管理できます。コンソールでは、データベース、テーブル、および関連するメタデータを参照および検索したり、メタデータ定義を作成、更新、削除したりできます。
- AWS Glue クローラー – クローラーは、データソースを自動的にスキャンし、データカタログにメタデータを入力するプログラムです。クローラーを作成して実行し、Amazon S3、Amazon RDS、Amazon DynamoDB、Amazon CloudWatch、および MySQL や PostgreSQL などの JDBC 準拠のリレーショナルデータベース、Snowflake や Google BigQuery などのいくつかの非 AWS ソースからのデータを検出してカタログ化できます。
- AWS Glue API – AWS Glue API を使用してプログラムでデータカタログにアクセスできます。これらの API を使用すると、プログラムでデータカタログを操作し、自動化と他のアプリケーションやサービスとの統合が可能になります。
- AWS Command Line Interface (AWS CLI) – AWS CLI を使用して、コマンドラインからデータカタログにアクセスして管理できます。CLI には、メタデータ定義の作成、更新、削除、メタデータ情報のクエリと取得のためのコマンドが用意されています。

- 他の AWS サービスとの統合 – データカタログは他のさまざまな AWS サービスと統合されるため、カタログに保存されているメタデータにアクセスして利用できます。たとえば、Amazon Athena を使用してデータカタログのメタデータを使用してデータソースをクエリし、AWS Lake Formation を使用してデータカタログリソースのデータアクセスとガバナンスを管理できます。

AWS Glue データカタログのベストプラクティス

このセクションでは、AWS Glue Data Catalog を効果的に管理および活用するためのベストプラクティスについて説明します。効率的なクローラーの使用、メタデータの整理、セキュリティ、パフォーマンスの最適化、自動化、データガバナンス、他の AWS サービスとの統合などのプラクティスについて重点的に説明します。

- クローラーを効果的に使用する – クローラーを定期的に行って、データソースの変更に応じてデータカタログを最新の状態に保ちます。頻繁に変化するデータソースには増分クローラーを使用してパフォーマンスを向上させます。変更が検出されたときに新しいパーティションを自動的に追加したり、スキーマを更新したりするようにクローラーを設定します。
- メタデータテーブルの整理と名前付け – データカタログのデータベースとテーブルについて一貫した命名規則を確立します。関連データソースを論理データベースまたはフォルダにグループ化して、より適切に整理します。各テーブルの目的と内容を伝えるわかりやすい名前を使用します。
- スキーマを効果的に管理する – AWS Glue クローラーのスキーマ推論機能を活用します。ダウンストリームアプリケーションが破損しないように、スキーマの変更を適用する前に確認および更新します。スキーマ進化機能を使用して、スキーマの変更を適切に処理します。
- データカタログの保護 – データカタログの保管中および転送中のデータ暗号化を有効にします。きめ細かなアクセスコントロールポリシーを実装し、機密データへのアクセスを制限します。データカタログのアクセス許可とアクティビティログを定期的に監査して確認します。
- 他の AWS サービスとの統合 - データカタログを Amazon Athena、Redshift Spectrum、AWS Lake Formation などのサービスを一元管理できるメタデータレイヤーとして使用します。AWS Glue ETL ジョブを活用して、データカタログのメタデータを維持しながら、データを変換してさまざまなデータストアにロードします。
- パフォーマンスのモニタリングと最適化 - Amazon CloudWatch メトリクスを使用してクローラーと ETL ジョブのパフォーマンスをモニタリングします。大きなデータセットをデータカタログにパーティション分割して、クエリのパフォーマンスを向上します。頻繁にアクセスされるメタデータのパフォーマンス最適化を実装します。
- AWS Glue ドキュメントとベストプラクティスの最新情報を入手 - AWS Glue ドキュメントと AWS Glue リソースで最新の更新、ベストプラクティス、推奨事項を定期的に確認します。AWS

Glue ウェビナー、ワークショップ、その他のイベントに参加して、エキスパートから学び、新機能について最新情報を得ます。

AWS Glue スキーマレジストリ

Note

AWS Glue スキーマレジストリは、アジアパシフィック (ジャカルタ) と中東 (UAE) リージョンの AWS Glue コンソールではサポートされていません。。

AWS Glue Schema Registry は、データストリームスキーマを一元的に検出、制御し、発展させることができる新機能です。スキーマは、データレコードの構造と形式を定義します。AWS Glue のスキーマレジストリを使用すると、Apache Kafka、[「Amazon Managed Streaming for Apache Kafka」](#)、[「Amazon Kinesis Data Streams」](#)、[「Amazon Managed Service for Apache Flink」](#)、[「AWS Lambda」](#) と便利な統合を使用するデータストリーミングアプリケーションでスキーマを管理して実施することができます。

AWS Glue の Schema Registry では、AVRO (v1.10.2) データ形式、[Everit ライブラリ](#) による JSON スキーマ検証を含むスキーマ (Draft-04、Draft-06、および Draft-07 仕様) 用の [JSON スキーマ形式](#) を使用する JSON データ形式、extensions と groups のサポートが含まれない Protocol Buffers (Protobuf) バージョンの proto2 と proto3、および Java 言語の利用などをサポートしており、今後は他のデータ形式や言語にも対応していきます。サポートされている機能には、互換性、メタデータを介したスキーマのソーシング、スキーマの自動登録、IAM 互換性が含まれます。さらにオプションで ZLIB 圧縮を使用することで、ストレージとデータ転送の削減が行えます。AWS Glue Schema Registry はサーバーレスで使用は無料です。

プロデューサーとコンシューマーの間のデータ形式契約としてスキーマを使用すると、データガバナンスおよびデータ品質の向上につながります。さらにデータのコンシューマーは、上流で行われる互換性に関する変更に対する復旧力を得ることができます。

Schema Registry を使用すると、異なるシステムで、シリアル化と非シリアル化用のスキーマを共有できます。例えば、データのプロデューサーとコンシューマーを考えてみます。データの公開時、プロデューサーはスキーマを把握しています。Schema Registry は、Amazon MSK や Apache Kafka などの特定のシステム用に、シリアル化と非シリアル化の機能を提供します。

詳しくは、[「Schema Registry のしくみ」](#) を参照してください。

トピック

- [スキーマ](#)
- [レジストリ](#)
- [スキーマのバージョニングと互換性](#)
- [オープンソース Serde ライブラリ](#)
- [Schema Registry のクォータ](#)
- [Schema Registry のしくみ](#)
- [Schema Registry の使用開始](#)
- [AWS Glue Schema Registry との統合](#)
- [サードパーティ製のスキーマレジストリから AWS Glue Schema Registry への移行](#)

スキーマ

スキーマは、データレコードの構造と形式を定義します。スキーマは、信頼性の高いデータの公開、利用、または保存のための仕様をバージョニングしたものです。

この例の Avro のスキーマでは、形式と構造はレイアウトとフィールド名によって定義され、フィールド名の形式はデータ型 (例:string、int) により定義されています

```
{
  "type": "record",
  "namespace": "ABC_Organization",
  "name": "Employee",
  "fields": [
    {
      "name": "Name",
      "type": "string"
    },
    {
      "name": "Age",
      "type": "int"
    },
    {
      "name": "address",
      "type": {
        "type": "record",
        "name": "addressRecord",
        "fields": [
```

```
    {
      "name": "street",
      "type": "string"
    },
    {
      "name": "zipcode",
      "type": "int"
    }
  ]
}
]
```

この例で使用している、JSON の JSON Schema Draft-07 では、形式を [JSON Schema organization](#) が定義しています。

```
{
  "$id": "https://example.com/person.schema.json",
  "$schema": "http://json-schema.org/draft-07/schema#",
  "title": "Person",
  "type": "object",
  "properties": {
    "firstName": {
      "type": "string",
      "description": "The person's first name."
    },
    "lastName": {
      "type": "string",
      "description": "The person's last name."
    },
    "age": {
      "description": "Age in years which must be equal to or greater than zero.",
      "type": "integer",
      "minimum": 0
    }
  }
}
```

この例の Protobuf では、データ形式を [Protocol Buffers 1言語のバージョン 2 \(proto2\)](#) で定義しています。

```
syntax = "proto2";
```

```
package tutorial;

option java_multiple_files = true;
option java_package = "com.example.tutorial.protos";
option java_outer_classname = "AddressBookProtos";

message Person {
  optional string name = 1;
  optional int32 id = 2;
  optional string email = 3;

  enum PhoneType {
    MOBILE = 0;
    HOME = 1;
    WORK = 2;
  }

  message PhoneNumber {
    optional string number = 1;
    optional PhoneType type = 2 [default = HOME];
  }

  repeated PhoneNumber phones = 4;
}

message AddressBook {
  repeated Person people = 1;
}
```

レジストリ

レジストリはスキーマの論理コンテナです。レジストリを使用すると、スキーマを整理したり、アプリケーションのアクセス制御を管理したりできます。レジストリは Amazon リソースネーム (ARN) を持っており、これにより、レジストリ内のスキーマ操作に対するさまざまなアクセス許可の整理と設定を可能にしています。

レジストリはデフォルトのまま使用することも可能ですが、必要な数の新しいレジストリを作成することもできます。

AWS Glue Schema Registry での階層

- RegistryName: [string]

- RegistryArn: [AWS ARN]
 - CreatedTime: [timestamp]
 - UpdatedTime: [timestamp]
- SchemaName: [string]
 - SchemaArn: [AWS ARN]
 - DataFormat: [Avro、Json、Protobuf]
 - Compatibility: [例、BACKWARD、BACKWARD_ALL、FORWARD、FORWARD_ALL、FULL、FULL_ALL、NONE、DISABLED]
 - Status: [例、PENDING、AVAILABLE、DELETING]
 - SchemaCheckpoint: [integer]
 - CreatedTime: [timestamp]
 - UpdatedTime: [timestamp]
- SchemaVersion: [string]
 - SchemaVersionNumber: [integer]
 - Status: [例、PENDING、AVAILABLE、DELETING、FAILURE]
 - SchemaDefinition: [string, Value: JSON]
 - CreatedTime: [timestamp]
- SchemaVersionMetadata: [list]
 - MetadataKey: [string]
 - MetadataInfo
 - MetadataValue: [string]
 - CreatedTime: [timestamp]

スキーマのバージョンニングと互換性

各スキーマは、複数のバージョン作成することができます。バージョンニングは、スキーマに適用される互換性ルールによって制御されます。新しいスキーマバージョンの登録に関するリクエストについては、この規則に対するチェックが Schema Registry により行われ後で処理が続行されます。

チェックポイントとしてマークされたスキーマバージョンは、スキーマの新しいバージョンを登録する際の互換性判断のために使用されます。スキーマが最初に作成されると、最初のバージョンがデフォルトのチェックポイントになります。スキーマのバージョンが増加した場合には、CLI/SDKにより、一連の制約に従う UpdateSchema API を使用しながら、チェックポイントに使用するスキーマのバージョンを変更できます。デフォルトでは、コンソールでスキーマの定義または互換モードを編集することで、チェックポイントが最新バージョンに変更されます。

互換モードを使用すると、スキーマを時間の経過とともにどのように増加させるか、あるいは増加させないかを制御できます。これらのモードは、データのプロデューサおよびコンシューマである、アプリケーション間の契約を形成します。スキーマの新しいバージョンがレジストリに送信されると、そのスキーマ名に適用されている互換性規則を使用して、新しいバージョンの受け付けが可能かどうか判断されます。互換モードには、NONE、DISABLED、BACKWARD、BACKWARD_ALL、FORWARD、FORWARD_ALL、FULL、FULLの8種類があります

Avro データ形式では、フィールドはオプションまたは必須の両方の場合があります。フィールドがオプションの場合、Type には null が含まれます。必須フィールドでは、Type は null に設定されません。

Protobuf のデータ形式の場合、proto2 構文ではフィールドを (繰り返しを含み) オプションまたは必須にすることができます。一方、proto3 構文ではすべてのフィールドが (繰り返しを含み) オプションとなります。すべての互換性ルールは、Protocol Buffers の仕様と、[Google の Protocol Buffers ドキュメント](#)を詳細に把握した上で定義されています。

- NONE: 互換モードは適用されません。このモードは、開発向けのシナリオか、スキーマに適用する互換モードがわからない場合に使用できます。新しいバージョンが追加されたとき、互換性チェックを受けずに受け入れられます。
- DISABLED: 互換性についてこの選択をすると、特定のスキーマのバージョン管理が不要になります。新しいバージョンを追加することはできません。
- BACKWARD: これは推奨の互換性モードで、コンシューマは、現在および過去のスキーマバージョンの両方からの読み取りが可能です。このモードを使用すると、フィールドの削除や任意のフィールドの追加を行う際に、以前のスキーマバージョンとの互換性をチェックできます。アプリケーションが最新のスキーマ用に作成された場合などが、BACKWARD の一般的なユースケースです。

AVRO

例えば、名前 (必須)、姓 (必須)、E メール (必須)、電話番号 (オプション) で定義されているスキーマがあるとします。

次のスキーマバージョンで、必須である E メールフィールドが削除されたとしても、これは正常に登録されます。後方互換性のために、コンシューマは、現在および 1 つ前のスキーマバージョンを読み取れることが必要です。古いメッセージの余分な E メールフィールドが無視されるため、コンシューマは新しいスキーマを読み取ることができます。

仮に、郵便番号などのフィールドを必須として追加する新しいスキーマバージョンが提案された場合、BACKWARD の互換モードでは正常に登録されません。新しいバージョンのコンシューマは、必須な郵便番号フィールドを含まない、スキーマが変更される前の古いメッセージを読み取ることができないためです。ただし、新しいスキーマで郵便番号フィールドがオプションとして設定されていれば、提案されたバージョンは正常に登録されます。郵便番号フィールドはオプションなので、コンシューマはこのフィールドなしでも古いスキーマを読み取れるためです。

JSON

例えば、あるスキーマバージョンが、名 (オプション)、姓 (オプション)、E メール (オプション)、電話番号 (オプション) で定義されているとします。

仮に、次に続くスキーマバージョンで電話番号プロパティがオプションとして追加された場合でも、元のスキーマバージョンで `additionalProperties` フィールドを `false` に設定しプロパティの追加を禁止しているのであれば、この新しいバージョンは正しく登録されます。後方互換性のために、コンシューマは、現在および 1 つ前のスキーマバージョンを読み取れることが必要です。コンシューマは、電話番号プロパティが存在しない元のスキーマで生成されたデータを読み取ることができます。

BACKWARD の互換モードで、オプションの電話番号プロパティを追加する新しいスキーマバージョンが提案された際に、元のスキーマバージョンで `additionalProperties` フィールドを `true` に設定し任意のプロパティの追加を許可していると、このバージョンは正しく登録されません。新しいバージョンのコンシューマは、スキーマが変更される前の古いメッセージを読み取ることができません。これは、別の型 (数ではなく文字列など) で記述された電話番号プロパティデータを、読み取ることができないためです。

PROTOBUF

例えば、first name (必須)、last name (必須)、email (必須) phone number (オプション) フィールドを持つ、proto2 構文の Message Person で定義されているスキーマバージョンを考えます。

AVRO の場合と同様に、次に続くスキーマバージョンで、必須である email フィールドが削除されたとしても、これは正常に登録されます。後方互換性のために、コンシューマは、現在および 1 つ前のスキーマバージョンを読み取ることが必要です。古いメッセージでは、余分な email フィールドが無視されるため、コンシューマは新しいスキーマを読み取ることができます。

BACKWARD の互換モードでは、仮に zip code などのフィールドを必須として追加する新しいスキーマバージョンが提案された場合、これは正常に登録されません。スキーマが変更される前の古いメッセージには、必須な zip code フィールドが含まれていないため、新しいバージョンのコンシューマは、そのメッセージを読み取ることができないためです。ただし、新しいスキーマで zip code フィールドがオプションとして設定されていれば、提案されたバージョンは正常に登録されます。zip code フィールドはオプションなので、コンシューマはこのフィールドなしでも古いスキーマを読み取れます。

gRPC ユースケースの場合、新しい RPC サービスまたは RPC メソッドを追加することで、後方互換性が変更されます。例えば、2 つの RPC メソッド Foo および Bar を使用する RPC サービス MyService で定義されている、スキーマバージョンがあるとします。

次に続くスキーマバージョンで、Baz という新しい RPC メソッドが追加された場合、これは正常に登録されます。新しく追加された RPC メソッド Baz はオプションのため、コンシューマは元のスキーマで生成されたデータを、BACKWARD 互換性に従って読み取ることができます。

提案された新しいスキーマバージョンで、既存の RPC メソッド Foo を削除する場合、BACKWARD 互換性では正常に登録されません。新しいバージョンのコンシューマは、スキーマが変更される前の古いメッセージを読み取ることができません。これは、gRPC アプリケーション内に RPC メソッド Foo が存在せず、データを理解して読み取れないためです。

- BACKWARD_ALL: この互換モードでは、コンシューマは現在および過去のすべてのスキーマバージョンの読み取りが可能です。このモードを使用すると、フィールドを削除したり任意のフィールドを追加する際に、以前のすべてのスキーマバージョンとの互換性をチェックできます。
- FORWARD: この互換モードにより、コンシューマは現在とその次のスキーマバージョンの両方を読み取ることができますが、それより後に続くバージョンを必ずしも読み取れる訳ではありません。このオプションを使用すると、フィールドを追加したり、任意のフィールドを削除したりする際に、最新のスキーマバージョンとの互換性をチェックできます。過去のスキーマ用に作成された

アプリケーションが、より新しいスキーマを処理できる必要がある場合などが、FORWARD の一般的なユースケースです。

AVRO

例として、名前 (必須)、姓 (必須)、E メール (オプション) で定義されているスキーマバージョンを考えます。

この場合、電話番号などのフィールドを必須として追加する新しいスキーマバージョンは正常に登録されます。FORWARD 互換モードでは、新しいスキーマで生成されたデータを、コンシューマーが以前のバージョンを使用して読み取れることが必要です。

必須である名 (ファーストネーム) のフィールドを削除するスキーマバージョンが提案されとしても、FORWARD 互換モードでは正常に登録されません。名の必須フィールドがないので、以前のバージョンのコンシューマーが提案されたスキーマを読み取ることができなくなるためです。しかし、名 (ファーストネーム) のフィールドが本来オプションであれば、オプションの名のフィールドを持たない新しいスキーマに基づいたデータをコンシューマーが読み取りできるため、提案された新しいスキーマは正常に登録されます。

JSON

例えば、あるスキーマバージョンが、名 (オプション)、姓 (オプション)、E メール (オプション)、電話番号 (オプション) で定義されているとします。

オプションの電話番号プロパティを削除する新しいスキーマバージョンの場合、このスキーマバージョンで `additionalProperties` フィールドを `false` に設定しプロパティの追加を禁止しているのであれば、正しく登録されます。FORWARD 互換モードでは、新しいスキーマで生成されたデータを、コンシューマーが以前のバージョンを使用して読み取れることが必要です。

オプションの電話番号プロパティを削除するスキーマバージョンが提案された際に、この新しいスキーマバージョンで `additionalProperties` フィールドを `true` に設定し任意のプロパティの追加を許可していると、FORWARD の互換モードでは正しく登録されません。提案されたスキーマが、異なる型 (数ではなく文字列など) の電話番号プロパティを持つ場合には、以前のバージョンを使用するコンシューマーで、このスキーマの読み取りが不可能になるためです。

PROTOBUF

例えば、`first name` (必須)、`last name` (必須)、`email` (オプション) のフィールドを持つ、`proto2` 構文の `Message Person` で定義されているスキーマバージョンを考えます。

AVRO の場合と同様に、phone number などの必須フィールドを追加する新しいスキーマバージョンの場合、これは正常に登録されます。FORWARD 互換モードでは、新しいスキーマで生成されたデータを、コンシューマーが以前のバージョンを使用して読み取れることが必要です。

必須である名 first name フィールドを削除するスキーマバージョンが提案されとしても、FORWARD 互換モードでは正常に登録されません。必須である first name フィールドがないので、以前のバージョンのコンシューマーは、提案されたスキーマを読み取ることができなくなるためです。しかし、first name フィールドが本来オプションであれば、新しいスキーマに first name フィールドがないとしてもこれはオプションであるため、コンシューマーはこのスキーマに基づいたデータを読み取ることが可能であり、したがって提案された新しいスキーマは正常に登録されます。

gRPC ユースケースの場合、RPC サービスまたは RPC メソッドを削除すると、前方互換性が変更されます。例えば、2 つの RPC メソッド Foo および Bar を使用する RPC サービス MyService で定義されている、スキーマバージョンがあるとします。

次に続くスキーマバージョンが、Foo という名前の既存の RPC メソッドを削除する場合も、これは FORWARD 互換性に従って正常に登録されます。コンシューマーは、新しいスキーマで生成されたデータを、以前のバージョンを使用して読み取ることが可能です。RPC メソッド Baz を追加するスキーマバージョンが提案され場合、これは FORWARD 互換モードでは正常に登録されません。RPC メソッド Baz がないので、以前のバージョンのコンシューマーは、この提案されたスキーマを読み取ることができないためです。

- FORWARD_ALL: この互換モードでは、任意の新しい登録済みスキーマのプロデューサーによって書き込まれたデータを読み取ることを、コンシューマーに許可します。この選択は、フィールドを追加したり、オプションフィールドを削除したり、以前のすべてのスキーマバージョンとの互換性を確認したりする必要がある場合に使用できます。
- FULL: この互換モードにより、コンシューマーは、1 つ前または次のバージョンのスキーマを使用するプロデューサーによって書き込まれたデータを読み取ることができますが、それ以前またはそれ以降のバージョンについては読み取れません。このモードでは、任意のフィールドを追加または削除する際に、最新のスキーマバージョンとの互換性をチェックできます。
- FULL_ALL: この互換モードにより、コンシューマーは、過去の任意のスキーマバージョンを使用するプロデューサーによって書き込まれたデータを読み取ることができます。このモードは、オプションのフィールドを追加または削除する際に、過去のすべてのスキーマバージョンとの互換性をチェックするために使用できます。

オープンソース Serde ライブラリ

AWS では、データをシリアル化および非シリアル化するためのフレームワークとして、オープンソースの Serde ライブラリが用意されています。このオープンソースなライブラリ設計により、一般的なオープンソースのアプリケーションやフレームワークが、それぞれのプロジェクトでこれらのライブラリをサポートできるようにしています。

Serde ライブラリの仕組みの詳細については、「[Schema Registry のしくみ](#)」を参照してください。

Schema Registry のクォータ

AWS のクォータ (制限とも呼ばれます) は、AWS アカウントのリソース、アクション、および制限の最大値です。以下に、AWS Glue の Schema Registry におけるソフト制限を示します。

スキーマバージョンのメタデータでのキー値ペアの数。

AWS リージョンごとの各 SchemaVersion において、最大 10 個のキーと値のペアを使用できます。

メタデータでのキーと値のペアは、[QuerySchemaVersionMetadata アクション \(Python: query_schema_version_metadata\)](#) または [PutSchemaVersionMetadata アクション \(Python: put_schema_version_metadata\)](#) API を使用して表示または設定できます。

AWS Glue の Schema Registry でのハード制限は次のとおりです。

レジストリ

このアカウントでは、AWS リージョンごとに最大 100 のレジストリを使用できます。

SchemaVersion

このアカウントでは、AWS リージョンごとに最大 10000 のスキーマバージョンを使用できます。

各新しいスキーマは新しいスキーマバージョンを作成します。したがって、各スキーマに 1 つのバージョンしかない場合、リージョンごとの各アカウントには理論的に最大 10000 のスキーマを使用できます。

スキーマペイロード

スキーマペイロードには、170 KB のサイズ制限があります。

Schema Registry のしくみ

このセクションでは、Schema Registry でのシリアル化および非シリアル化プロセスの仕組みについて説明します。

1. スキーマの登録: スキーマがレジストリにまだ存在しない場合、スキーマを保存先の名前と同じスキーマ名で登録することができます (test_topic、test_stream、prod_firehose など)。あるいは、プロデューサーがスキーマのカスタム名を指定することも可能です。プロデューサーは、source: msk_kafka_topic_A などのメタデータとしてスキーマにキーと値のペアを追加したり、スキーマ作成時に AWS タグをそのスキーマに追加したりできます。スキーマが登録されると、Schema Registry はスキーマのバージョン ID をシリアライザに返します。スキーマが存在するものの、シリアライザが存在しない新しいバージョンを使用している場合、それを新しいバージョンとして登録する前に、Schema Registry はスキーマ参照で互換性ルールを確認して、新しいバージョンに関する互換性を確保します。

スキーマの登録には、手動登録と自動登録の 2 つの方法があります。スキーマを手動で登録するには、AWS Glue コンソールまたは CLI/SDK を使用します。

シリアライザ設定で自動登録が有効化されている場合は、スキーマの自動登録が実行されます。プロデューサー設定で REGISTRY_NAME が指定されていない場合、自動登録はデフォルトのレジストリ (default-registry) の下に新しいスキーマバージョンを登録します。自動登録プロパティの指定については、「[SerDe ライブラリのインストール](#)」を参照してください。

2. シリアライザがスキーマに関してデータレコードを検証: データを生成するアプリケーションがそのスキーマを登録した場合、Schema Registry のシリアライザは、アプリケーションによって生成されるレコードが、登録されたスキーマに一致するフィールドとデータ型で構造化されていることを検証します。レコードのスキーマが登録されたスキーマと一致しない場合、シリアライザは例外を出力し、アプリケーションは、そのレコードを送信先にデリバリーすることに失敗します。

スキーマが存在せず、プロデューサー設定を介してスキーマ名が指定されてもいない場合は、スキーマがトピック名 (Apache Kafka または Amazon MSK の場合) あるいはストリーム名 (Kinesis Data Streams の場合) と同じ名前で作成されます。

すべてのレコードは、スキーマ定義とデータを含みます。スキーマ定義は、Schema Registry 内の既存のスキーマおよびバージョンに関して照会されます。

デフォルトでは、プロデューサーは、登録済みスキーマのスキーマ定義とスキーマバージョン ID をキャッシュします。レコードのスキーマバージョン定義がキャッシュ内に保存されたものと一致

しない場合、プロデューサーは Schema Registry を使用してスキーマの検証を試みます。スキーマのバージョンが有効と認められた場合、そのバージョン ID と定義はプロデューサーにローカルにキャッシュされます。

[SerDe ライブラリのインストール](#) のステップ 3 のように、オプションのプロデューサープロパティにより、デフォルトのキャッシュ期間 (24 時間) を調整することが可能です。

3. レコードのシリアル化と配信: レコードがスキーマに準拠している場合、シリアライザは各レコードをスキーマのバージョン ID で装飾し、選択したデータ形式 (AVRO、JSON、Protobuf 他の形式は追加予定) に基づいてシリアル化した上で、そのレコードを (プロデューサーのオプション設定により) 圧縮し送信先に送ります。
4. コンシューマーがデータを非シリアル化: このデータを読み取るコンシューマーは、レコードペイロードからのスキーマバージョン ID を解析する、Schema Registry のデシリアライザ用ライブラリを使用します。
5. デシリアライザが Schema Registry からのスキーマをリクエスト: デシリアライザが特定のスキーマバージョン ID を持つレコードを初めて認識すると、そのデシリアライザは、スキーマバージョン ID を使用してスキーマレジストリからスキーマを要求します。その後、そのスキーマをコンシューマー上でローカルにキャッシュします。Schema Registry でレコードの非シリアル化が不可能な場合、コンシューマーはレコードからのデータをログに記録し、処理を続行するかアプリケーションを停止できます。
6. デシリアライザがスキーマを使用してレコードを非シリアル化: デシリアライザは、Schema Registry からスキーマバージョン ID を取得する際に (プロデューサーによって送信されたレコードが圧縮されている場合は) レコードを解凍し、スキーマを使用してレコードを非シリアル化します。これで、アプリケーションはレコードを処理できます。

Note

暗号化: クライアントは、HTTPS 上での TLS 暗号化を使用して転送中のデータを暗号化する API 呼び出しを介して Schema Registry と通信します。Schema Registry に保存されているスキーマは、サービス管理の AWS Key Management Service (AWS KMS) キーにより常に暗号化されます。

Note

ユーザー認証: Schema Registry は、アイデンティティベースの IAM ポリシーをサポートします。

Schema Registry の使用開始

以下のセクションでは、Schema Registry のセットアップと使用に関する概要とチュートリアルを示します。Schema Registry の概念と各コンポーネントについては、「[AWS Glue スキーマレジストリ](#)」を参照してください。

トピック

- [SerDe ライブラリのインストール](#)
- [AWS Glue Schema Registry API のために AWS CLI を使用する](#)
- [レジストリを作成する](#)
- [JSON の特定のレコード \(JAVA POJO\) 処理する](#)
- [スキーマの作成](#)
- [スキーマまたはレジストリの更新](#)
- [スキーマまたはレジストリの削除](#)
- [シリアライザ用の IAM の例](#)
- [デシリアライザー用の IAM の例](#)
- [AWS PrivateLink を使用したプライベート接続](#)
- [Amazon CloudWatch メトリクスへのアクセス](#)
- [Schema Registry の AWS CloudFormation テンプレート例](#)

SerDe ライブラリのインストール

Note

前提条件: 次のステップを実行する前に、Amazon Managed Streaming for Apache Kafka (Amazon MSK) または Apache Kafka クラスターを起動しておく必要があります。使用するプロデューサーとコンシューマーは、Java 8 以上で実行する必要があります。

SerDe ライブラリは、データのシリアル化と非シリアル化のためのフレームワークを提供します。

データを生成するアプリケーション (総称してシリアライザ) 用の、オープンソースのシリアライザをインストールします。シリアライザは、シリアル化、圧縮、および Schema Registry とのやり取りを処理します。シリアライザは、Schema Registry 対応の送信先 (Amazon MSK など) に書き込まれるレコードから、スキーマを自動的に抽出します。同様に、データを利用するアプリケーションには、オープンソースのデシリアライザをインストールします。

プロデューサとコンシューマにライブラリをインストールするには

1. プロデューサとコンシューマ両方の pom.xml ファイルの中で、以下のコードにより依存関係を追加します。

```
<dependency>
  <groupId>software.amazon.glue</groupId>
  <artifactId>schema-registry-serde</artifactId>
  <version>1.1.5</version>
</dependency>
```

または、[AWS Glue Schema Registry GitHub リポジトリ](#) からクローンを作成することもできます。

2. 次の必須プロパティを使用してプロデューサをセットアップします。

```
props.put(ProducerConfig.KEY_SERIALIZER_CLASS_CONFIG,
StringSerializer.class.getName()); // Can replace StringSerializer.class.getName()
with any other key serializer that you may use
props.put(ProducerConfig.VALUE_SERIALIZER_CLASS_CONFIG,
GlueSchemaRegistryKafkaSerializer.class.getName());
props.put(AWSSchemaRegistryConstants.AWS_REGION, "us-east-2");
properties.put(AWSSchemaRegistryConstants.DATA_FORMAT, "JSON"); // OR "AVRO"
```

既存のスキーマがない場合は、自動登録を有効にします (次のステップ)。適用できるスキーマがある場合は、「my-schema」をそのスキーマ名に置き換えます。また、スキーマの自動登録が無効になっている場合は、「registry-name」を指定する必要があります。スキーマが「default-registry」の下に作成されている場合は、レジストリ名を省略できます。

3. (オプション) これらのオプションのプロデューサープロパティのいずれかを設定します。プロパティの詳細については、[ReadMe ファイル](#)を参照してください。

```
props.put(AWSSchemaRegistryConstants.SCHEMA_AUTO_REGISTRATION_SETTING, "true"); // If
not passed, uses "false"
props.put(AWSSchemaRegistryConstants.SCHEMA_NAME, "my-schema"); // If not passed,
uses transport name (topic name in case of Kafka, or stream name in case of Kinesis
Data Streams)
props.put(AWSSchemaRegistryConstants.REGISTRY_NAME, "my-registry"); // If not passed,
uses "default-registry"
props.put(AWSSchemaRegistryConstants.CACHE_TIME_TO_LIVE_MILLIS, "86400000"); // If
not passed, uses 86400000 (24 Hours)
props.put(AWSSchemaRegistryConstants.CACHE_SIZE, "10"); // default value is 200
props.put(AWSSchemaRegistryConstants.COMPATIBILITY_SETTING, Compatibility.FULL); //
Pass a compatibility mode. If not passed, uses Compatibility.BACKWARD
props.put(AWSSchemaRegistryConstants.DESCRPTION, "This registry is used for several
purposes."); // If not passed, constructs a description
props.put(AWSSchemaRegistryConstants.COMPRESSION_TYPE,
AWSSchemaRegistryConstants.COMPRESSION.ZLIB); // If not passed, records are sent
uncompressed
```

自動登録では、スキーマのバージョンがデフォルトのレジストリ (default-registry) に登録されま
す。SCHEMA_NAME が前のステップで指定されていない場合、トピック名は SCHEMA_NAME とし
て推定されます。

互換モードの詳細については、「[スキーマのバージョンングと互換性](#)」を参照してください。

4. 以下の必須プロパティを使用してコンシューマを設定します。

```
props.put(ConsumerConfig.KEY_DESERIALIZER_CLASS_CONFIG,
StringDeserializer.class.getName());
props.put(ConsumerConfig.VALUE_DESERIALIZER_CLASS_CONFIG,
GlueSchemaRegistryKafkaDeserializer.class.getName());
props.put(AWSSchemaRegistryConstants.AWS_REGION, "us-east-2"); // Pass an AWS #####
props.put(AWSSchemaRegistryConstants.AVRO_RECORD_TYPE,
AvroRecordType.GENERIC_RECORD.getName()); // Only required for AVRO data format
```

5. (オプション) これらのオプションのコンシューマプロパティを設定します。プロパティの詳細に ついては、[ReadMe ファイル](#)を参照してください。

```
properties.put(AWSSchemaRegistryConstants.CACHE_TIME_TO_LIVE_MILLIS, "86400000"); //
If not passed, uses 86400000
props.put(AWSSchemaRegistryConstants.CACHE_SIZE, "10"); // default value is 200
```

```
props.put(AWSSchemaRegistryConstants.SECONDARY_DESERIALIZER,  
"com.amazonaws.services.schemaregistry.deserializers.external.ThirdPartyDeserializer"); //  
For migration fall back scenario
```

AWS Glue Schema Registry API のために AWS CLI を使用する

AWS CLI で AWS Glue Schema Registry API を使用するには、最新バージョンの AWS CLI 使用する必要が有ります。

レジストリを作成する

AWS Glue API または AWS Glue コンソールを使用して、デフォルトのレジストリを使用することも、必要な数の新しいレジストリを作成することもできます。

AWS Glue API

ここでの手順により、AWS Glue API を使用しながら対象のタスクを実行できます。

新しいレジストリを追加するには、[CreateRegistry アクション \(Python: create_registry\)](#) API を使用します。RegistryName では、作成するレジストリの名前を指定します。この文字数は最大 255 まで、文字、数字、ハイフン、アンダースコア、ドル記号、およびハッシュ記号のみ使用できます。

2,048 バイト以下で「[URI アドレスの複数行の文字列パターン](#)」に一致する文字列として Description を指定します。

オプションで、キーと値のペアのマップ配列として、1 つ以上の Tags をレジストリに指定します。

```
aws glue create-registry --registry-name registryName1 --description description
```

作成されたレジストリには、Amazon リソースネーム (ARN) が割り当てられます。これは、RegistryArn API 応答により表示することが可能です。この段階でレジストリ作成が完了しているため、そのレジストリのために 1 つ以上のスキーマを作成します。

AWS Glue コンソール

AWS Glue コンソールで新しいレジストリを追加するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。

- ナビゲーションペインの [Data catalog] (データカタログ) で、[Schema registries] (スキーマレジストリ) をクリックします。
- [Add registry] (レジストリを追加) をクリックします。
- [Registry name] (レジストリ名) に、文字、数字、ハイフン、アンダースコアを含むレジストリの名前を入力します。この名前は変更できません。
- レジストリの [Description] (説明) を入力します (オプション)。
- オプションで、1 つ以上のタグをレジストリに適用します。[Add new tag] (新しいタグを追加) を選択し、[Tag key] (タグキー) とオプションの [Tag value] (タグ値) を指定します。
- [Add registry] (レジストリを追加) をクリックします。

AWS Glue

Data catalog

Databases

Tables

Connections

Crawlers

Classifiers

Schema registries

Schemas

Settings

ETL

AWS Glue Studio

New

Workflows

Jobs

ML Transforms

Triggers

Dev endpoints

Notebooks

Schema registries > Add registry

Add a new schema registry

Add a schema registry to store one or multiple new related schemas.

Registry name
Name can't be changed post creation.

Enter a registry name...

Only letters (A-Z), numbers (0-9), hyphens (-), underscores (_), dollar signs (\$), or hash marks (#) allowed. 255 characters maximum.

Description - optional

Enter a registry description...

2048 characters maximum.

Registry tags - optional
No tags defined.

Add new tag

You can add up to 50 more tags.

Cancel Add registry

作成されたレジストリには、Amazon リソースネーム (ARN) が割り当てられます。これは、[Schema registries] (スキーマレジストリ) のリストから選択して表示することが可能です。この段階でレジストリ作成が完了しているため、そのレジストリのために 1 つ以上のスキーマを作成します。

JSON の特定のレコード (JAVA POJO) 処理する

従来の単純な Java オブジェクト (POJO) を使用して、オブジェクトをレコードとして渡すことができます。これは、AVRO における特定のレコードの概念と類似しています。[mbknor-jackson-jsonschema](https://github.com/mbknor-jackson-jsonschema) により、渡された POJO の JSON スキーマを生成できます。また、このライブラリでは、JSON スキーマに追加情報を挿入することもできます。

AWS Glue Schema Registry ライブラリは、スキーマに注入された「className」フィールドを使用して、完全に分類されたクラス名を提供します。「className」フィールドは、そのクラスのオブジェクト内での非シリアル化のために、デシリアライザによって使用されます。

Example class :

```
@JsonSchemaDescription("This is a car")
@JsonSchemaTitle("Simple Car Schema")
@Builder
@AllArgsConstructor
@EqualsAndHashCode
// Fully qualified class name to be added to an additionally injected property
// called className for deserializer to determine which class to deserialize
// the bytes into
@JsonSchemaInject(
    strings = {@JsonSchemaString(path = "className",
        value =
            "com.amazonaws.services.schemaregistry.integrationtests.generators.Car")}
)
// List of annotations to help infer JSON Schema are defined by https://github.com/
mbknor/mbknor-jackson-jsonSchema
public class Car {
    @JsonProperty(required = true)
    private String make;

    @JsonProperty(required = true)
    private String model;

    @JsonSchemaDefault("true")
    @JsonProperty
    public boolean used;

    @JsonSchemaInject(ints = {@JsonSchemaInt(path = "multipleOf", value = 1000)})
    @Max(200000)
    @JsonProperty
    private int miles;
```

```
@Min(2000)
@JsonProperty
private int year;

@JsonProperty
private Date purchaseDate;

@JsonProperty
@JsonFormat(shape = JsonFormat.Shape.NUMBER)
private Date listedDate;

@JsonProperty
private String[] owners;

@JsonProperty
private Collection<Float> serviceChecks;

// Empty constructor is required by Jackson to deserialize bytes
// into an Object of this class
public Car() {}
}
```

スキーマの作成

スキーマは、AWS Glue API または AWS Glue コンソールを使用して作成できます。

AWS Glue API

ここでの手順により、AWS Glue API を使用しながら対象のタスクを実行できます。

新しいスキーマを追加するには [CreateSchema アクション \(Python: create_schema\)](#) API を使用します。

RegistryId 構造体を使用して、スキーマのレジストリを指定します。または、RegistryId を省略すると、デフォルトのレジストリが使用されます。

文字、数字、ハイフン、アンダースコアを使用し、DataFormat として **AVRO** または **JSON** を設定しながら SchemaName を指定します。一度スキーマに設定された DataFormat は変更できません。

Compatibility モードを指定する。

- Backward (推奨) – コンシューマは、現在のバージョンと 1 つ前のバージョンの両方を読み取ることができます。
- Backward all – コンシューマは、現在のバージョンと過去のすべてのバージョンを読み取ることができます。
- Forward – コンシューマは、現在のバージョンと次に続くバージョンの両方を読み取ることができます。
- Forward all – コンシューマは、現在のバージョンと後続のすべてのバージョンの両方を読み取ることができます。
- Full – Backward all と Forward all を組み合わせたモードです。
- Full all – Backward all と Forward all を組み合わせたモードです。
- None – 互換性チェックは実行されません。
- Disabled – このスキーマのバージョニングを抑制します。

オプションで、スキーマに Tags を指定します。

SchemaDefinition を指定することで、スキーマのデータ形式を Avro、JSON、もしくは Protobuf として定義します。例を参照してください。

Avro データ形式の場合:

```
aws glue create-schema --registry-id RegistryName="registryName1" --schema-name
testschema --compatibility NONE --data-format AVRO --schema-definition "{\"type\":
\\\"record\\\", \\\"name\\\": \\\"r1\\\", \\\"fields\\\": [ {\\\"name\\\": \\\"f1\\\", \\\"type\\\": \\\"int\\\"},
{\\\"name\\\": \\\"f2\\\", \\\"type\\\": \\\"string\\\"} ]}"
```

```
aws glue create-schema --registry-id RegistryArn="arn:aws:glue:us-
east-2:901234567890:registry/registryName1" --schema-name testschema --compatibility
NONE --data-format AVRO --schema-definition "{\"type\": \\\"record\\\", \\\"name\\\": \\\"r1\\\",
\\\"fields\\\": [ {\\\"name\\\": \\\"f1\\\", \\\"type\\\": \\\"int\\\"}, {\\\"name\\\": \\\"f2\\\", \\\"type\\\":
\\\"string\\\"} ]}"
```

JSON データ形式の場合:

```
aws glue create-schema --registry-id RegistryName="registryName" --schema-name
testSchemaJson --compatibility NONE --data-format JSON --schema-definition "{\"$schema
\\\": \\\"http://json-schema.org/draft-07/schema#\\\", \\\"type\\\": \\\"object\\\", \\\"properties\\\":
{\\\"f1\\\": {\\\"type\\\": \\\"string\\\"}}}"
```

```
aws glue create-schema --registry-id RegistryArn="arn:aws:glue:us-east-2:901234567890:registry/registryName" --schema-name testSchemaJson --compatibility NONE --data-format JSON --schema-definition "{\"$schema\": \"http://json-schema.org/draft-07/schema#\", \"type\": \"object\", \"properties\": {\"f1\": {\"type\": \"string\"}}}"
```

Protobuf データ形式の場合:

```
aws glue create-schema --registry-id RegistryName="registryName" --schema-name testSchemaProtobuf --compatibility NONE --data-format PROTOBUF --schema-definition "syntax = \"proto2\"; package org.test; message Basic { optional int32 basic = 1;}"
```

```
aws glue create-schema --registry-id RegistryArn="arn:aws:glue:us-east-2:901234567890:registry/registryName" --schema-name testSchemaProtobuf --compatibility NONE --data-format PROTOBUF --schema-definition "syntax = \"proto2\"; package org.test; message Basic { optional int32 basic = 1;}"
```

AWS Glue コンソール

AWS Glue コンソールを使用して新しいスキーマを追加するには

1. AWS マネジメントコンソールにサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインの [Data catalog] (データカタログ) で、[Schemas] (スキーマ) をクリックします。
3. [Add schema] (スキーマを追加) をクリックします。
4. [Schema name] (スキーマ名) に、文字、数字、ハイフン、アンダースコア、ドル記号、ハッシュマークで構成された名前を入力します。この名前は変更できません。
5. [Registry] (レジストリ) ドロップダウンメニューから、スキーマの保存先となるレジストリを選択します。作成後は、親レジストリを変更することはできません。
6. [Data format] (データ形式) は、「Apache Avro」または「JSON」のままにしておきます。この形式は、このスキーマのすべてのバージョンに適用されます。
7. [Compatibility mode] (互換モード) をクリックします。
 - Backward (推奨) – レシーバーは、現在のバージョンと 1 つ前のバージョンの両方を読み取ることができます。
 - Backward All – レシーバーは、現在および過去のすべてのバージョンを読み取ることができます。

- Forward – センダーは、現在のバージョンと 1 つ前のバージョンの両方に書き込むことができます。
 - Forward All – センダーは、現在のバージョンと過去のすべてのバージョンに書き込むことができます。
 - Full – Backward と Forward を組み合わせたモードです。
 - Full All – Backward All と Forward All を組み合わせたモードです。
 - None – 互換性チェックは実行されません。
 - Disabled – このスキーマのバージョンングを抑制します。
8. [Description] (説明) に、レジストリのための説明 (オプション) を、最大 250 文字で入力します。

AWS Glue

Data catalog

Databases

Tables

Connections

Crawlers

Classifiers

Schema registries

Schemas

Settings

ETL

AWS Glue Studio

New

Workflows

Jobs

ML Transforms

Triggers

Dev endpoints

Notebooks

Security

Security configurations

Tutorials

Add crawler

Explore table

Add job

Resources What's new 

Schemas > Add schema

Add a new schema

Specify your new schema name, properties, and schema definition.

Schema name

Name can't be changed post creation.

Only letters (A-Z), numbers (0-9), hyphens (-), underscores (_), dollar signs (\$), or hash marks (#) allowed. 255 characters maximum.

Registry

Parent registry can't be changed post creation.

[Add new registry](#)

Data format

Glue schemas only support Apache Avro for now, which offers the compatibility options below. [Learn more](#) 

Compatibility mode

Compatibility may be changed post creation and affects data senders and/or receivers.

**Backward compatibility** [Learn more](#) 

This compatibility choice allows consumers to read both the current and the previous schema version. This means that for instance, a new schema version cannot drop data fields or change the type of these fields, so they can't be read by consumers using the previous version.

Description - optional

2048 characters maximum.

9. オプションで、1 つ以上のタグをスキーマに適用します。[Add new tag] (新しいタグを追加) を選択し、[Tag key] (タグキー) とオプションの [Tag value] (タグ値) を指定します。
- 10 [First schema バージョニング] (最初のスキーマバージョン) ボックスに、初期スキーマを入力するか貼り付けます。

Avro 形式については「[Avro データ形式での作業](#)」を参照

JSON 形式については「[JSON データ形式での操作](#)」を参照

11必要に応じて、[Add metadata] (メタデータを追加) をクリックして、スキーマバージョンの注釈付けや分類を行うためのバージョンメタデータを追加します。

12[Create schema and version] (スキーマとバージョンを作成する) をクリックします。

AWS Glue

Data catalog

Databases

Tables

Connections

Crawlers

Classifiers

Schema registries

Schemas

Settings

ETL

AWS Glue Studio

New

Blueprints

Workflows

Jobs

ML Transforms

Triggers

Dev endpoints

Notebooks

Security

Security configurations

Tutorials

Add crawler

Explore table

Add job

Resources

Schema tags - optional

No tags defined.

Add new tag

You can add up to 50 more tags.

First schema version

Please specify the initial definition of your schema below, so that it can be used in your applications or within Amazon Glue. You may change your schema definition by registering new versions at any point later. Please enter Apache Avro schema below. [Learn more](#)

1

Version metadata - optional

No metadata key-value pairs.

Add metadata

You can add 10 more metadata key-value pairs.

Cancel Create schema and version

スキーマが作成され、[Schemas] (スキーマ) の下に一覧表示されます。

Avro データ形式での作業

Avro では、データのシリアル化とデータ交換サービスが利用できます。Avro は、可読性と解釈しやすさのために、データ定義が JSON 形式で格納されます。データ自体はバイナリ形式で保存されます。

Apache Avro スキーマの定義については、「[Apache Avro specification](#)」を参照してください。

JSON データ形式での操作

JSON 形式では、データをシリアル化できます。JSON スキーマ形式の標準は、「[JSON Schema format](#)」で定義されています。

スキーマまたはレジストリの更新

作成したスキーマ、スキーマバージョン、またはレジストリは、編集することができます。

レジストリの更新

レジストリは、AWS Glue API または AWS Glue コンソールを使用して更新できます。既存のレジストリの名前は変更できません。レジストリの説明は編集が可能です。

AWS Glue API

既存のレジストリを更新するには、[UpdateRegistry アクション \(Python: update_registry\)](#) API を使用します。

RegistryId 構造体を使用して、更新するレジストリを指定します。レジストリの変更するには、Description を渡します。

```
aws glue update-registry --description updatedDescription --registry-id
RegistryArn="arn:aws:glue:us-east-2:901234567890:registry/registryName1"
```

AWS Glue コンソール

AWS Glue コンソールを使用してレジストリを更新するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインの [Data catalog] (データカタログ) で、[Schema registries] (スキーマレジストリ) をクリックします。

3. レジストリの一覧から、そのチェックボックスをオンにして、レジストリを選択します。
4. [Action] (アクション) メニューで、[Edit registry] (レジストリの編集) を選択します。

スキーマの更新

スキーマでは、説明または互換性設定の更新が行えます。

既存のスキーマを更新するには、[UpdateSchema アクション \(Python: update_schema\)](#) API を使用します。

SchemaId 構造体を使用して、更新するスキーマを指定します。VersionNumber または Compatibility のいずれかを指定する必要があります。

コード例 11:

```
aws glue update-schema --description testDescription --schema-id
  SchemaName="testSchema1",RegistryName="registryName1" --schema-version-number
  LatestVersion=true --compatibility NONE
```

```
aws glue update-schema --description testDescription --schema-id
  SchemaArn="arn:aws:glue:us-east-2:901234567890:schema/registryName1/testSchema1" --
  schema-version-number LatestVersion=true --compatibility NONE
```

スキーマバージョンの追加。

スキーマバージョンを追加する際は、そのバージョンを比較して、新しいスキーマが受け入れられることを確認する必要があります。

既存のスキーマに新しいバージョンを追加するには、[RegisterSchemaVersion アクション \(Python: register_schema_version\)](#) API を使用します。

SchemaId 構造体を使用して、バージョンを追加するスキーマを指定し、SchemaDefinition によりスキーマを定義します。

コード例 12:

```
aws glue register-schema-version --schema-definition '{"type\": \"record\", \"name\":
  \"r1\", \"fields\": [ {\"name\": \"f1\", \"type\": \"int\"}, {\"name\": \"f2\", \"type
  \": \"string\"} ]}' --schema-id SchemaArn="arn:aws:glue:us-east-1:901234567890:schema/
  registryName/testschema"
```

```
aws glue register-schema-version --schema-definition "{\"type\": \"record\", \"name\": \"r1\", \"fields\": [ {\"name\": \"f1\", \"type\": \"int\"}, {\"name\": \"f2\", \"type\": \"string\"} ]}" --schema-id SchemaName="testschema",RegistryName="testregistry"
```

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインの [Data catalog] (データカタログ) で、[Schemas] (スキーマ) をクリックします。
3. チェックボックスをオンにして、スキーマのリストからスキーマを選択します。
4. チェックボックスをクリックして、リストから 1 つ以上のスキーマを選択します。
5. [Action] (アクションメニュー) から、[Register new version] (新しいバージョンを登録) を選択します。
6. [New version] (新しいバージョン) ボックスに、新しいスキーマを入力または貼り付けます。
7. [Compare with previous version] (過去のバージョンと比較) をクリックして、以前のスキーマのバージョンとの違いを確認します。
8. 必要に応じて、[Add metadata] (メタデータを追加) をクリックして、スキーマバージョンの注釈付けや分類を行うためのバージョンメタデータを追加します。[Key] (キー) および (オプションの) [Value] (値) を入力します。
9. [Register version] (バージョンの登録) をクリックします。

AWS Glue

Data catalog

Databases

Tables

Connections

Crawlers

Classifiers

Schema registries

Schemas

Settings

ETL

AWS Glue Studio

New

Blueprints

Workflows

Jobs

ML Transforms

Triggers

Dev endpoints

Notebooks

Security

Security configurations

Tutorials

Add crawler

Explore table

Add job

Schemas > test-1 > Register version

Register a new schema version

Register version 4 to your schema.

Schema name	test-1
Data format	Apache Avro
Compatibility mode	Backward compatibility
Schema tags	No tags defined.

New Version 4

This is a copy of version 1's schema definition. A schema definition not associated with any existing schema versions must be defined in order to register a new schema version.

```
1  {  
2    "type": "record",  
3    "name": "r0",  
4    "fields": [  
5      {  
6        "name": "f1",  
7        "type": "int"  
8      }  
9    ]  
10 }
```

[Compare with previous version](#)

Version metadata - optional

No metadata key-value pairs.

[Add metadata](#)

You can add 10 more metadata key-value pairs.

[Cancel](#)[Register version](#)

バージョンの一覧の中に、スキーマのバージョンが表示されます。バージョンで互換モードが変更された場合、そのバージョンはチェックポイントとしてマークされます。

スキーマのバージョン比較の例。

[Compare with previous version] (過去のバージョンと比較) をクリックすると、以前のバージョンと新しいバージョンと一緒に表示されます。変更された情報は、次のように強調表示されます。

- 黄色: 変更された情報を示します。

- 緑: 最新バージョンで追加されたコンテンツを示します。
- 赤: 最新バージョンで削除されたコンテンツを示します。

より古いバージョンと比較することも可能です。

Schema test-1 Compatibility Mode Backward compatibility

Version 1 (latest a... ▼) Version 4 (new) ▼

```

1 {
2   "type": "record",
3-  "name": " r 0 ",
4   "fields": [
5     {
6       "name": "f1",
7       "type": "int"
8     }
9   ]
10 }

```

```

1 {
2   "type": "record",
3+  "name": "use r .record ",
4+  "aliases": "userInfo",
5   "fields": [
6     {
7       "name": "f1",
8       "type": "int"
9     }
10  ]
11 }

```

Registered Thu, 01 Oct 2020 17:37:19 GMT Registered -

Metadata - Metadata -

[Close](#)

スキーマまたはレジストリの削除

スキーマ、スキーマバージョン、またはレジストリの削除は永続的な操作であり、元に戻すことはできません。

スキーマの削除

レジストリ内で使用する必要がなくなったスキーマは、AWS Management Console または [DeleteSchema アクション \(Python: delete_schema\)](#) API を使用して削除することができます。。

1 つ以上のスキーマを削除することは永続的なアクションであり、元に戻すことはできません。(1 つあるいは複数の) スキーマが不要になったことを確認します。

レジストリからスキーマを削除するには、SchemaId 構造体により対象のスキーマを特定しながら [DeleteSchema アクション \(Python: delete_schema\)](#) API を呼び出します。

例:

```
aws glue delete-schema --schema-id SchemaArn="arn:aws:glue:us-east-2:901234567890:schema/registryName1/schemaname"
```

```
aws glue delete-schema --schema-id SchemaName="TestSchema6-deleteschemabyname",RegistryName="default-registry"
```

AWS Glue コンソール

AWS Glue コンソールからスキーマを削除するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインの [Data catalog] (データカタログ) で、[Schema registries] (スキーマレジストリ) をクリックします。
3. レジストリのリストから、自分のスキーマを含むレジストリを選択します。
4. チェックボックスをクリックして、リストから 1 つ以上のスキーマを選択します。
5. [Action] (アクション) メニューで、[Delete schema] (スキーマの削除) をクリックします。
6. フィールドに「**Delete**」というテキストを入力して、削除を確定します。
7. [削除] を選択します。

指定した (1 つ以上の) スキーマがレジストリから削除されます。

スキーマバージョンの削除

スキーマはレジストリに蓄積されるので、不要なスキーマバージョンは、AWS Management Console または [DeleteSchemaVersions アクション \(Python: delete_schema_versions\)](#) API を使用して削除できます。1 つ以上のスキーマバージョンを削除することは永続的なアクションであり、元に戻すことはできません。そのスキーマバージョンが不要であることを確認します。

スキーマのバージョンを削除する場合は、以下の制約に注意してください。

- チェックポイントとなっているバージョンを削除することはできません。

- 25 を超えて連続するバージョンの範囲を削除することはできません。
- 最新のスキーマバージョンが保留状態にある場合は、削除は行えません。

SchemaId 構造体を使用してスキーマを指定し、削除するバージョンの範囲を Versions で指定します。バージョンまたはバージョンの範囲の指定の詳細については、「[DeleteRegistry アクション \(Python: delete_registry\)](#)」を参照してください。指定したスキーマバージョンがレジストリから削除されます。

この呼び出しの後に [ListSchemaVersions アクション \(Python: list_schema_versions\)](#) API を呼び出すと、削除されたバージョンのステータスが一覧表示されます。

例:

```
aws glue delete-schema-versions --schema-id
  SchemaName="TestSchema6",RegistryName="default-registry" --versions "1-1"
```

```
aws glue delete-schema-versions --schema-id SchemaArn="arn:aws:glue:us-
east-2:901234567890:schema/default-registry/TestSchema6-NON-Existent" --versions "1-1"
```

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインの [Data catalog] (データカタログ) で、[Schema registries] (スキーマレジストリ) をクリックします。
3. レジストリのリストから、自分のスキーマを含むレジストリを選択します。
4. チェックボックスをクリックして、リストから 1 つ以上のスキーマを選択します。
5. [Action] (アクション) メニューで、[Delete schema] (スキーマの削除) をクリックします。
6. フィールドに「**Delete**」というテキストを入力して、削除を確定します。
7. [削除] を選択します。

指定したスキーマバージョンがレジストリから削除されます。

レジストリの削除

レジストリに含まれるスキーマの整理の必要性がなくなった場合は、そのレジストリを削除することができます。これらのスキーマは、別のレジストリに再割り当てする必要があります。

1 つ以上のレジストリを削除することは永続的なアクションであり、元に戻すことはできません。(1 つもしくは複数の) レジストリが不要であることを確認します。

デフォルトのレジストリは、AWS CLI を使用して削除できます。

AWS Glue API

レジストリ全体を、登録されたスキーマとそのすべてのバージョンとともに削除するには、[DeleteRegistry アクション \(Python: delete_registry\)](#) API を呼び出します。RegistryId 構造体を使用し、レジストリを特定します。

例:

```
aws glue delete-registry --registry-id RegistryArn="arn:aws:glue:us-east-2:901234567890:registry/registryName1"
```

```
aws glue delete-registry --registry-id RegistryName="TestRegistry-deletebyname"
```

削除オペレーションのステータスを取得するには、非同期呼び出し後に GetRegistry API を呼び出します。

AWS Glue コンソール

AWS Glue コンソールからレジストリの削除を行うには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインの [Data catalog] (データカタログ) で、[Schema registries] (スキーマレジストリ) をクリックします。
3. チェックボックスをオンにして、リストからレジストリを選択します。
4. [Action] (アクション) メニューで、[Delete registry] (レジストリの削除) を選択します。
5. フィールドに「**Delete**」というテキストを入力して、削除を確定します。
6. [削除] を選択します。

選択したレジストリが AWS Glue から削除されます。

シリアライザ用の IAM の例

Note

AWS 管理ポリシーは、一般的ユースケースに必要なアクセス許可を付与します。管理ポリシーを使用してスキーマのレジストリを管理する方法については、「[AWS の管理 \(定義済み\) ポリシー AWS Glue](#)」を参照してください。

シリアライザの場合、以下と同様の最小限のポリシーを作成して、特定のスキーマ定義のために `schemaVersionId` を検索できるようにします。レジストリ内のスキーマを読み取るには、そのレジストリに対する読み取り許可が必要であることに注意してください。読み取り可能なレジストリは、`Resource` 句を使用して制限できます。

コード例 13:

```
{
  "Sid" : "GetSchemaByDefinition",
  "Effect" : "Allow",
  "Action" :
  [
    "glue:GetSchemaByDefinition"
  ],
  "Resource" : ["arn:aws:glue:us-east-2:012345678:registry/registryname-1",
                "arn:aws:glue:us-east-2:012345678:schema/registryname-1/
schemaname-1",
                "arn:aws:glue:us-east-2:012345678:schema/registryname-1/
schemaname-2"
              ]
}
```

さらに、以下の新しいメソッドを追加して、プロデューサに対し新しいスキーマとバージョンの作成を許可することもできます。レジストリ内でスキーマを追加/削除/拡大させるためには、そのレジストリを調査できることが必要です。調査することが可能なレジストリは、`Resource` 句を使用して制限できます。

コード例 14:

```
{
  "Sid" : "RegisterSchemaWithMetadata",
```

```
"Effect" : "Allow",
"Action" :
[
    "glue:GetSchemaByDefinition",
    "glue:CreateSchema",
    "glue:RegisterSchemaVersion",
    "glue:PutSchemaVersionMetadata",
],
"Resource" : ["arn:aws:glue:aws-region:123456789012:registry/registryname-1",
              "arn:aws:glue:aws-region:123456789012:schema/registryname-1/
schemaname-1",
              "arn:aws:glue:aws-region:123456789012:schema/registryname-1/
schemaname-2"
              ]
}
```

デシリアライザー用の IAM の例

デシリアライザ (コンシューマ側) の場合、以下のようなポリシーを作成する必要があります。これにより、非シリアル化のために Schema Registry からスキーマを取得することを、デシリアライザに対し許可します。レジストリ内のスキーマを取得するためには、そのレジストリを調査することが許可されている必要があります。

コード例 15:

```
{
  "Sid" : "GetSchemaVersion",
  "Effect" : "Allow",
  "Action" :
  [
    "glue:GetSchemaVersion"
  ],
  "Resource" : ["*"]
}
```

AWS PrivateLink を使用したプライベート接続

AWS Glue でインターフェイス VPC エンドポイントを定義しながら AWS PrivateLinkを使用すると、データのプロデューサの VPC を AWS Glue に接続することができます。VPC インターフェイスエンドポイントにより、AWS ネットワーク内全体で VPC と AWS Glue 間の通信を処理します。詳細については、「[Using AWS Glue with VPC Endpoints](#)」を参照してください。

Amazon CloudWatch メトリクスへのアクセス

Amazon CloudWatch メトリクスは、CloudWatch の無料利用枠の一部として利用できます。これらのメトリクスには、CloudWatch コンソールからのアクセスが可能です。APIレベルのメトリクスとしては、CreateSchema (Success および Latency)、GetSchemaByDefinition (Success および Latency)、GetSchemaVersion (Success および Latency)、RegisterSchemaVersion (Success および Latency)、PutSchemaVersionMetadata (Success および Latency) があります。リソースレベルのメトリクスには、Registry.ThrottledByLimit、SchemaVersion.ThrottledByLimit、SchemaVersion.Size があります。

Schema Registry の AWS CloudFormation テンプレート例

以下に、AWS CloudFormation で Schema Registry リソースを作成するためのテンプレート例を示します。アカウントにこのスタックを作成するには、上記のテンプレートを SampleTemplate.yaml ファイルにコピーした上で、次のコマンドを実行します。

```
aws cloudformation create-stack --stack-name ABCSchemaRegistryStack --template-body
''cat SampleTemplate.yaml''
```

この例では、レジストリを作成するために AWS::Glue::Registry を、スキーマを作成するために AWS::Glue::Schema を、スキーマバージョンを作成するために AWS::Glue::SchemaVersion を使用し、AWS::Glue::SchemaVersionMetadata によりスキーマバージョンのメタデータを記述しています。

```
Description: "A sample CloudFormation template for creating Schema Registry resources."
Resources:
  ABCRegistry:
    Type: "AWS::Glue::Registry"
    Properties:
      Name: "ABCSchemaRegistry"
      Description: "ABC Corp. Schema Registry"
      Tags:
        - Key: "Project"
          Value: "Foo"
  ABCSchema:
    Type: "AWS::Glue::Schema"
    Properties:
      Registry:
        Arn: !Ref ABCRegistry
      Name: "TestSchema"
      Compatibility: "NONE"
```

```
DataFormat: "AVRO"
SchemaDefinition: >
  {"namespace":"foo.avro","type":"record","name":"user","fields":
[{"name":"name","type":"string"}, {"name":"favorite_number","type":"int"}]}
Tags:
  - Key: "Project"
    Value: "Foo"
SecondSchemaVersion:
  Type: "AWS::Glue::SchemaVersion"
Properties:
  Schema:
    SchemaArn: !Ref ABCSchema
    SchemaDefinition: >
      {"namespace":"foo.avro","type":"record","name":"user","fields":
[{"name":"status","type":"string", "default":"ON"}, {"name":"name","type":"string"},
{"name":"favorite_number","type":"int"}]}
FirstSchemaVersionMetadata:
  Type: "AWS::Glue::SchemaVersionMetadata"
  Properties:
    SchemaVersionId: !GetAtt ABCSchema.InitialSchemaVersionId
    Key: "Application"
    Value: "Kinesis"
SecondSchemaVersionMetadata:
  Type: "AWS::Glue::SchemaVersionMetadata"
  Properties:
    SchemaVersionId: !Ref SecondSchemaVersion
    Key: "Application"
    Value: "Kinesis"
```

AWS Glue Schema Registry との統合

以下のセクションで、AWS Glue Schema Registry との統合について説明します。このセクションの例では、AVRO データ形式のスキーマを使用します。JSON データ形式のスキーマなど、その他の例については、「[AWS Glue Schema Registry open source repository](#)」で、統合テストと README に関する情報をご覧ください。

トピック

- [ユースケース: Schema Registry を Amazon MSK または Apache Kafka に接続する](#)
- [ユースケース: Amazon Kinesis Data Streams と AWS Glue Schema Registry との統合](#)
- [Amazon Managed Service for Apache Flink のユースケース](#)

- [ユースケース: AWS Lambda との統合](#)
- [ユースケース: AWS Glue Data Catalog](#)
- [ユースケース: AWS Glue ストリーミング](#)
- [ユースケース: Apache Kafka ストリーム](#)
- [ユースケース: Apache Kafka Connect](#)

ユースケース: Schema Registry を Amazon MSK または Apache Kafka に接続する

Apache Kafka トピックにデータを書き込む場合には、以下の手順に従い作業を開始します。

1. Amazon Managed Streaming for Apache Kafka(Amazon MSK) または Apache Kafka のクラスターを作成し、少なくとも 1 つのトピックを含めます。Amazon MSK クラスターを作成する場合は、AWS Management Console を使用します。Amazon Managed Streaming for Apache Kafka デベロッパーガイドの「[Getting Started Using Amazon MSK](#)」にある手順に従います。
2. 上記の [SerDe ライブラリのインストール](#) ステップを実行します。
3. スキーマのレジストリ、スキーマ、またはスキーマバージョンを作成するには、このドキュメントにある [Schema Registry の使用開始](#) セクションの手順に従います。
4. Amazon MSK または Apache Kafka のトピックとの間で、レコードの書き込みや読み取りを行うために、Schema Registry を使用してプロデューサーとコンシューマーを起動します。プロデューサーとコンシューマーのコード例は、Serdeライブラリの [ReadMe ファイル](#) から入手できます。プロデューサーの Schema Registry ライブラリは、レコードを自動的にシリアル化し、スキーマバージョン ID でそのレコードを修飾します。
5. このレコードにスキーマが入力済みの場合、または自動登録が有効になっている場合には、スキーマが Schema Registry に登録されます。
6. AWS Glue Schema Registry ライブラリを使用して、Amazon MSK または Apache Kafka のトピックからスキーマの読み取りを行うコンシューマーは、自動的に Schema Registry からスキーマを検索します。

ユースケース: Amazon Kinesis Data Streams と AWS Glue Schema Registry との統合

この統合には、既存の Amazon Kinesis データストリームが必要です。詳細については、Amazon Kinesis Data Streams デベロッパーガイドの「[Getting Started with Amazon Kinesis Data Streams](#)」を参照してください。

Kinesis データストリームでは、データの操作用に以下の 2 つの方法があります。

- Java の Kinesis Producer Library (KPL) および Kinesis Client Library (KCL) ライブラリを使用します。多言語サポートは提供されていません。
- AWS SDK for Java に用意されている PutRecords、PutRecord、および GetRecords Kinesis Data Streams API を使用します。

現在、KPL/KCL ライブラリを使用中であれば、そのメソッドを引き続き使用することをお勧めします。ここでの例に示すように、Schema Registry が統合済みの、更新された KCL および KPL バージョンを使用できます。それ以外で、KDS API を直接使用している場合には、サンプルコードを通じて AWS Glue Schema Registry を利用します。

Schema Registry との統合は、KPL v0.14.2 以降と KCL v2.3 以降でのみ使用できます。JSON データ形式による Schema Registry との統合は、KPL v0.14.8 以降および KCL v2.3.6 以降で使用できます。

Kinesis SDK V2 を使用したデータの操作

このセクションでは、Kinesis SDK V2 による Kinesis の操作について説明します。

```
// Example JSON Record, you can construct a AVRO record also
private static final JsonDataWithSchema record =
    JsonDataWithSchema.builder(schemaString, payloadString);
private static final DataFormat dataFormat = DataFormat.JSON;

//Configurations for Schema Registry
GlueSchemaRegistryConfiguration gsrConfig = new GlueSchemaRegistryConfiguration("us-
east-1");

GlueSchemaRegistrySerializer glueSchemaRegistrySerializer =
    new GlueSchemaRegistrySerializerImpl(awsCredentialsProvider, gsrConfig);
GlueSchemaRegistryDataFormatSerializer dataFormatSerializer =
    new GlueSchemaRegistrySerializerFactory().getInstance(dataFormat, gsrConfig);

Schema gsrSchema =
    new Schema(dataFormatSerializer.getSchemaDefinition(record), dataFormat.name(),
        "MySchema");

byte[] serializedBytes = dataFormatSerializer.serialize(record);

byte[] gsrEncodedBytes = glueSchemaRegistrySerializer.encode(streamName, gsrSchema,
    serializedBytes);
```

```
PutRecordRequest putRecordRequest = PutRecordRequest.builder()
    .streamName(streamName)
    .partitionKey("partitionKey")
    .data(SdkBytes.fromByteArray(gsrEncodedBytes))
    .build();
shardId = kinesisClient.putRecord(putRecordRequest)
    .get()
    .shardId();

GlueSchemaRegistryDeserializer glueSchemaRegistryDeserializer = new
    GlueSchemaRegistryDeserializerImpl(awsCredentialsProvider, gsrConfig);

GlueSchemaRegistryDataFormatDeserializer gsrDataFormatDeserializer =
    glueSchemaRegistryDeserializerFactory.getInstance(dataFormat, gsrConfig);

GetShardIteratorRequest getShardIteratorRequest = GetShardIteratorRequest.builder()
    .streamName(streamName)
    .shardId(shardId)
    .shardIteratorType(ShardIteratorType.TRIM_HORIZON)
    .build();

String shardIterator = kinesisClient.getShardIterator(getShardIteratorRequest)
    .get()
    .shardIterator();

GetRecordsRequest getRecordRequest = GetRecordsRequest.builder()
    .shardIterator(shardIterator)
    .build();
GetRecordsResponse recordsResponse = kinesisClient.getRecords(getRecordRequest)
    .get();

List<Object> consumerRecords = new ArrayList<>();
List<Record> recordsFromKinesis = recordsResponse.records();

for (int i = 0; i < recordsFromKinesis.size(); i++) {
    byte[] consumedBytes = recordsFromKinesis.get(i)
        .data()
        .asByteArray();

    Schema gsrSchema = glueSchemaRegistryDeserializer.getSchema(consumedBytes);
    Object decodedRecord =
    gsrDataFormatDeserializer.deserialize(ByteBuffer.wrap(consumedBytes),

    gsrSchema.getSchemaDefinition());
```

```
consumerRecords.add(decodedRecord);  
}
```

KPL/KCL ライブラリを使用したデータの操作

このセクションでは、KPL/KCL ライブラリを使用する際の Kinesis Data Streams と Schema Registry の統合について説明します。KPL/KCL の使用方法については、Amazon Kinesis Data Streams デベロッパーガイドの「[Developing Producers Using the Amazon Kinesis Producer Library](#)」を参照してください。

KPL で Schema Registry を設定する

1. AWS Glue Schema Registry で作成したデータ、データ形式、スキーマ名のスキーマ定義を行います。
2. 必要に応じて、GlueSchemaRegistryConfiguration オブジェクトも構成します。
3. addUserRecord API にスキーマオブジェクトを渡します。

```
private static final String SCHEMA_DEFINITION = "{\"namespace\": \"example.avro\",\\n\"  
+ \" \"type\": \"record\",\\n\"  
+ \" \"name\": \"User\",\\n\"  
+ \" \"fields\": [\\n\"  
+ \" {\"name\": \"name\", \"type\": \"string\"},\\n\"  
+ \" {\"name\": \"favorite_number\", \"type\": [\"int\", \"null\"]},\\n\"  
+ \" {\"name\": \"favorite_color\", \"type\": [\"string\", \"null\"]}\\n\"  
+ \" ]\\n\"  
+ \"}\";
```

```
KinesisProducerConfiguration config = new KinesisProducerConfiguration();  
config.setRegion("us-west-1")
```

```
//[Optional] configuration for Schema Registry.
```

```
GlueSchemaRegistryConfiguration schemaRegistryConfig =  
new GlueSchemaRegistryConfiguration("us-west-1");
```

```
schemaRegistryConfig.setCompression(true);
```

```
config.setGlueSchemaRegistryConfiguration(schemaRegistryConfig);
```

```
///Optional configuration ends.
```

```
final KinesisProducer producer =
```

```
new KinesisProducer(config);

final ByteBuffer data = getDataToSend();

com.amazonaws.services.schemaregistry.common.Schema gsrSchema =
new Schema(SCHEMA_DEFINITION, DataFormat.AVRO.toString(), "demoSchema");

ListenableFuture<UserRecordResult> f = producer.addUserRecord(
config.getStreamName(), TIMESTAMP, Utils.randomExplicitHashKey(), data, gsrSchema);

private static ByteBuffer getDataToSend() {
    org.apache.avro.Schema avroSchema =
        new org.apache.avro.Schema.Parser().parse(SCHEMA_DEFINITION);

    GenericRecord user = new GenericData.Record(avroSchema);
    user.put("name", "Emily");
    user.put("favorite_number", 32);
    user.put("favorite_color", "green");

    ByteArrayOutputStream outBytes = new ByteArrayOutputStream();
    Encoder encoder = EncoderFactory.get().directBinaryEncoder(outBytes, null);
    new GenericDatumWriter<>(avroSchema).write(user, encoder);
    encoder.flush();
    return ByteBuffer.wrap(outBytes.toByteArray());
}
```

Kinesis Client Library のセットアップ

Kinesis Client Library コンシューマーを、Java により構築します。詳細については、Amazon Kinesis Data Streams デベロッパーガイドの「[Developing a Kinesis Client Library Consumer in Java](#)」を参照してください。

1. `GlueSchemaRegistryConfiguration` オブジェクトを渡すことで `GlueSchemaRegistryDeserializer` インスタンスを作成します。
2. `GlueSchemaRegistryDeserializer` を `retrievalConfig.glueSchemaRegistryDeserializer` に渡します。
3. `kinesisClientRecord.getSchema()` を呼び出して、受信メッセージのスキーマにアクセスします。

```
GlueSchemaRegistryConfiguration schemaRegistryConfig =
new GlueSchemaRegistryConfiguration(this.region.toString());
```

```

GlueSchemaRegistryDeserializer glueSchemaRegistryDeserializer =
    new
GlueSchemaRegistryDeserializerImpl(DefaultCredentialsProvider.builder().build(),
schemaRegistryConfig);

RetrievalConfig retrievalConfig =
configsBuilder.retrievalConfig().retrievalSpecificConfig(new
PollingConfig(streamName, kinesisClient));
retrievalConfig.glueSchemaRegistryDeserializer(glueSchemaRegistryDeserializer);

Scheduler scheduler = new Scheduler(
    configsBuilder.checkpointConfig(),
    configsBuilder.coordinatorConfig(),
    configsBuilder.leaseManagementConfig(),
    configsBuilder.lifecycleConfig(),
    configsBuilder.metricsConfig(),
    configsBuilder.processorConfig(),
    retrievalConfig
);

public void processRecords(ProcessRecordsInput processRecordsInput) {
    MDC.put(SHARD_ID_MDC_KEY, shardId);
    try {
        log.info("Processing {} record(s)",
            processRecordsInput.records().size());
        processRecordsInput.records()
            .forEach(
                r ->
                    log.info("Processed record pk: {} -- Seq: {} : data {} with
schema: {}",
                        r.partitionKey(),
r.sequenceNumber(), recordToAvroObj(r).toString(), r.getSchema());
            } catch (Throwable t) {
                log.error("Caught throwable while processing records. Aborting.");
                Runtime.getRuntime().halt(1);
            } finally {
                MDC.remove(SHARD_ID_MDC_KEY);
            }
    }

private GenericRecord recordToAvroObj(KinesisClientRecord r) {
    byte[] data = new byte[r.data().remaining()];
    r.data().get(data, 0, data.length);
}

```

```
org.apache.avro.Schema schema = new
org.apache.avro.Schema.Parser().parse(r.schema().getSchemaDefinition());
DatumReader datumReader = new GenericDatumReader<>(schema);

BinaryDecoder binaryDecoder = DecoderFactory.get().binaryDecoder(data, 0,
data.length, null);
return (GenericRecord) datumReader.read(null, binaryDecoder);
}
```

Kinesis Data Streams API を使用したデータの操作

このセクションでは、Kinesis Data Streams API を使用しての、Kinesis Data Streams と Schema Registry の統合について説明します。

1. Maven の以下の依存関係を更新します。

```
<dependencyManagement>
  <dependencies>
    <dependency>
      <groupId>com.amazonaws</groupId>
      <artifactId>aws-java-sdk-bom</artifactId>
      <version>1.11.884</version>
      <type>pom</type>
      <scope>import</scope>
    </dependency>
  </dependencies>
</dependencyManagement>

<dependencies>
  <dependency>
    <groupId>com.amazonaws</groupId>
    <artifactId>aws-java-sdk-kinesis</artifactId>
  </dependency>

  <dependency>
    <groupId>software.amazon.glue</groupId>
    <artifactId>schema-registry-serde</artifactId>
    <version>1.1.5</version>
  </dependency>

  <dependency>
    <groupId>com.fasterxml.jackson.dataformat</groupId>
```

```

        <artifactId>jackson-dataformat-cbor</artifactId>
        <version>2.11.3</version>
    </dependency>
</dependencies>

```

- PutRecords または Kinesis Data Streams の PutRecord API を使用しながら、プロデューサー内にスキーマヘッダー情報を追加します。

```

//The following lines add a Schema Header to the record
    com.amazonaws.services.schemaregistry.common.Schema awsSchema =
        new com.amazonaws.services.schemaregistry.common.Schema(schemaDefinition,
DataFormat.AVRO.name(),
            schemaName);
    GlueSchemaRegistrySerializerImpl glueSchemaRegistrySerializer =
        new
GlueSchemaRegistrySerializerImpl(DefaultCredentialsProvider.builder().build(), new
GlueSchemaRegistryConfiguration(getConfigs()));
    byte[] recordWithSchemaHeader =
        glueSchemaRegistrySerializer.encode(streamName, awsSchema,
recordAsBytes);

```

- プロデューサー内で PutRecords または PutRecord API を使用して、レコードをデータストリームに配置します。
- コンシューマ内で、ヘッダーからスキーマレコードを削除し、Avro スキーマレコードをシリアル化します。

```

//The following lines remove Schema Header from record
    GlueSchemaRegistryDeserializerImpl glueSchemaRegistryDeserializer =
        new
GlueSchemaRegistryDeserializerImpl(DefaultCredentialsProvider.builder().build(),
getConfigs());
    byte[] recordWithSchemaHeaderBytes = new
byte[recordWithSchemaHeader.remaining()];
    recordWithSchemaHeader.get(recordWithSchemaHeaderBytes, 0,
recordWithSchemaHeaderBytes.length);
    com.amazonaws.services.schemaregistry.common.Schema awsSchema =
        glueSchemaRegistryDeserializer.getSchema(recordWithSchemaHeaderBytes);
    byte[] record =
glueSchemaRegistryDeserializer.getData(recordWithSchemaHeaderBytes);

//The following lines serialize an AVRO schema record
if (DataFormat.AVRO.name().equals(awsSchema.getDataFormat())) {

```

```
        Schema avroSchema = new
org.apache.avro.Schema.Parser().parse(awsSchema.getSchemaDefinition());
        Object genericRecord = convertBytesToRecord(avroSchema, record);
        System.out.println(genericRecord);
    }
```

Kinesis Data Streams API を使用したデータの操作

以下に、PutRecords および GetRecords API を使用するコード例を示します。

```
//Full sample code
import
    com.amazonaws.services.schemaregistry.deserializers.GlueSchemaRegistryDeserializerImpl;
import
    com.amazonaws.services.schemaregistry.serializers.GlueSchemaRegistrySerializerImpl;
import com.amazonaws.services.schemaregistry.utils.AVROUtils;
import com.amazonaws.services.schemaregistry.utils.AWSSchemaRegistryConstants;
import org.apache.avro.Schema;
import org.apache.avro.generic.GenericData;
import org.apache.avro.generic.GenericDatumReader;
import org.apache.avro.generic.GenericDatumWriter;
import org.apache.avro.generic.GenericRecord;
import org.apache.avro.io.Decoder;
import org.apache.avro.io.DecoderFactory;
import org.apache.avro.io.Encoder;
import org.apache.avro.io.EncoderFactory;
import software.amazon.awssdk.auth.credentials.DefaultCredentialsProvider;
import software.amazon.awssdk.services.glue.model.DataFormat;

import java.io.ByteArrayOutputStream;
import java.io.File;
import java.io.IOException;
import java.nio.ByteBuffer;
import java.util.Collections;
import java.util.HashMap;
import java.util.Map;

public class PutAndGetExampleWithEncodedData {
    static final String regionName = "us-east-2";
    static final String streamName = "testStream1";
    static final String schemaName = "User-Topic";
    static final String AVRO_USER_SCHEMA_FILE = "src/main/resources/user.avsc";
```

```
KinesisApi kinesisApi = new KinesisApi();

void runSampleForPutRecord() throws IOException {
    Object testRecord = getTestRecord();
    byte[] recordAsBytes = convertRecordToBytes(testRecord);
    String schemaDefinition =
AVROUtils.getInstance().getSchemaDefinition(testRecord);

    //The following lines add a Schema Header to a record
    com.amazonaws.services.schemaregistry.common.Schema awsSchema =
        new com.amazonaws.services.schemaregistry.common.Schema(schemaDefinition,
DataFormat.AVRO.name(),
            schemaName);
    GlueSchemaRegistrySerializerImpl glueSchemaRegistrySerializer =
        new
GlueSchemaRegistrySerializerImpl(DefaultCredentialsProvider.builder().build(), new
GlueSchemaRegistryConfiguration(regionName));
    byte[] recordWithSchemaHeader =
        glueSchemaRegistrySerializer.encode(streamName, awsSchema, recordAsBytes);

    //Use PutRecords api to pass a list of records
    kinesisApi.putRecords(Collections.singletonList(recordWithSchemaHeader),
streamName, regionName);

    //OR
    //Use PutRecord api to pass single record
    //kinesisApi.putRecord(recordWithSchemaHeader, streamName, regionName);
}

byte[] runSampleForGetRecord() throws IOException {
    ByteBuffer recordWithSchemaHeader = kinesisApi.getRecords(streamName,
regionName);

    //The following lines remove the schema registry header
    GlueSchemaRegistryDeserializerImpl glueSchemaRegistryDeserializer =
        new
GlueSchemaRegistryDeserializerImpl(DefaultCredentialsProvider.builder().build(), new
GlueSchemaRegistryConfiguration(regionName));
    byte[] recordWithSchemaHeaderBytes = new
byte[recordWithSchemaHeader.remaining()];
    recordWithSchemaHeader.get(recordWithSchemaHeaderBytes, 0,
recordWithSchemaHeaderBytes.length);

    com.amazonaws.services.schemaregistry.common.Schema awsSchema =
```

```
    glueSchemaRegistryDeserializer.getSchema(recordWithSchemaHeaderBytes);

    byte[] record =
glueSchemaRegistryDeserializer.getData(recordWithSchemaHeaderBytes);

    //The following lines serialize an AVRO schema record
    if (DataFormat.AVRO.name().equals(awsSchema.getDataFormat())) {
        Schema avroSchema = new
org.apache.avro.Schema.Parser().parse(awsSchema.getSchemaDefinition());
        Object genericRecord = convertBytesToRecord(avroSchema, record);
        System.out.println(genericRecord);
    }

    return record;
}

private byte[] convertRecordToBytes(final Object record) throws IOException {
    ByteArrayOutputStream recordAsBytes = new ByteArrayOutputStream();
    Encoder encoder = EncoderFactory.get().directBinaryEncoder(recordAsBytes,
null);
    GenericDatumWriter datumWriter = new
GenericDatumWriter<>(AVROUtils.getInstance().getSchema(record));
    datumWriter.write(record, encoder);
    encoder.flush();
    return recordAsBytes.toByteArray();
}

private GenericRecord convertBytesToRecord(Schema avroSchema, byte[] record) throws
IOException {
    final GenericDatumReader<GenericRecord> datumReader = new
GenericDatumReader<>(avroSchema);
    Decoder decoder = DecoderFactory.get().binaryDecoder(record, null);
    GenericRecord genericRecord = datumReader.read(null, decoder);
    return genericRecord;
}

private Map<String, String> getMetadata() {
    Map<String, String> metadata = new HashMap<>();
    metadata.put("event-source-1", "topic1");
    metadata.put("event-source-2", "topic2");
    metadata.put("event-source-3", "topic3");
    metadata.put("event-source-4", "topic4");
    metadata.put("event-source-5", "topic5");
    return metadata;
}
```

```
    }

    private GlueSchemaRegistryConfiguration getConfigs() {
        GlueSchemaRegistryConfiguration configs = new
GlueSchemaRegistryConfiguration(regionName);
        configs.setSchemaName(schemaName);
        configs.setAutoRegistration(true);
        configs.setMetadata(getMetadata());
        return configs;
    }

    private Object getTestRecord() throws IOException {
        GenericRecord genericRecord;
        Schema.Parser parser = new Schema.Parser();
        Schema avroSchema = parser.parse(new File(AVRO_USER_SCHEMA_FILE));

        genericRecord = new GenericData.Record(avroSchema);
        genericRecord.put("name", "testName");
        genericRecord.put("favorite_number", 99);
        genericRecord.put("favorite_color", "red");

        return genericRecord;
    }
}
```

Amazon Managed Service for Apache Flink のユースケース

Apache Flinkは、無制限および制限付きのデータストリームに対するステートフルな計算に広く使用されている、オープンソースフレームワークの分散処理エンジンです。Amazon Managed Service for Apache Flink は、ストリーミングデータを処理するため、Apache Flink アプリケーションを構築して管理できるようにする完全マネージド型の AWS サービスです。

オープンソースの Apache Flink では、多数のソースとシンクを利用できます。例えば、事前定義済みのデータソースには、ファイル、ディレクトリ、およびソケットからの読み込みや、コレクションとイテレータからのデータの取り込みなどが含まれています。Apache Flink DataStream Connector は、Apache Flinkが、Apache Kafka や Kinesis などの各種サードパーティー製システムと、ソースおよび/またはシンクとしてインターフェースするためのコードを提供します。

詳細については、[Amazon Kinesis Data Analytics デベロッパーガイド](#)を参照してください。

Apache Flink Kafka Connector

Apache Flinkは、Kafka のトピックに対するデータの読み取りおよび書き込みを、正確に一度で行えるようにするための、Apache Kafka データストリームのコネクタを提供します。Flink の Kafka コンシューマ `FlinkKafkaConsumer` では、1 つ以上の Kafka トピックから読み取りを行うアクセスが提供されます。Apache Flink の Kafka プロデューサ `FlinkKafkaProducer` では、1 つ以上の Kafka トピックに対しレコードのストリームを書き込むことができます。詳細については、「[Apache Kafka Connector](#)」を参照してください。

Apache Flink Kinesis Streams Connector

Kinesis データストリームのコネクタは、Amazon Kinesis Data Streams へのアクセスを提供します。並列ストリーミングデータソース `FlinkKinesisConsumer` は、同じAWS のサービスリージョン内で複数の Kinesis ストリームにサブスクライブされ、ジョブの実行中にストリームの再シャーディングを透過的に (確実に 1 回で) 処理できます。コンシューマーの各サブタスクが、複数の Kinesis シャードからのデータレコードの取得を受け持ちます。各サブタスクによって取得されるシャードの数は、Kinesis によってシャードが閉じられ、また作成されるたびに変化します。`FlinkKinesisProducer` は Kinesis Producer Library (KPL) を使用して、Kinesis ストリーム内に Apache Flink ストリームからのデータを配置します。詳細については、「[Amazon Kinesis Streams Connector](#)」を参照してください。

詳細については、[AWS Glue のGitHub リポジトリ](#)を参照してください。

Apache Flink との統合

Schema Registry で提供されている SerDes ライブラリは、Apache Flink と統合されています。Apache Flink を使用するには、[SerializationSchema](#) および [DeserializationSchema](#) のインターフェイス (`GlueSchemaRegistryAvroSerializationSchema` および `GlueSchemaRegistryAvroDeserializationSchema`) を実装する必要があります。これらは、Apache Flink コネクタにプラグインして使用します。

Apache Flink アプリケーションへの AWS Glue Schema Registry の依存関係の追加

Apache Flink アプリケーション内で AWS Glue Schema Registry との統合の依存関係をセットアップするには

1. `pom.xml` ファイルに依存関係を追加します。

```
<dependency>
```

```

<groupId>software.amazon.glue</groupId>
<artifactId>schema-registry-flink-serde</artifactId>
<version>1.0.0</version>
</dependency>

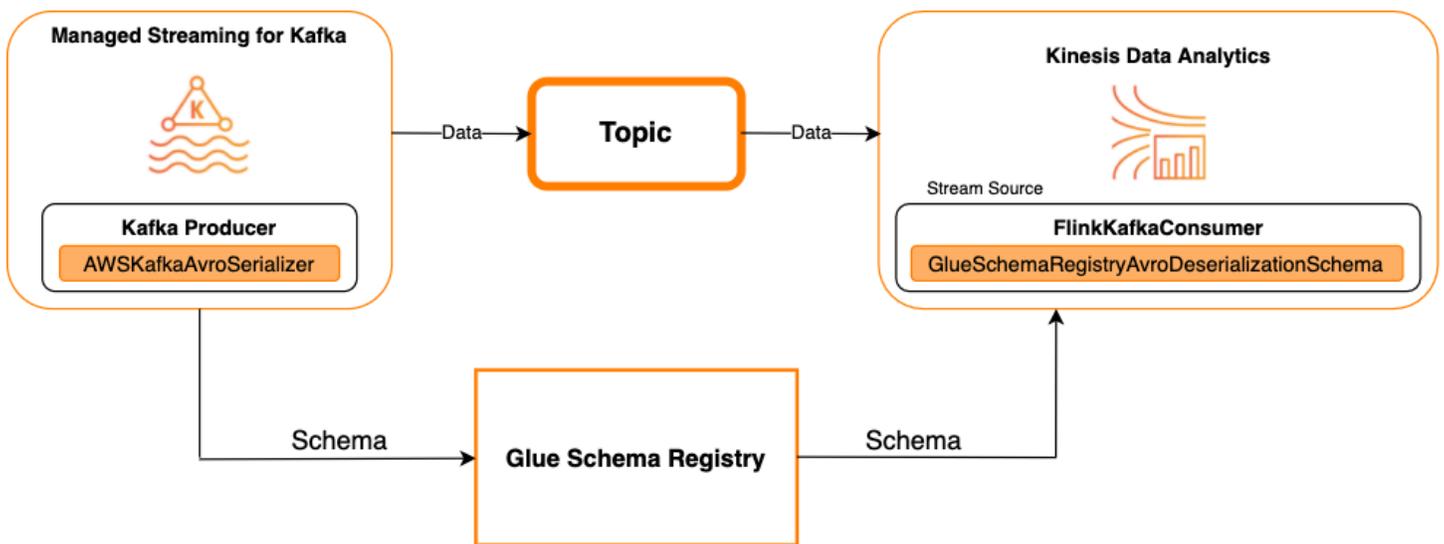
```

Kafka または Amazon MSK を Apache Flink と統合する

Kafka をソースまたはシンクとしながら、Apache Flink 対応の Managed Service for Apache Flink を使用できます。

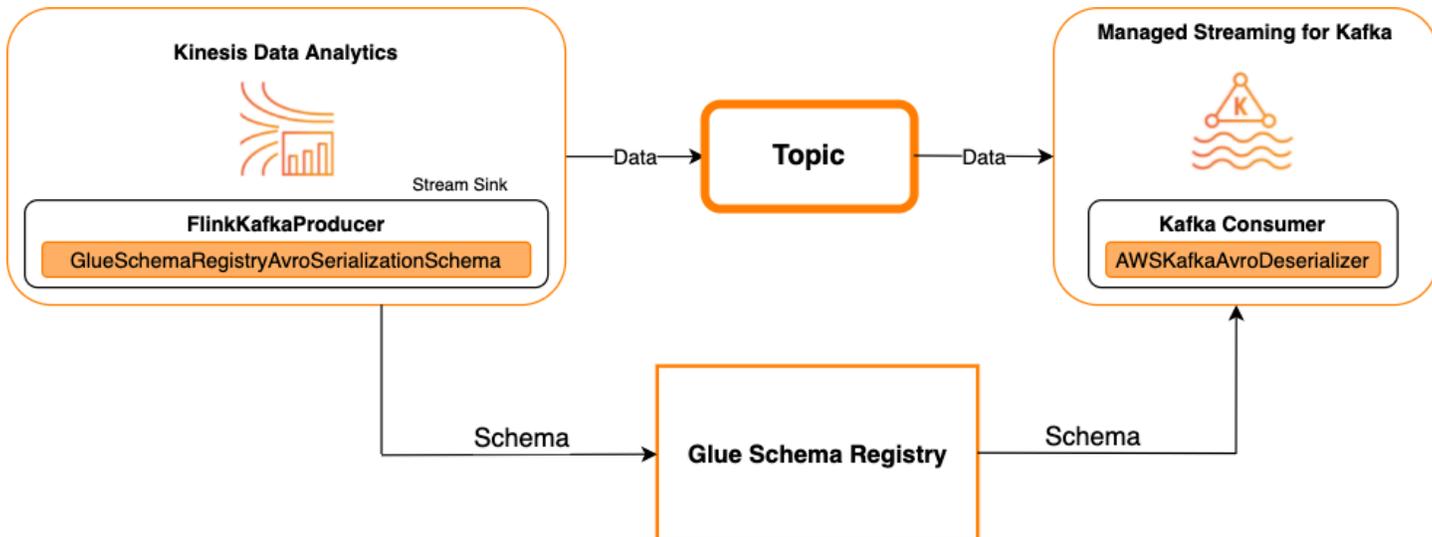
Kafka をソースとする場合

次の図は、Kafka をソースとしながら、Kinesis Data Streams と Apache Flink 対応の Managed Service for Apache Flink を統合した様子です。



Kafka をシンクとする場合

次の図は、Kafka をシンクとしながら、Kinesis Data Streams と Apache Flink 対応の Managed Service for Apache Flink を統合した様子です。



Kafka をソースまたはシンクとしながら、Kafka (または Amazon MSK) を Apache Flink 対応の Managed Service for Apache Flink と統合するには、以下のコード変更を行います。太字で示されたコードブロックを、類似するセクション内の対応するコードにそれぞれ追加します。

Kafka をソースとする場合は、デシリアライザ用コード (ブロック 2) を使用します。Kafka をシンクとする場合は、シリアライザコード (ブロック 3) を使用します。

```

StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment();

String topic = "topic";
Properties properties = new Properties();
properties.setProperty("bootstrap.servers", "localhost:9092");
properties.setProperty("group.id", "test");

// block 1
Map<String, Object> configs = new HashMap<>();
configs.put(AWSSchemaRegistryConstants.AWS_REGION, "aws-region");
configs.put(AWSSchemaRegistryConstants.SCHEMA_AUTO_REGISTRATION_SETTING, true);
configs.put(AWSSchemaRegistryConstants.AVRO_RECORD_TYPE,
    AvroRecordType.GENERIC_RECORD.getName());

FlinkKafkaConsumer<GenericRecord> consumer = new FlinkKafkaConsumer<>(
    topic,
    // block 2
    GlueSchemaRegistryAvroDeserializationSchema.forGeneric(schema, configs),
    properties);

FlinkKafkaProducer<GenericRecord> producer = new FlinkKafkaProducer<>(

```

```
topic,
// block 3
GlueSchemaRegistryAvroSerializationSchema.forGeneric(schema, topic, configs),
properties);
```

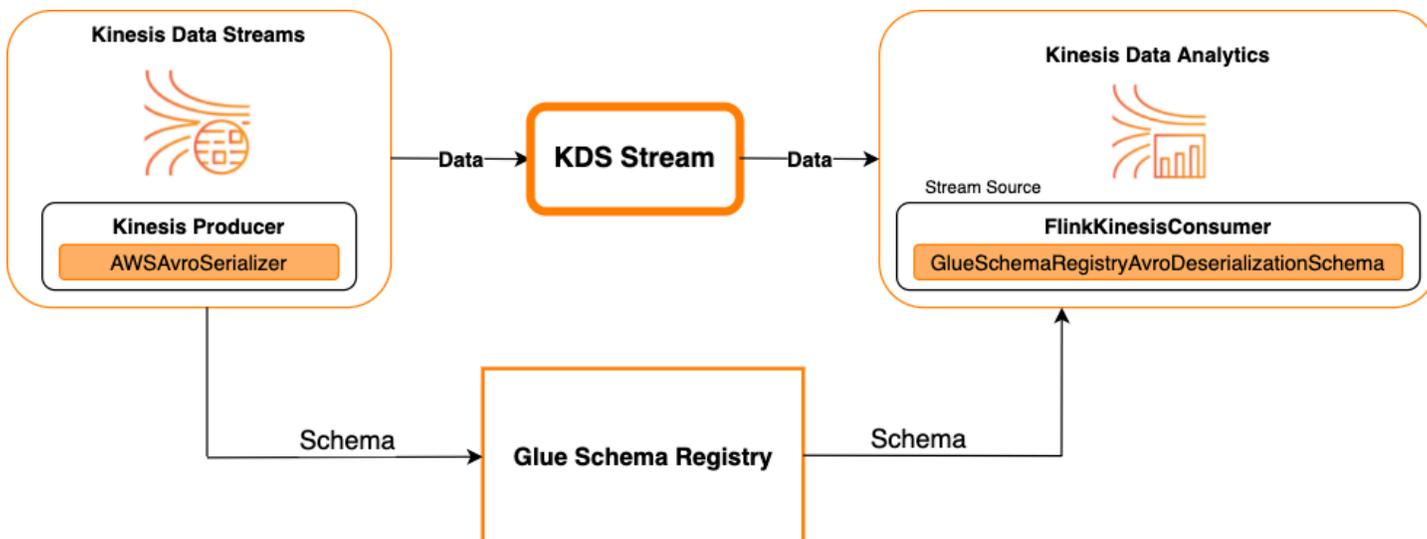
```
DataStream<GenericRecord> stream = env.addSource(consumer);
stream.addSink(producer);
env.execute();
```

Kinesis Data Streams と Apache Flink との統合

Kinesis Data Streams をソースまたはシンクとしながら、Apache Flink 対応の Managed Service for Apache Flink を使用できます。

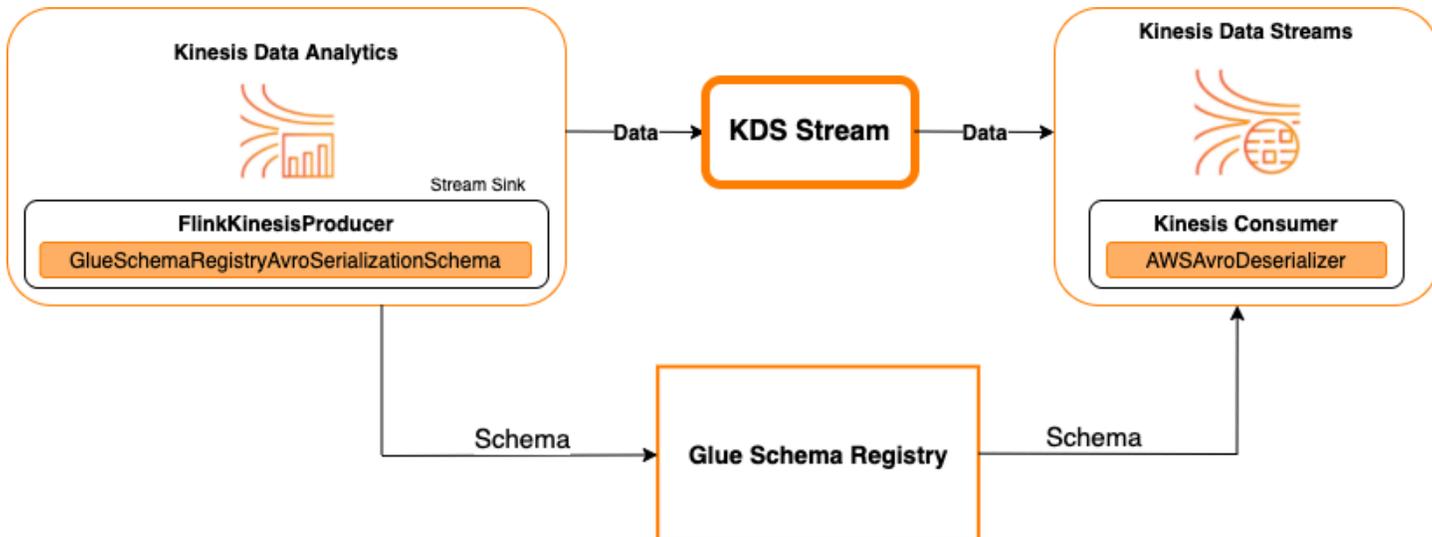
Kinesis Data Streams をソースとする場合

次の図は、Kinesis Data Streams をソースとしながら、Kinesis Data Streams と Apache Flink 対応の Managed Service for Apache Flink を統合した様子です。



Kinesis Data Streams をシンクとする場合

次の図は、Kinesis Data Streams をシンクとしながら、Kinesis Data Streams と Apache Flink 対応の Managed Service for Apache Flink を統合した様子です。



Kinesis Data Streams をソースまたはシンクとしながら、Kinesis Data Streams と Apache Flink 対応の Managed Service for Apache Flink を統合するには、以下のコード変更を行います。太字で示されたコードブロックを、類似するセクション内の対応するコードにそれぞれ追加します。

Kinesis Data Streams をソースとする場合は、デシリアライザコード (ブロック 2) を使用します。Kinesis Data Streams をシンクとする場合は、シリアライザコード (ブロック 3) を使用します。

```

StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment();

String streamName = "stream";
Properties consumerConfig = new Properties();
consumerConfig.put(AWSConfigConstants.AWS_REGION, "aws-region");
consumerConfig.put(AWSConfigConstants.AWS_ACCESS_KEY_ID, "aws_access_key_id");
consumerConfig.put(AWSConfigConstants.AWS_SECRET_ACCESS_KEY, "aws_secret_access_key");
consumerConfig.put(ConsumerConfigConstants.STREAM_INITIAL_POSITION, "LATEST");

// block 1
Map<String, Object> configs = new HashMap<>();
configs.put(AWSSchemaRegistryConstants.AWS_REGION, "aws-region");
configs.put(AWSSchemaRegistryConstants.SCHEMA_AUTO_REGISTRATION_SETTING, true);
configs.put(AWSSchemaRegistryConstants.AVRO_RECORD_TYPE,
AvroRecordType.GENERIC_RECORD.getName());

FlinkKinesisConsumer<GenericRecord> consumer = new FlinkKinesisConsumer<>(
    streamName,
    // block 2
GlueSchemaRegistryAvroDeserializationSchema.forGeneric(schema, configs),
    properties);
  
```

```
FlinkKinesisProducer<GenericRecord> producer = new FlinkKinesisProducer<>(  
    // block 3  
    GlueSchemaRegistryAvroSerializationSchema.forGeneric(schema, topic, configs),  
    properties);  
producer.setDefaultStream(streamName);  
producer.setDefaultPartition("0");  
  
DataStream<GenericRecord> stream = env.addSource(consumer);  
stream.addSink(producer);  
env.execute();
```

ユースケース: AWS Lambda との統合

AWS Lambda 関数を Apache Kafka/Amazon MSK のコンシューマーとして使用し、Avro でエンコードされたメッセージを AWS Glue Schema Registry により非シリアル化するには、[MSK ラボのページ](#)をご覧ください。

ユースケース: AWS Glue Data Catalog

AWS Glue テーブルは、手動による指定、または AWS Glue Schema Registry への参照によって指定できる、スキーマをサポートしています。AWS Glue テーブルまたは Data Catalog のパーティションを作成または更新する際に、オプションで Schema Registry に格納されているスキーマを使用できるように、Schema Registry には Data Catalog が統合されています。Schema Registry のスキーマ定義を特定するには、少なくとも、対象となるスキーマの ARN を知る必要があります。スキーマ定義を含むスキーマのスキーマバージョンは、UUID またはバージョン番号により参照が可能です。スキーマバージョンの中でも「最新」バージョンについては、バージョン番号または UUID を把握しなくても常に参照することができます。

CreateTable または UpdateTable オペレーションの呼び出し時は、Schema Registry 内の既存のスキーマに対する TableInput を指定するために SchemaReference 構造体 (StorageDescriptor を含む) を渡します。同様に、GetTable または GetPartition API を呼び出す場合は、そのレスポンスにスキーマと SchemaReference が含まれます。スキーマ参照を使用してテーブルまたはパーティションが作成されると、Data Catalog はこのスキーマ参照のスキーマ取得を試みます。Schema Registry 内にスキーマが見つからない場合は、GetTable レスポンスで空のスキーマを返します。それ以外では、このレスポンスにスキーマとスキーマ参照の両方が出力されます。

また、AWS Glue コンソールからアクションを実行することも可能です。

これらのオペレーションを実行し、スキーマ情報を作成、更新、表示するには、呼び出しユーザーに、GetSchemaVersion API へのアクセス権限を付与する、IAM ロールを付与する必要があります。

テーブルの追加またはテーブルのスキーマの更新

既存のスキーマから新しいテーブルを追加すると、そのテーブルは特定のスキーマバージョンにバインドされます。新しいスキーマバージョンの登録が完了すると、このテーブル定義が、AWS Glue コンソールの [View table] (テーブルの表示) ページ、もしくは [UpdateTable アクション \(Python: update_table\)](#) API を使用して更新できるようになります。

既存のスキーマからのテーブルの追加

AWS Glue コンソールまたは CreateTable API を使用して、レジストリ内のスキーマバージョンから AWS Glue を作成できます。

AWS Glue API

CreateTable API を呼び出す際に、(StorageDescriptor に SchemaReference が指定されている) TableInput を、スキーマレジストリの既存のスキーマに追加します。

AWS Glue コンソール

AWS Glue コンソールを使用してテーブルを作成するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインの [Data catalog] (データカタログ) で、[Tables] (テーブル) をクリックします。
3. [Add Tables] (テーブルの追加) メニューで、[Add table from existing schema] (既存のスキーマからテーブルを追加する) をクリックします。
4. テーブルのプロパティとデータストアを、AWS Glue デベロッパーガイド に沿って設定します。
5. [Choose a Glue schema] (Glue スキーマの選択) ページで、スキーマが置かれている [Registry] (レジストリ) を選択します。
6. [Schema name] (スキーマ名) をクリックし、適用するスキーマの [Version] (バージョン) を選択します。
7. スキーマのプレビューを確認し、[Next] (次へ) をクリックします。

8. テーブルを確認し、作成します。

作成したテーブルに適用されたスキーマとバージョンは、テーブルの一覧内で [Glue schema] (Glue スキーマ) 列に表示されます。テーブルを表示すると、さらに詳細を確認できます。

テーブルのスキーマの更新

新しいスキーマバージョンが使用可能になったら、テーブルのスキーマを [UpdateTable アクション \(Python: update_table\)](#) API または AWS Glue コンソールにより更新することができます

Important

手動で指定された AWS Glue スキーマを含む既存のテーブル用にスキーマを更新する場合、Schema Registry で参照される新しいスキーマは互換性を持たない可能性があります。この場合、ジョブが失敗することがあります。

AWS Glue API

UpdateTable API を呼び出す際に、(StorageDescriptor に SchemaReference が指定されている) TableInput を、スキーマレジストリの既存のスキーマに追加します。

AWS Glue コンソール

AWS Glue コンソールからテーブルのスキーマを更新するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインの [Data catalog] (データカタログ) で、[Tables] (テーブル) をクリックします。
3. テーブルの一覧でテーブルを表示します。
4. 新しいバージョンの情報が表示されたボックスで、[Update schema] (スキーマの更新) をクリックします。
5. 現在のスキーマと更新後のスキーマの違いを確認します。
6. さらに詳細を表示するには、[Show all schema differences] (スキーマの違いをすべて表示) をクリックします。
7. [Save table] (テーブルを保存) をクリックし、新しいバージョンを受け入れます。

ユースケース: AWS Glue ストリーミング

AWS Glue ストリーミングは、ストリーミングソースからのデータを消費し、出力シンクに書き込む前に ETL オペレーションを実行します。入力ストリーミングソースは、データテーブルを使用して指定するか、ソース構成を指定して直接指定することができます。

AWS Glue ストリーミングは、AWS Glue スキーマレジストリに存在するスキーマで作成されたストリーミングソースのデータカタログテーブルをサポートします。AWS Glue スキーマレジストリにスキーマを作成し、そのスキーマを使用してストリーミングソースで AWS Glue テーブルを作成できます。この AWS Glue テーブルは、AWS Glue ストリーミングジョブへの入力として使用し、入力ストリーミングのデータを逆シリアル化することができます。

注意すべき点は、AWS Glue スキーマレジスト内のスキーマが変化した場合、AWS Glue ストリーミングジョブを再度開始して、スキーマの変更を反映させる必要があることです。

ユースケース: Apache Kafka ストリーム

Apache Kafka Streams APIは、Apache Kafka に格納されているデータを処理・分析するためのクライアントライブラリです。このセクションでは、Apache Kafka Streams と AWS Glue Schema Registry の統合について説明します。これにより、データストリーミングアプリケーションのスキーマを、管理および適用できるようになります。Apache Kafka Streams の詳細については、「[Apache Kafka Streams](#)」を参照してください。

SerDes ライブラリとの統合

GlueSchemaRegistryKafkaStreamsSerde クラスにより、Streams のアプリケーションを設定できます。

Kafka Streams アプリケーションのコード例

Apache Kafka Streams アプリケーション内で AWS Glue Schema Registry を使用するには

1. Kafka Streams アプリケーションを設定します。

```
final Properties props = new Properties();
    props.put(StreamsConfig.APPLICATION_ID_CONFIG, "avro-streams");
    props.put(StreamsConfig.BOOTSTRAP_SERVERS_CONFIG, "localhost:9092");
    props.put(StreamsConfig.CACHE_MAX_BYTES_BUFFERING_CONFIG, 0);
    props.put(StreamsConfig.DEFAULT_KEY_SERDE_CLASS_CONFIG,
Serdes.String().getClass().getName());
    props.put(StreamsConfig.DEFAULT_VALUE_SERDE_CLASS_CONFIG,
AWSKafkaAvroSerDe.class.getName());
```

```
props.put(ConsumerConfig.AUTO_OFFSET_RESET_CONFIG, "earliest");

props.put(AWSSchemaRegistryConstants.AWS_REGION, "aws-region");
props.put(AWSSchemaRegistryConstants.SCHEMA_AUTO_REGISTRATION_SETTING, true);
props.put(AWSSchemaRegistryConstants.AVRO_RECORD_TYPE,
AvroRecordType.GENERIC_RECORD.getName());
props.put(AWSSchemaRegistryConstants.DATA_FORMAT, DataFormat.AVRO.name());
```

2. トピック avro-input からストリームを作成します。

```
StreamsBuilder builder = new StreamsBuilder();
final KStream<String, GenericRecord> source = builder.stream("avro-input");
```

3. データレコードを処理します (favorite_color の値がピンクであるか、値が 15 となっているレコードの除外など)。

```
final KStream<String, GenericRecord> result = source
    .filter((key, value) -
> !"pink".equals(String.valueOf(value.get("favorite_color"))));
    .filter((key, value) -> !"15.0".equals(String.valueOf(value.get("amount"))));
```

4. トピック avro-output に結果を書き込みます。

```
result.to("avro-output");
```

5. Apache Kafka Streams アプリケーションを起動します。

```
KafkaStreams streams = new KafkaStreams(builder.build(), props);
streams.start();
```

実装結果

以下の結果は、ステップ 3 において (favorite_color が「pink」であるか値が「15.0」であるために) 除外されたレコードに関するフィルタリング処理を示しています。

フィルタリング前のレコード:

```
{"name": "Sansa", "favorite_number": 99, "favorite_color": "white"}
{"name": "Harry", "favorite_number": 10, "favorite_color": "black"}
{"name": "Hermione", "favorite_number": 1, "favorite_color": "red"}
{"name": "Ron", "favorite_number": 0, "favorite_color": "pink"}
{"name": "Jay", "favorite_number": 0, "favorite_color": "pink"}

{"id": "commute_1", "amount": 3.5}
{"id": "grocery_1", "amount": 25.5}
{"id": "entertainment_1", "amount": 19.2}
{"id": "entertainment_2", "amount": 105}
{"id": "commute_1", "amount": 15}
```

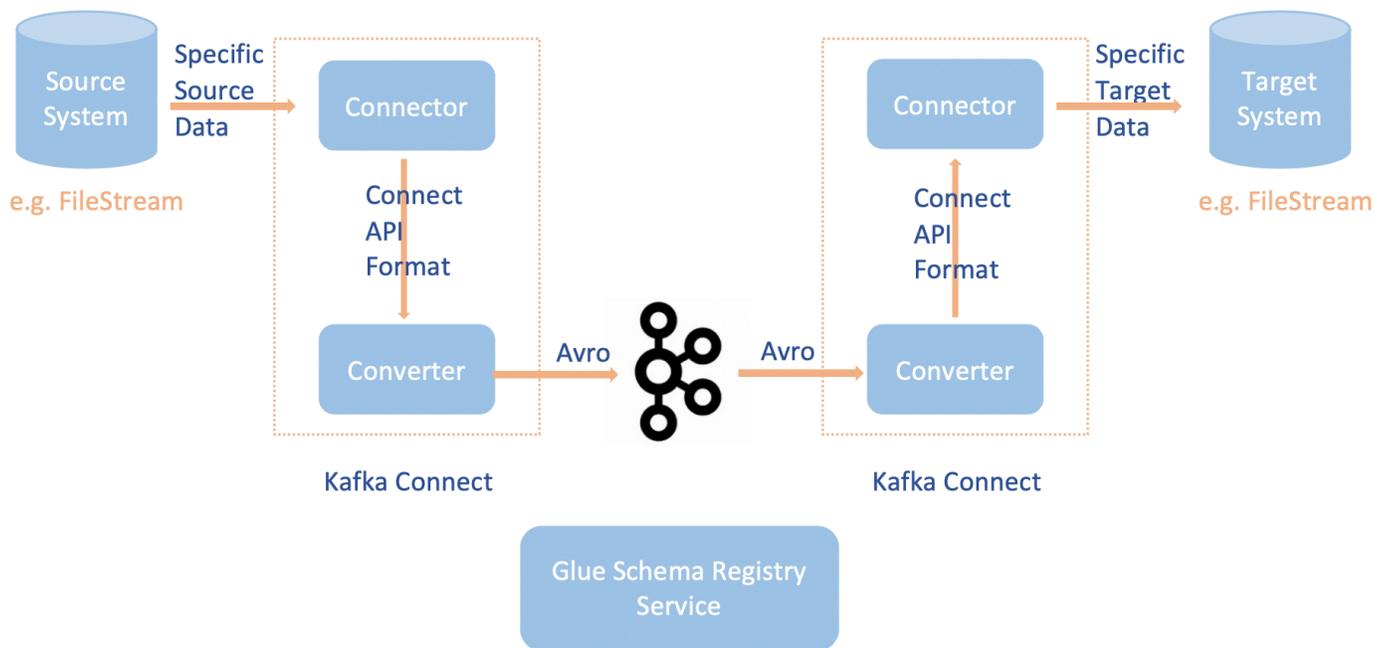
フィルタリング後のレコード:

```
{"name": "Sansa", "favorite_number": 99, "favorite_color": "white"}
{"name": "Harry", "favorite_number": 10, "favorite_color": "black"}
{"name": "Hermione", "favorite_number": 1, "favorite_color": "red"}
{"name": "Ron", "favorite_number": 0, "favorite_color": "pink"}

{"id": "commute_1", "amount": 3.5}
{"id": "grocery_1", "amount": 25.5}
{"id": "entertainment_1", "amount": 19.2}
{"id": "entertainment_2", "amount": 105}
```

ユースケース: Apache Kafka Connect

Apache Kafka Connect と AWS Glue Schema Registry を統合することで、コネクタからスキーマ情報を取得できるようになります。Apache Kafka のコンバータにより、Apache Kafka 内のデータ形式と、Apache Kafka Connect データへの変換方法を指定します。すべての Apache Kafka Connect ユーザーは、これらのコンバータを Apache Kafka との間でロードまたは保存する際に、データに適用する形式に基づいた設定を行う必要があります。これにより、Apache Kafka Connect データを AWS Glue Schema Registry で使用する型 (例: Avro) に変換する独自のコンバータを定義し、さらにシリアライザを使用してスキーマを登録しシリアル化を実行します。その後コンバータはデシリアライザを使用して、Apache Kafka から受信したデータを逆シリアル化し、元の Apache Kafka Connect データに変換することができます。ワークフローの例を以下の図に示します。



1. [AWS Glue Schema Registry 用 Githubリポジトリ](#) をクローンして、aws-glue-schema-registry プロジェクトをインストールします。

```
git clone git@github.com:aws-labs/aws-glue-schema-registry.git
cd aws-glue-schema-registry
mvn clean install
mvn dependency:copy-dependencies
```

2. Apache Kafka Connect を Standalone モードで使用する予定の場合、このステップで以下に示した手順を使用して、connect-standalone.properties を更新します。Apache Kafka Connect を Distributed モードで使用する予定の場合は、同じ手順により connect-avro-distributed.properties を更新します。

a. Apache Kafka の接続プロパティファイルにも、これらのプロパティを追加します。

```
key.converter.region=aws-region
value.converter.region=aws-region
key.converter.schemaAutoRegistrationEnabled=true
value.converter.schemaAutoRegistrationEnabled=true
key.converter.avroRecordType=GENERIC_RECORD
value.converter.avroRecordType=GENERIC_RECORD
```

- b. 以下のコマンドを、[kafka-run-class.sh] の下の [Launch mode] (起動モード) セクションに追加します。

```
-cp $CLASSPATH:"<your AWS GlueSchema Registry base directory>/target/dependency/*"
```

3. [kafka-run-class.sh] の下の [Launch mode] (起動モード) セクションに以下のコマンドを追加する

```
-cp $CLASSPATH:"<your AWS GlueSchema Registry base directory>/target/dependency/*"
```

プリンシパルは以下のようになります。

```
# Launch mode
if [ "x$DAEMON_MODE" = "xtrue" ]; then
  nohup "$JAVA" $KAFKA_HEAP_OPTS $KAFKA_JVM_PERFORMANCE_OPTS $KAFKA_GC_LOG_OPTS
  $KAFKA_JMX_OPTS $KAFKA_LOG4J_OPTS -cp $CLASSPATH:"/Users/johndoe/aws-glue-schema-
  registry/target/dependency/*" $KAFKA_OPTS "$@" > "$CONSOLE_OUTPUT_FILE" 2>&1 < /dev/
  null &
else
  exec "$JAVA" $KAFKA_HEAP_OPTS $KAFKA_JVM_PERFORMANCE_OPTS $KAFKA_GC_LOG_OPTS
  $KAFKA_JMX_OPTS $KAFKA_LOG4J_OPTS -cp $CLASSPATH:"/Users/johndoe/aws-glue-schema-
  registry/target/dependency/*" $KAFKA_OPTS "$@"
fi
```

4. bash を使用している場合は、以下のコマンドを実行して、bash_profile で CLASSPATH を設定します。他のシェルの場合は、それに応じて環境を更新します。

```
echo 'export GSR_LIB_BASE_DIR=<>' >> ~/.bash_profile
echo 'export GSR_LIB_VERSION=1.0.0' >> ~/.bash_profile
echo 'export KAFKA_HOME=<your Apache Kafka installation directory>' >> ~/.bash_profile
echo 'export CLASSPATH=$CLASSPATH:$GSR_LIB_BASE_DIR/avro-kafkaconnect-converter/
target/schema-registry-kafkaconnect-converter-$GSR_LIB_VERSION.jar:$GSR_LIB_BASE_DIR/
common/target/schema-registry-common-$GSR_LIB_VERSION.jar:$GSR_LIB_BASE_DIR/
avro-serializer-deserializer/target/schema-registry-serde-$GSR_LIB_VERSION.jar'
>> ~/.bash_profile
source ~/.bash_profile
```

5. (オプション) 単純なファイルをソースとして使用しテストを行う場合は、ファイルソースコネクタのクローンを作成します。

```
git clone https://github.com/mmolimar/kafka-connect-fs.git
```

```
cd kafka-connect-fs/
```

- a. ソースコネクタの設定で、データ形式を Avro に、ファイルリーダーを AvroFileReader に変更します。さらに、読み込んでいるファイルパスからサンプルの Avro オブジェクトを更新します。例:

```
vim config/kafka-connect-fs.properties
```

```
fs.uris=<path to a sample avro object>  
policy.regex=^.*\.avro$  
file_reader.class=com.github.mmolimar.kafka.connect.fs.file.reader.AvroFileReader
```

- b. ソースコネクタをインストールします。

```
mvn clean package  
echo "export CLASSPATH=\$CLASSPATH:\\"$(find target/ -type f -name '*.jar'| grep  
'\-package' | tr '\n' ':')\\"" >> ~/.bash_profile  
source ~/.bash_profile
```

- c. *<your Apache Kafka installation directory>*/config/connect-file-sink.properties のシンクのプロパティを更新し、トピック名と出力ファイル名を更新します。

```
file=<output file full path>  
topics=<my topic>
```

6. Source Connector (この例では、ソースファイルのコネクタ) を起動します。

```
$KAFKA_HOME/bin/connect-standalone.sh $KAFKA_HOME/config/connect-standalone.properties config/kafka-connect-fs.properties
```

7. Sink Connector (この例では、シンクファイルのコネクタ) を実行します。

```
$KAFKA_HOME/bin/connect-standalone.sh $KAFKA_HOME/config/connect-standalone.properties $KAFKA_HOME/config/connect-file-sink.properties
```

Kafka Connect の使用例については、[AWS Glue Schema Registry 用 Githubリポジトリ](#)の、integration-tests フォルダにある run-local-tests.sh スクリプトでご確認ください。

サードパーティー製のスキーマレジストリから AWS Glue Schema Registry への移行

サードパーティー製スキーマレジストリから AWS Glue Schema Registry への移行作業は、その時点での既存のサードパーティー製スキーマレジストリにより異なります。サードパーティーのスキーマレジストリを使用して送信された Apache Kafka トピック内にレコードがある場合、コンシューマはサードパーティーのスキーマレジストリを使用して、それらのレコードを非シリアル化する必要があります。AWSKafkaAvroDeserializer には、セカンダリのデシリアライザクラスを定義する機能があります。このクラスは、サードパーティーのデシリアライザを特定し、対象のレコードを非シリアル化するために使用できます、

サードパーティー製スキーマには、使用停止に関して 2 つの基準があります。1 つ目の基準では、サードパーティー製スキーマレジストリを使用する Apache Kafka トピックのレコードを必要とするコンシューマが、なくなった後にのみ使用停止となります。2 つ目の基準では、トピックに指定された保持期間に応じ、Apache Kafka トピックの試用期間が終了することで使用停止となります。無制限の保持期間を持つトピックの場合は、AWS Glue Schema Registry への移行は可能ですが、サードパーティー製スキーマレジストリを使用停止にすることはできません。この回避策としては、アプリケーションまたは Mirror Maker 2 により現在のトピックを読み取り、AWS Glue Schema Registry を使用する新しいトピックとして作成し直すことができます。

サードパーティー製スキーマレジストリから AWS Glue Schema Registry へ移行するには

1. AWS Glue Schema Registry にレジストリを作成するか、デフォルトのレジストリを使用します。
2. コンシューマを停止します。AWS Glue Schema Registry をプライマリのデシリアライザとして含めるように、コンシューマを変更します、サードパーティーのスキーマレジストリをセカンダリに変更します。
 - コンシューマのプロパティを設定します。この例では、secondary_deccerializer は別のデシリアライザに設定されています。動作は次のとおりです。コンシューマは Amazon MSK からレコードを取得し、最初に AWSKafkaAvroDeserializer の使用を試みます。AWS Glue Schema Registry のスキーマの、Avro スキーマ ID を含むマジックバイトを読み取ることができない場合、AWSKafkaAvroDeserializer は、secondary_deserializer で指定されるデシリアライザクラスの使用を試みます。セカンダリのデシリアライザに固有のプロパティは、以下に示すように、schema_registry_url_config や specific_avro_reader_config などのコンシューマプロパティでも指定する必要があります。

```
consumerProps.setProperty(ConsumerConfig.KEY_DESERIALIZER_CLASS_CONFIG,
    StringDeserializer.class.getName());
consumerProps.setProperty(ConsumerConfig.VALUE_DESERIALIZER_CLASS_CONFIG,
    AWSKafkaAvroDeserializer.class.getName());
consumerProps.setProperty(AWSSchemaRegistryConstants.AWS_REGION,
    KafkaClickstreamConsumer.gsrRegion);
consumerProps.setProperty(AWSSchemaRegistryConstants.SECONDARY_DESERIALIZER,
    KafkaAvroDeserializer.class.getName());
consumerProps.setProperty(KafkaAvroDeserializerConfig.SCHEMA_REGISTRY_URL_CONFIG,
    "URL for third-party schema registry");
consumerProps.setProperty(KafkaAvroDeserializerConfig.SPECIFIC_AVRO_READER_CONFIG,
    "true");
```

3. コンシューマを再起動します。
4. プロデューサを停止し、そのプロデューサで AWS Glue Schema Registry を指定します。
 - a. プロデューサのプロパティを設定します。この例のプロデューサでは、デフォルトのレジストリと自動登録スキーマバージョンを使用しています。

```
producerProps.setProperty(ProducerConfig.KEY_SERIALIZER_CLASS_CONFIG,
    StringSerializer.class.getName());
producerProps.setProperty(ProducerConfig.VALUE_SERIALIZER_CLASS_CONFIG,
    AWSKafkaAvroSerializer.class.getName());
producerProps.setProperty(AWSSchemaRegistryConstants.AWS_REGION, "us-east-2");
producerProps.setProperty(AWSSchemaRegistryConstants.AVRO_RECORD_TYPE,
    AvroRecordType.SPECIFIC_RECORD.getName());
producerProps.setProperty(AWSSchemaRegistryConstants.SCHEMA_AUTO_REGISTRATION_SETTING,
    "true");
```

5. (オプション) 既存のスキーマおよびスキーマバージョンを、現在のサードパーティー製スキーマレジストリから AWS Glue Schema Registry (AWS Glue Schema Registry のデフォルトレジストリ、または AWS Glue Schema Registry 内の特定の非デフォルトレジストリ) に手動で移行します。これには、サードパーティーのスキーマレジストリから JSON 形式でスキーマをエクスポートし、AWS Management Console または AWS CLI を使用しながら AWS Glue Schema Registry の新しいスキーマを作成します。

このステップは、AWS CLI および AWS Management Console を使用して、新しく作成されたスキーマバージョンについて、以前のスキーマバージョンとの互換性チェックを有効にする必要がある場合、またはスキーマバージョンの自動登録を有効にした新しいスキーマで、プロデューサがメッセージを送信する場合に重要です。

6. プロデューサを起動します。

データへの接続

AWS Glue 接続は、特定のデータストアのログイン認証情報、URI 文字列、仮想プライベートクラウド (VPC) 情報などを保存する Data Catalog オブジェクトです。AWS Glue クローラー、ジョブ、開発エンドポイントは、特定のタイプのデータストアにアクセスするために接続を使用します。ソースとターゲットの両方に接続を使用したり、複数のクローラーまたは抽出、変換、ロード (ETL) ジョブで同じ接続を再利用したりできます。

AWS Glue では、次の接続タイプがサポートされています。

- Amazon DocumentDB
- AWS Glue for Spark で使用する Amazon OpenSearch Service。
- Amazon Redshift
- Azure Cosmos、AWS Glue ETL ジョブで Azure Cosmos DB for NoSQL を使用
- AWS Glue for Spark で使用する Azure SQL。
- for Spark BigQueryで使用する AWS Glue Google 。
- JDBC
- Kafka
- MongoDB
- MongoDB Atlas
- Salesforce
- for Spark で使用する AWS Glue SAP HANA。
- for Spark で使用する AWS Glue Snowflake。
- for Spark を使用する場合 AWS Glue の Teradata Vantage。
- AWS Glue for Spark で使用するための Vertica。
- さまざまな Amazon Relational Database Service (Amazon RDS) オフアリング。
- Network (Amazon Virtual Private Cloud (Amazon VPC) 内のデータソースへの接続を指定します)
- Aurora (ネイティブ JDBC ドライバーを使用している場合はサポートされます。すべてのドライバー機能を利用できるわけではありません)

AWS Glue Studio を使用して、コネクタの接続を作成することもできます。コネクタとは AWS Glue Studio内でデータストアに対するアクセスを支援するための、オプションのコードパッケージです。

詳細については、「[Using connectors and connections with AWS Glue Studio](#)」を参照してください。

オンプレミスデータベースに接続する方法については、AWS Big Data Blog [ウェブサイトの「を使用してオンプレミスのデータストアにアクセスして分析する方法AWS Glue」](#)を参照してください。

このセクションには、AWS Glue 接続の使用に役立つ次のトピックが含まれています。

- [AWS Glue 接続プロパティ](#)
- [AWS Secrets Manager への接続認証情報の保存](#)
- [AWS Glue 接続の追加](#)
- [AWS Glue 接続のテスト](#)
- [すべての AWS コールを VPC を経由するように設定する](#)
- [VPC の JDBC データストアに接続する](#)
- [MongoDB または MongoDB Atlas 接続を使用する](#)
- [VPC エンドポイントを使用して Amazon S3 データストアをクローलする](#)
- [AWS Glue での接続の問題のトラブルシューティング](#)
- [チュートリアル: Elasticsearch 向けの AWS Glue コネクタを使用する](#)

AWS Glue 接続プロパティ

このトピックには、AWS Glue 接続のプロパティに関する情報が含まれています。

トピック

- [必要な接続プロパティ](#)
- [AWS Glue JDBC 接続プロパティ](#)
- [AWS Glue MongoDB と MongoDB Atlas の接続プロパティ](#)
- [Salesforce 接続プロパティ](#)
- [Snowflake 接続](#)
- [Vertica 接続](#)
- [SAP HANA 接続](#)
- [Azure SQL 接続](#)
- [Teradata Vantage 接続](#)

- [OpenSearch サービス接続](#)
- [Azure Cosmos 接続](#)
- [AWS Glue SSL 接続プロパティ](#)
- [クライアント認証用の Apache Kafka 接続プロパティ](#)
- [Google BigQuery 接続](#)
- [Vertica 接続](#)

必要な接続プロパティ

AWS Glue コンソールで設定を定義する場合は、次のプロパティに値を指定する必要があります。

接続名

接続用の一意の名前を入力します。

接続タイプ

JDBC またはいずれかの接続タイプを選択します。

JDBC 接続タイプの詳細については、「[the section called “JDBC 接続プロパティ”](#)」を参照してください。

Amazon Virtual Private Cloud 環境 (Amazon VPC) 内のデータソースへ接続するには、Network を選択します。

選択したタイプに応じて、AWS Glue コンソールにその他の必須フィールドが表示されます。例えば、Amazon RDS を選択した場合は、データベースエンジンを選択する必要があります。

要求される SSL 接続

このオプションを選択すると、AWS Glue では、データストア接続が信頼できる Secure Sockets Layer (SSL) 経由であることの確認が必要です。

このオプションを選択したときに利用できる追加オプションなどの詳細については、「[the section called “SSL 接続プロパティ”](#)」を参照してください。

MSK クラスターを選択 (Amazon Managed Streaming for Apache Kafka (MSK) のみ)

別の AWS アカウントの MSK クラスターを指定します。

Kafka ブートストラップサーバーの URL (Kafka のみ)

ブートストラップサーバーの URL のカンマ区切りリストを指定します。ポート番号を含めます。例: b-1.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094, b-2.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094, b-3.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094

AWS Glue JDBC 接続プロパティ

AWS Glue は、JDBC 接続を介して次のデータストアに接続できます。

- Amazon Redshift
- Amazon Aurora
- Microsoft SQL Server
- MySQL
- Oracle
- PostgreSQL
- AWS Glue クローラーを使用する場合の Snowflake。
- Aurora (ネイティブ JDBC ドライバーを使用している場合はサポートされます。すべてのドライバー機能を利用できるわけではありません)
- Amazon RDS for MariaDB

Important

現在、ETL ジョブは 1 つのサブネット内ではしか JDBC 接続を使用できません。1 つのジョブに複数のデータストアがある場合は、同じサブネットにあるか、サブネットからアクセス可能である必要があります。

AWS Glue クローラー用に独自の JDBC ドライバーバージョンを導入する場合、クローラーは AWS Glue ジョブと Amazon S3 のリソースを使用して、用意したドライバーが自分の環境で実行されるようにします。リソースの追加使用量はアカウントに反映されます。さらに、独自の JDBC ドライバーを用意しても、クローラーがドライバーの機能をすべて活用できるわけではありません。ドライバーは、「[Data Catalog での接続の定義](#)」に記載されているプロパティに限定されます。

JDBC 接続タイプの追加プロパティを次に示します。

JDBC URL

JDBC データストアの URL を入力します。ほとんどのデータベースエンジンの場合、このフィールドは次の形式になります。この形式では、*protocol*、*host*、*port*、および *db_name* を独自の情報に置き換えます。

```
jdbc:protocol://host:port/db_name
```

データベースエンジンに応じて、別の JDBC URL の形式が必要な場合があります。この形式では、コロン (:) とスラッシュ (/) の使用方法が若干異なるか、データベースを指定するためのキーワードが異なる場合があります。

JDBC をデータストアに接続するためには、データストアの *db_name* が必要です。*db_name* は、指定された *username* と *password* を使用してネットワーク接続を確立するために使用されます。接続すると、AWS Glue はデータストア内の他のデータベースにアクセスして、クローラーを実行したり ETL ジョブを実行したりできます。

次の JDBC URL の例は、いくつかのデータベースエンジンの構文を示しています。

- dev データベースを使用して Amazon Redshift クラスターデータストアに接続するには

```
jdbc:redshift://xxx.us-east-1.redshift.amazonaws.com:8192/dev
```

- employee データベースを使用して Amazon RDS for MySQL データストアに接続するには

```
jdbc:mysql://xxx-cluster.cluster-xxx.us-east-1.rds.amazonaws.com:3306/  
employee
```

- employee データベースを使用して Amazon RDS for PostgreSQL データストアに接続するには

```
jdbc:postgresql://xxx-cluster.cluster-xxx.us-  
east-1.rds.amazonaws.com:5432/employee
```

- employee サービス名を使用して Amazon RDS for Oracle データストアに接続するには

```
jdbc:oracle:thin://@xxx-cluster.cluster-xxx.us-  
east-1.rds.amazonaws.com:1521/employee
```

Amazon RDS for Oracle の構文は次のパターンに従います。これらのパターンでは、*host*、*port*、*service_name*、*SID* を独自の情報に置き換えます。

- `jdbc:oracle:thin://@host:port/service_name`
- `jdbc:oracle:thin://@host:port:SID`
- employee データベースを使用して Amazon RDS for Microsoft SQL Server データストアに接続するには

```
jdbc:sqlserver://xxx-cluster.cluster-xxx.us-east-1.rds.amazonaws.com:1433;databaseName=employee
```

Amazon RDS for SQL Server の構文は次のパターンに従います。これらのパターンでは、`server_name`、`port`、および `db_name` を独自の情報に置き換えます。

- `jdbc:sqlserver://server_name:port;database=db_name`
- `jdbc:sqlserver://server_name:port;databaseName=db_name`
- employee データベースの Amazon Aurora PostgreSQL インスタンスに接続するには、データベースインスタンスのエンドポイント、ポート、データベース名を指定します。

```
jdbc:postgresql://employee_instance_1.xxxxxxxxxxxxxx.us-east-2.rds.amazonaws.com:5432/employee
```

- employee データベースを使用して Amazon RDS for MariaDB データストアに接続するには、データベースインスタンスのエンドポイント、ポート、データベース名を指定します。

```
jdbc:mysql://xxx-cluster.cluster-xxx.aws-region.rds.amazonaws.com:3306/employee
```

Warning

Snowflake JDBC 接続は AWS Glue クローラーでのみサポートされています。AWS Glue ジョブで Snowflake コネクタを使用する場合は、Snowflake 接続タイプを使用します。

sample データベースのインスタンスに接続するには、Snowflake インスタンスのエンドポイント、ユーザー、データベース名、およびロール名を指定します。必要に応じて warehouse パラメータを追加できます。

```
jdbc:snowflake://account_name.snowflakecomputing.com/?user=user_name&db=sample&role=role_name&warehouse=warehouse_name
```

⚠ Important

JDBC 経由の Snowflake 接続では、URL 内のパラメータは、`user`、`db`、`role_name`、`warehouse` の順序にする必要があります。

- AWS プライベートリンクを使用して sample データベースの Snowflake インスタンスに接続するには、次のように Snowflake JDBC URL を指定します。

```
jdbc:snowflake://account_name.region.privatelink.snowflakecomputing.com/?  
user=user_name&db=sample&role=role_name&warehouse=warehouse_name
```

ユーザーネーム

i Note

ユーザー名とパスワードを直接入力するのではなく、AWS シークレットを使用して接続認証情報を保存することをお勧めします。詳細については、「[AWS Secrets Manager への接続認証情報の保存](#)」を参照してください。

JDBC データストアにアクセスする権限を持つユーザー名を指定します。

パスワード

JDBC データストアへのアクセス権限を持つユーザー名のパスワードを入力します。

[ポート]

JDBC URL で使用するポートを入力し Amazon RDS Oracle インスタンスに接続します。このフィールドは、Amazon RDS Oracle インスタンスに対して [Require SSL connection] (SSL 接続が必要) を選択したときにしか表示されません。

VPC

データストアを含む Virtual Private Cloud (VPC) の名前を選択します。AWS Glue コンソールには、現在のリージョンの VPC がすべて表示されます。

⚠ Important

Snowflake からのデータなど AWS、 からホストされている JDBC 接続で作業する場合、VPC にはトラフィックをパブリックサブネットとプライベートサブネットに分割す

る NAT ゲートウェイが必要です。パブリックサブネットは外部ソースへの接続に使用され、内部サブネットは による処理に使用されます AWS Glue。Amazon VPC を外部接続用に設定する方法については、「[NAT デバイスを使用してインターネットまたは他のネットワークに接続する](#)」と「[AWS Glue から Amazon RDS データストアに JDBC 接続するための Amazon VPC の設定](#)」を参照してください。

サブネット

データストアを含む VPC 内のサブネットを選択します。AWS Glue コンソールには VPC 内のデータストアのすべてのサブネットが一覧表示されます。

セキュリティグループ

データストアに関連付けられているセキュリティグループを選択します。AWS Glue には、AWS Glue の接続を許可するインバウンドソースルールを持つ 1 つ以上のセキュリティグループが必要です。AWS Glue コンソールには、VPC へのインバウンドアクセスが許可されているすべてのセキュリティグループが一覧表示されます。AWS Glue は、VPC サブネットにアタッチされている Elastic Network Interface にこれらのセキュリティグループを関連付けます。

JDBC ドライバークラス名 - オプション

次のカスタム JDBC ドライバークラス名を指定します。

- Postgres – org.postgresql.Driver
- MySQL – com.mysql.jdbc.Driver、com.mysql.cj.jdbc.Driver
- Redshift – com.amazon.redshift.jdbc.Driver、com.amazon.redshift.jdbc42.Driver
- Oracle – oracle.jdbc.driver。OracleDriver
- SQL Server – com.microsoft.sqlserver.jdbc.SQLServerDriver

JDBC ドライバース3パス - オプション

カスタム JDBC ドライバース3パスに対して Amazon S3 の場所を指定します。これは.jar ファイルへの絶対パスです。クローラーがサポートするデータベースのデータソースに接続するための独自の JDBC ドライバース3パスを用意したい場合は、customJdbcDriverS3Path パラメータと customJdbcDriverClassName パラメータの値を指定できます。

お客様が用意した JDBC ドライバーの使用は、必要な [必要な接続プロパティ](#) に限られます。

AWS Glue MongoDB と MongoDB Atlas の接続プロパティ

MongoDB または MongoDB Atlas 接続タイプの追加のプロパティを次に示します。

MongoDB URL

お使いの MongoDB または MongoDB Atlas データストアの URL を入力します。

- MongoDB の場合: `mongodb://host:port/database` ホストは、ホスト名、IP アドレス、UNIX ドメインソケットのいずれかにすることができます。接続文字列でポートが指定されていない場合は、デフォルトの MongoDB ポート、27017 が使用されます。
- MongoDB Atlas: `mongodb+srv://server.example.com/database` ホストは、DNS SRV レコードに対応するホスト名にすることができます。SRV 形式ではポートは不要であり、デフォルトの MongoDB ポート、27017 が使用されます。

ユーザーネーム

Note

ユーザー名とパスワードを直接入力するのではなく、AWS シークレットを使用して接続認証情報を保存することをお勧めします。詳細については、「[AWS Secrets Manager への接続認証情報の保存](#)」を参照してください。

JDBC データストアにアクセスする権限を持つユーザー名を指定します。

パスワード

MongoDB または MongoDB Atlas へのアクセス権限を持つユーザー名のパスワードを入力します。

Salesforce 接続プロパティ

Salesforce 接続タイプの追加のプロパティを次に示します。

- ENTITY_NAME (文字列) - (必須) 読み取り/書き込みに使用されます。Salesforce のオブジェクトの名前。

- `API_VERSION` (文字列) - (必須) 読み取り/書き込みに使用されます。使用する Salesforce Rest API バージョン。
- `SELECTED_FIELDS`(List<String>) - デフォルト: empty(SELECT *)。読み込みに使用されます。オブジェクトに選択する列。
- `FILTER_PREDICATE` (文字列) - デフォルト: 空。読み込みに使用されます。Spark SQL 形式である必要があります。
- `QUERY` (文字列) - デフォルト: 空。読み込みに使用されます。完全な Spark SQL クエリ。
- `PARTITION_FIELD` (文字列) - 読み取りに使用されます。クエリのパーティション化に使用されるフィールド。
- `LOWER_BOUND` (文字列) - 読み取りに使用されます。選択したパーティションフィールドの包括的な下限値。
- `UPPER_BOUND` (文字列) - 読み取りに使用されます。選択したパーティションフィールドの排他的上限値。
- `NUM_PARTITIONS` (整数) - デフォルト: 1。読み込みに使用されます。読み取り用のパーティションの数。
- `IMPORT_DELETED_RECORDS` (文字列) - デフォルト: FALSE。読み取りに使用されます。クエリ中に削除レコードを取得するには。
- `WRITE_OPERATION` (文字列) - デフォルト: INSERT。書き込みに使用されます。値は INSERT、UPDATE、UPSERT、DELETE である必要があります。
- `ID_FIELD_NAMES` (文字列) - デフォルト: null。UPSERT でのみ使用されます。

Snowflake 接続

AWS Glue ETL ジョブで使用される Snowflake 接続をセットアップするには、次のプロパティを使用します。Snowflake をクローリングする場合は、JDBC 接続を使用します。

Snowflake URL

Snowflake エンドポイントの URL。Snowflake エンドポイントの URL の詳細については、Snowflake ドキュメントの「[アカウントへの接続](#)」を参照してください。

AWS シークレット

のシークレットのシークレット名 AWS Secrets Manager。AWS Glue は、シークレットの `sfUser` および `sfPassword` キーを使用して Snowflake に接続します。

Snowflake ロール (オプション)

Snowflake セキュリティロール AWS Glue は、接続時に を使用します。

AWS PrivateLinkを使用して、Amazon VPC でホストされている Snowflake エンドポイントへの接続を設定する場合は、次のプロパティを使用します。

VPC

データストアを含むVirtual Private Cloud (VPC) の名前を選択します。AWS Glue コンソールには、現在のリージョンの VPC がすべて表示されます。

サブネット

データストアを含む VPC 内のサブネットを選択します。AWS Glue コンソールには VPC 内のデータストアのすべてのサブネットが一覧表示されます。

セキュリティグループ

データストアに関連付けられているセキュリティグループを選択します。AWS Glue には、AWS Glue の接続を許可するインバウンドソースルールを持つ 1 つ以上のセキュリティグループが必要です。AWS Glue コンソールには、VPC へのインバウンドアクセスが許可されているすべてのセキュリティグループが一覧表示されます。AWS Glue は、VPC サブネットにアタッチされている Elastic Network Interface にこれらのセキュリティグループを関連付けます。

Vertica 接続

次のプロパティを使用して、AWS Glue ETL ジョブの Vertica 接続を設定します。

Vertica ホスト

Vertica インストールのホスト名。

Vertica ポート

Vertica インストールに使用できるポート。

AWS シークレット

のシークレットのシークレット名 AWS Secrets Manager。AWS Glue は、シークレットのキーを使用して Vertica に接続します。

Amazon VPC でホストされている Vertica エンドポイントへの接続を設定する場合は、次のプロパティを使用します。

VPC

データストアを含むVirtual Private Cloud (VPC) の名前を選択します。AWS Glue コンソールには、現在のリージョンの VPC がすべて表示されます。

サブネット

データストアを含む VPC 内のサブネットを選択します。AWS Glue コンソールには VPC 内のデータストアのすべてのサブネットが一覧表示されます。

セキュリティグループ

データストアに関連付けられているセキュリティグループを選択します。AWS Glue には、AWS Glue の接続を許可するインバウンドソースルールを持つ 1 つ以上のセキュリティグループが必要です。AWS Glue コンソールには、VPC へのインバウンドアクセスが許可されているすべてのセキュリティグループが一覧表示されます。AWS Glue は、VPC サブネットにアタッチされている Elastic Network Interface にこれらのセキュリティグループを関連付けます。

SAP HANA 接続

AWS Glue ETL ジョブの SAP HANA 接続を設定するには、次のプロパティを使用します。

SAP HANA URL

SAP JDBC URL。

SAP HANA JDBC URL の形式は

`jdbc:sap://saphanaHostname:saphanaPort?databaseName=saphanaDBname,ParameterName`
です

AWS Glue には、次の JDBC URL パラメータが必要です。

- `databaseName` – 接続先の SAP HANA のデフォルトデータベース。

AWS シークレット

のシークレットのシークレット名 AWS Secrets Manager。AWS Glue は、シークレットのキーを使用して SAP HANA に接続します。

Amazon VPC でホストされている SAP HANA エンドポイントへの接続を設定する場合は、次のプロパティを使用します。

VPC

データストアを含むVirtual Private Cloud (VPC) の名前を選択します。AWS Glue コンソールには、現在のリージョンの VPC がすべて表示されます。

サブネット

データストアを含む VPC 内のサブネットを選択します。AWS Glue コンソールには VPC 内のデータストアのすべてのサブネットが一覧表示されます。

セキュリティグループ

データストアに関連付けられているセキュリティグループを選択します。AWS Glue には、AWS Glue の接続を許可するインバウンドソースルールを持つ 1 つ以上のセキュリティグループが必要です。AWS Glue コンソールには、VPC へのインバウンドアクセスが許可されているすべてのセキュリティグループが一覧表示されます。AWS Glue は、VPC サブネットにアタッチされている Elastic Network Interface にこれらのセキュリティグループを関連付けます。

Azure SQL 接続

次のプロパティを使用して、AWS Glue ETL ジョブの Azure SQL 接続を設定します。

Azure SQL URL

Azure SQL エンドポイントの JDBC URL。

URL は、次のような形式になります:

```
jdbc:sqlserver://databaseServerName:databasePort;databaseName=azuresqlDBname;
```

AWS Glue には、次の URL プロパティが必要です。

- `databaseName` – 接続先の Azure SQL のデフォルトデータベース。

Azure SQL Managed Instances の JDBC URL の詳細については、[Microsoft のドキュメント](#)を参照してください。

AWS シークレット

のシークレットのシークレット名 AWS Secrets Manager。AWS Glue はシークレットのキーを使用して Azure SQL に接続します。

Teradata Vantage 接続

次のプロパティを使用して、AWS Glue ETL ジョブの Teradata Vantage 接続を設定します。

Teradata URL

Teradata インスタンスに接続するには、データベースインスタンスのホスト名と、関連する Teradata パラメータを指定します。

```
jdbc:teradata://teradataHostname/ParameterName=ParameterValue,ParameterName=Pa
```

AWS Glue は、次の JDBC URL パラメータをサポートします。

- DATABASE_NAME – 接続先の Teradata のデフォルトデータベース。
- DBS_PORT – 標準でない場合は、Teradata ポートを指定します。

AWS シークレット

のシークレットのシークレット名 AWS Secrets Manager。AWS Glue は、シークレットのキーを使用して Teradata Vantage に接続します。

Amazon VPC でホストされている Teradata Vantage エンドポイントへの接続を設定する場合は、次のプロパティを使用します。

VPC

データストアを含む Virtual Private Cloud (VPC) の名前を選択します。AWS Glue コンソールには、現在のリージョンの VPC がすべて表示されます。

サブネット

データストアを含む VPC 内のサブネットを選択します。AWS Glue コンソールには VPC 内のデータストアのすべてのサブネットが一覧表示されます。

セキュリティグループ

データストアに関連付けられているセキュリティグループを選択します。AWS Glue には、AWS Glue の接続を許可するインバウンドソースルールを持つ 1 つ以上のセキュリティグループが必要です。AWS Glue コンソールには、VPC へのインバウンドアクセスが許可されているすべてのセキュリティグループが一覧表示されます。AWS Glue は、VPC サブネットにアタッチされている Elastic Network Interface にこれらのセキュリティグループを関連付けます。

OpenSearch サービス接続

AWS Glue ETL ジョブのサービス OpenSearch 接続を設定するには、次のプロパティを使用します。

ドメインエンドポイント

Amazon OpenSearch Service ドメインエンドポイントのデフォルト形式は、`https://search-domainName-unstructuredIdContent.region` . です。 `es.amazonaws.com` ドメインエンドポイントの識別の詳細については、[Amazon OpenSearch Service ドキュメントの「Amazon Service ドメインの作成と管理 OpenSearch」](#)を参照してください。

[ポート]

エンドポイントでポートが開きます。

AWS シークレット

のシークレットのシークレット名 AWS Secrets Manager。AWS Glue は、シークレットのキーを使用して OpenSearch サービスに接続します。

Amazon VPC でホストされている OpenSearch サービスエンドポイントへの接続を設定するときは、次のプロパティを使用します。

VPC

データストアを含むVirtual Private Cloud (VPC) の名前を選択します。AWS Glue コンソールには、現在のリージョンの VPC がすべて表示されます。

サブネット

データストアを含む VPC 内のサブネットを選択します。AWS Glue コンソールには VPC 内のデータストアのすべてのサブネットが一覧表示されます。

セキュリティグループ

データストアに関連付けられているセキュリティグループを選択します。AWS Glue には、AWS Glue の接続を許可するインバウンドソースルールを持つ 1 つ以上のセキュリティグループが必要です。AWS Glue コンソールには、VPC へのインバウンドアクセスが許可されているすべてのセキュリティグループが一覧表示されます。AWS Glue は、VPC サブネットにアタッチされている Elastic Network Interface にこれらのセキュリティグループを関連付けます。

Azure Cosmos 接続

次のプロパティを使用して、AWS Glue ETL ジョブの Azure Cosmos 接続を設定します。

Azure Cosmos DB アカウントのエンドポイント URI

Azure Cosmos への接続に使用するエンドポイント。詳細については、[Azure のドキュメント](#)を参照してください。

AWS シークレット

のシークレットのシークレット名 AWS Secrets Manager。AWS Glue は、シークレットのキーを使用して Azure Cosmos に接続します。

AWS Glue SSL 接続プロパティ

[Require SSL connection] (SSL 接続が必要) プロパティの詳細は、次の通りです。

SSL 接続が必要でない場合、AWS Glue では SSL を使用してデータストアへの接続を暗号化できなくても無視されます。設定の手順については、データストアのドキュメントを参照してください。このオプションを選択すると、AWS Glue が接続できない場合、開発エンドポイントでのジョブ実行、クローラー、または ETL のステートメントは失敗します。

Note

Snowflake はデフォルトで SSL 接続をサポートしているため、このプロパティは Snowflake には適用されません。

このオプションは、AWS Glue クライアント側で検証されます。JDBC 接続の場合、AWS Glue は、証明書とホスト名の検証を使用する SSL 経由でのみ接続を行います。SSL 接続サポートは、次の場合に利用できます。

- Oracle Database
- Microsoft SQL Server
- PostgreSQL
- Amazon Redshift
- MySQL (Amazon RDS インスタンスのみ)

- Amazon Aurora MySQL (Amazon RDS インスタンスのみ)
- Amazon Aurora PostgreSQL (Amazon RDS インスタンスのみ)
- Kafka、以下が含まれます。 Amazon Managed Streaming for Apache Kafka
- MongoDB

Note

Amazon RDS Oracle データストアが [Require SSL connection] (SSL 接続が必要) を使用できるようにするためには、Oracle インスタンスへのオプショングループを作成してアタッチする必要があります。

1. にサインイン AWS Management Console し、<https://console.aws.amazon.com/rds/> で Amazon RDS コンソールを開きます。
2. オプショングループを Amazon RDS Oracle インスタンスに追加します。Amazon RDS コンソールで、オプショングループを追加する方法についての詳細は、「[オプショングループを作成する](#)」を参照してください。
3. オプションを SSL のオプショングループに追加します。SSL に指定したポートは、後で Amazon RDS Oracle インスタンスの AWS Glue JDBC 接続 URL を作成するときに使用されます。Amazon RDS コンソールでオプションを追加する方法の詳細については、Amazon RDS ユーザーガイドの「[オプショングループにオプションを追加する](#)」を参照してください。の。Oracle SSL のオプションの詳細については、Amazon RDS ユーザーガイドの「[Oracle SSL](#)」を参照してください。
4. AWS Glue コンソールで、Amazon RDS Oracle インスタンスへの接続を作成します。接続定義で、[Require SSL connection] を選択します。Amazon RDS Oracle SSL オプションで使用しているポートを、求められれば入力します。

接続に [Require SSL connection] (SSL 接続が必要) が選択されている場合は、次の追加のオプションプロパティを使用できます。

S3 のカスタム JDBC 証明書

オンプレミスまたはクラウドデータベースとの SSL 通信に現在使用している証明書がある場合は、この証明書を AWS Glue データソースまたはターゲットへの SSL 接続に使用できます。カスタムルート証明書が含まれている Amazon Simple Storage Service (Amazon S3) の場所を入力します。AWS Glue はこの証明書を使用してデータベースへの SSL 接続を確立します。AWS

Glue は X.509 証明書のみを処理します。証明書は DER エンコードし、Base64 エンコード PEM 形式で指定する必要があります。

このフィールドを空白のままにすると、デフォルトの証明書が使用されます。

カスタム JDBC 証明書文字列

JDBC データベース固有の証明書情報を入力します。これは、ドメインの一致または識別名 (DN) の一致に使用される文字列です。Oracle Database の場合、この文字列は `tnsnames.ora` ファイルの `security` セクションの `SSL_SERVER_CERT_DN` パラメータにマッピングされます。Microsoft SQL Server の場合、この文字列は `hostNameInCertificate` として使用されま

す。Oracle Database の `SSL_SERVER_CERT_DN` パラメータの例を次に示します。

```
cn=sales,cn=OracleContext,dc=us,dc=example,dc=com
```

Kafka のプライベート CA 証明書の場所

Kafka データストアとの SSL 通信に現在使用している証明書がある場合は、この証明書を AWS Glue 接続に使用できません。このオプションは、Kafka データストアでは必須で、Amazon Managed Streaming for Apache Kafka データストアではオプションです。カスタムルート証明書が含まれている Amazon Simple Storage Service (Amazon S3) の場所を入力します。AWS Glue はこの証明書を使用して Kafka データストアへの SSL 接続を確立します。AWS Glue は X.509 証明書のみを処理します。証明書は DER エンコードし、Base64 エンコード PEM 形式で指定する必要があります。

証明書の検証をスキップ

「証明書の検証をスキップ」チェックボックスをオンにして、AWS Glue によるカスタム証明書の検証をスキップします。検証することを選択すると、AWS Glue は証明書の署名アルゴリズムとサブジェクトパブリックキーアルゴリズムを検証します。証明書の検証に失敗すると、この接続を使用する ETL ジョブまたはクローラーはすべて失敗します。

許可される署名アルゴリズムは、SHA256withRSA、SHA384withRSA、または SHA512withRSA のみです。サブジェクトパブリックキーアルゴリズムの場合、キーの長さは 2048 以上にする必要があります。

Kafka クライアントのキーストアの場合

Kafka クライアント側認証用のクライアントキーストアファイルの Amazon S3 の場所。パスは `s3://bucket/prefix/filename.jks` の形式である必要があります。ファイル名と `.jks` 拡張子で終わる必要があります。

Kafka クライアントキーストアのパスワード (オプション)

指定されたキーストアにアクセスするためのパスワード。

Kafka クライアントキーのパスワード (オプション)

キーストアは、複数のキーで構成することができるので、これは、Kafkaサーバー側のキーで使用されるクライアントキーにアクセスするためのパスワードです。

クライアント認証用の Apache Kafka 接続プロパティ

AWS Glue は、Apache Kafka 接続を作成するときの認証用の Simple Authentication and Security Layer (SASL) フレームワークをサポートしています。SASL フレームワークはさまざまな認証メカニズムをサポートしており、AWS Glue は SCRAM プロトコル (ユーザー名およびパスワード)、GSSAPI プロトコル (Kerberos プロトコル)、PLAIN プロトコルを提供します。

AWS Glue Studio を使用して、次のいずれかのクライアント認証方法を設定します。詳細については、ユーザーガイドの [「コネクタの接続の作成」](#) を参照してください。AWS Glue Studio

- None - 認証なし。これは、テスト目的で接続する場合に便利です。
- SASL/SCRAM-SHA-512 - この認証方法を選択すると、認証情報を指定することができます。2つのオプションがあります。
 - AWS Secrets Manager を使用する (推奨) - このオプションを選択すると、ユーザー名とパスワードを AWS Secrets Manager に保存し、必要に応じて AWS Glue からアクセスさせることができます。SSL または SASL 認証の認証情報を格納するシークレットを指定します。詳細については、[「AWS Secrets Manager への接続認証情報の保存」](#) を参照してください。
 - 有効なユーザー名とパスワードを直接指定します。
- SASL/GSSAPI (Kerberos) - このオプションを選択すると、キータブファイル、`krb5.conf` ファイルの場所を選択して、Kerberos プリンシパル名と Kerberos サービス名を入力することができます。キータブファイルと `krb5.conf` ファイルの場所は、Simple Storage Service (Amazon S3) がある場所の中にする必要があります。MSK は SASL/GSSAPI をまだサポートしていないため、このオプションは Customer Managed Apache Kafka クラスターでのみ使用できます。詳細については、[MIT Kerberos ドキュメント: キータブ](#) を参照してください。

- SASL/PLAIN - 認証用の認証情報を指定するためにこの認証方法を選択します。2つのオプションがあります。
- AWS Secrets Manager を使用する (推奨) - このオプションを選択すると、認証情報を AWS Secrets Manager に保存し、必要に応じて情報 AWS Glue へのアクセスを許可できます。SSL または SASL 認証の認証情報を格納するシークレットを指定します。
- ユーザー名とパスワードを直接指定します。
- SSL クライアント認証 - このオプションを選択すると、Simple Storage Service (Amazon S3) を参照することで Kafka クライアントキーストアの場所を選択できます。オプションで、Kafka クライアントキーストアのパスワードと Kafka クライアントキーのパスワードを入力できます。

Google BigQuery 接続

AWS Glue ETL ジョブで使用される Google BigQuery 接続をセットアップするには、次のプロパティを使用します。詳細については、「[the section called “BigQuery 接続”](#)」を参照してください。

AWS シークレット

AWS Secrets Manager. AWS Glue ETL ジョブのシークレットのシークレット名は、シークレットの `credentials` キー BigQuery を使用して Google に接続します。

Vertica 接続

AWS Glue ETL ジョブで使用される Vertica 接続をセットアップするには、次のプロパティを使用します。詳細については、「[the section called “Vertica 接続”](#)」を参照してください。

AWS Secrets Manager への接続認証情報の保存

AWS Secrets Manager を使用してデータストアの接続認証情報を入力することをお勧めします。このように Secrets Manager を使用すれば、AWS Glue は ETL ジョブのランタイムやクローラー実行の際にシークレットにアクセスでき、認証情報を安全に保つことができます。

前提条件

AWS Glue で Secrets Manager を使用するには、[AWS Glue の IAM ロール](#)にシークレット値を取得するアクセス許可を付与する必要があります。AWS 管理ポリシー `AWSGlueServiceRole` には AWS Secrets Manager へのアクセス許可は含まれません。IAM ポリシーの例については、「AWS

Secrets Manager ユーザーガイド」の「[例: シークレット値を取得するアクセス許可](#)」を参照してください。

ネットワークの設定によっては、VPC エンドポイントを作成して VPC と Secrets Manager 間のプライベート接続を確立する必要がある場合もあります。詳細については、「[AWS Secrets Manager VPC エンドポイントの使用](#)」を参照してください。

AWS Glue のシークレットを作成するには

1. 「AWS Secrets Manager ユーザーガイド」の「[Create and manage secrets](#)」(シークレットの作成と管理) の手順に従ってください。次の JSON の例は、AWS Glue のシークレットを作成する際に、[Plaintext] タブで認証情報を指定する方法を示しています。

```
{
  "username": "EXAMPLE-USERNAME",
  "password": "EXAMPLE-PASSWORD"
}
```

2. AWS Glue Studio インターフェイスを使用して、シークレットを接続に関連付けます。手順の詳細については、「AWS Glue Studio ユーザーガイド」の「[コネクタ用の接続を作成する](#)」を参照してください。

AWS Glue 接続の追加

AWS Glue for Spark のデータソースにはプログラマ的に接続できます。詳細については、「[AWS Glue for Spark での ETL の接続タイプとオプション](#)」を参照してください。

また、AWS Glue コンソールを使用して、接続の追加、編集、削除、およびテストを行うことができます。AWS Glue 接続の詳細については、「[データへの接続](#)」を参照してください。

AWS Glue 接続を追加するには

1. にサインイン AWS Management Console し、<https://console.aws.amazon.com/glue/> で AWS Glue コンソールを開きます。
2. ナビゲーションペインの [Data catalog] で [Connections] (接続) を選択します。
3. [接続の追加] を選択してウィザードを完了し、「[the section called "AWS Glue 接続プロパティ"](#)」の説明に従って接続プロパティを入力します。

での Amazon Redshift への接続 AWS Glue Studio

Note

AWS Glue for Spark を使用して、の外部にある Amazon Redshift データベースのテーブルに対する読み取りと書き込みを行うことができます AWS Glue Studio。AWS Glue ジョブ Amazon Redshift でプログラマ的にを設定するには、「」を参照してください [Redshift 接続](#)。

AWS Glue は、の組み込みサポートを提供します Amazon Redshift。AWS Glue Studioは、に接続し Amazon Redshift、データ統合ジョブを作成し、AWS Glue Studioサーバーレス Spark ランタイムで実行するためのビジュアルインターフェイスを提供します。

トピック

- [Amazon Redshift 接続の作成](#)
- [Amazon Redshift ソースノードの作成](#)
- [Amazon Redshift ターゲットノードの作成](#)
- [詳細オプション](#)

Amazon Redshift 接続の作成

必要となる許可

Amazon Redshift Amazon Redshift クラスターとサーバーレス環境を使用するには、追加の権限が必要です。ETL ジョブに許可を追加する方法の詳細については、「[ETL ジョブに必要な IAM アクセス許可を確認する](#)」を参照してください。

- 赤方偏移:DescribeClusters
- レッドシフト-サーバーレス:ListWorkgroups
- レッドシフト-サーバーレス:ListNamespaces

概要

Amazon Redshift 接続を追加するときは、Amazon Redshift 既存の接続を選択するか、データソース-Redshift ノードを追加するときに新しい接続を作成できます。AWS Glue Studio

AWS Glue Amazon Redshift Amazon Redshift クラスターとサーバーレス環境の両方をサポートします。接続を作成すると、Amazon Redshift サーバーレス環境では接続オプションの横にサーバーレスラベルが表示されます。

Amazon Redshift 接続の作成方法については、「[データ間の移動](#)」を参照してください。
Amazon Redshift

Amazon Redshift ソースノードの作成

必要となる許可

Amazon Redshift データソースを使用する AWS Glue Studio ジョブには追加の許可が必要です。ETL ジョブに許可を追加する方法の詳細については、「[ETL ジョブに必要な IAM アクセス許可を確認する](#)」を参照してください。

Amazon Redshift 接続を使用するには、次の許可が必要です。

- redshift-data:ListSchemas
- redshift-data:ListTables
- redshift-data:DescribeTable
- redshift-data:ExecuteStatement
- redshift-data:DescribeStatement
- redshift-data:GetStatementResult

Amazon Redshift データソースの追加

データソース — Amazon Redshift ノードを追加するには

1. Amazon Redshift アクセスタイプを選択します。
 - 直接データ接続 (推奨) — Amazon Redshift データに直接アクセスする場合は、このオプションを選択してください。これは推奨されるオプションで、デフォルトでもあります。
 - Data Catalog tables — 使用したいデータカタログテーブルがある場合は、このオプションを選択してください。
2. [直接データ接続] を選択した場合は、Amazon Redshift データソースへの接続を選択します。既に接続が存在していて、それらの接続から選択できることが前提です。接続の作成が必要な場合は、[Redshift 接続の作成] を選択します。詳細については、「[コネクタと接続の使用に関する概要](#)」を参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。URL、セキュリティグループ、サブネット、アベイラビリティーゾーン、説明、作成時刻 (UTC) と最終更新時刻 (UTC) のタイムスタンプなど、接続についての情報が表示されます。

3. Amazon Redshift のソースオプションを選択してください。

- [Choose a single table] - これは、単一の Amazon Redshift テーブルからアクセスするデータを含むテーブルです。
- カスタムクエリを入力 — カスタムクエリに基づき、Amazon Redshift の複数のテーブルからデータセットにアクセスできます。

4. 単一のテーブルを選択した場合は、Amazon Redshift のスキーマを選択します。選択したテーブルによって、使えるスキーマのリストが決まります。

または、[カスタムクエリを入力] を選択します。このオプションを選択すると、Amazon Redshift の複数のテーブルから、カスタムデータセットにアクセスできます。このオプションを選択した場合は、Amazon Redshift クエリを入力します。

Amazon Redshift サーバーレス環境に接続するときに、カスタムクエリに次の許可を追加します。

```
GRANT SELECT ON ALL TABLES IN <schema> TO PUBLIC
```

[スキーマを推測] を選択すると、入力したクエリに基づいてスキーマを読み取ることができません。[Redshift クエリエディタを開く] を選択して、Amazon Redshift クエリを入力することもできます。詳細については、「[クエリエディタを使用してデータベースのクエリを実行する](#)」を参照してください。

5. [パフォーマンスとセキュリティ] で、Amazon S3 のステージングディレクトリと IAM ロールを選択します。

- Amazon S3 ステージングディレクトリ — データを一時的にステージングする、Amazon S3 の場所を選択します。
- IAM ロール — 選択した Amazon S3 に書き込める IAM ロールを選択します。

6. [カスタム Redshift パラメータ - オプション] に、パラメータと値を入力します。

Amazon Redshift ターゲットノードの作成

必要となる許可

AWS Glue Studio Amazon Redshift データターゲットを使用するジョブには追加の権限が必要です。ETL ジョブに許可を追加する方法の詳細については、「[ETL ジョブに必要な IAM アクセス許可を確認する](#)」を参照してください。

Amazon Redshift 接続を使用するには以下の権限が必要です。

- 赤方偏移データ:ListSchemas
- 赤方偏移データ:ListTables

ターゲットノードの追加 Amazon Redshift

Amazon Redshift ターゲットノードを作成するには:

1. Amazon Redshift 既存のテーブルをターゲットとして選択するか、新しいテーブル名を入力します。
2. データターゲット - Redshift ターゲットノードを使用する場合、次のオプションから選択できます。
 - APPEND — テーブルがすでに存在している場合は、新しいデータをテーブルにすべて挿入してダンプします。テーブルが存在しない場合は、新規にテーブルを作成し、そこに新しいデータをすべて挿入します。

またターゲットのテーブルで、既存のレコードを更新 (UPSERT) する場合は、このチェックボックスにチェックを入れてください。まずテーブルが存在している必要があります。存在しない場合には、操作は失敗します。

- MERGE — 指定した条件に基づいて、AWS Glue がターゲットとなるテーブルのデータを更新するか、またはデータを追加します。

Note

でマージアクションを使用するにはAWS Glue、 Amazon Redshift マージ機能を有効にする必要があります。Amazon Redshift インスタンスのマージを有効にする方法については、「[MERGE \(プレビュー\)](#)」を参照してください。

[オプション] を選択します。

- キーと簡単なアクションの選択 — ソースデータとターゲットデータセットとの、マッチングキーとして使用する列を選択します。
 - 一致した場合、次のオプションを指定します。
 - ターゲットのデータセットにあるレコードを、ソースのデータで更新します。
 - ターゲットのデータセットにあるレコードを削除します。
 - 一致しない場合、次のオプションを指定します。
 - ターゲットのデータセットに、新しい行としてソースデータを挿入します。
 - 何もしない。
- カスタム MERGE ステートメントの入力 — その後 [MERGE ステートメントの検証] を選択し、ステートメントが有効か、無効かを検証できます。
- TRUNCATE — 既にテーブルが存在している場合は、ターゲットのテーブルの内容を削除してから、ターゲットのテーブルを削除します。削除が成功してから、すべてのデータを挿入します。テーブルが存在していない場合、テーブルを作成し、すべてのデータを挿入します。削除が成功しなかった場合、操作は失敗します。
- DROP — 既にテーブルが存在している場合は、テーブルのメタデータとデータを削除します。削除が成功してから、すべてのデータを挿入します。テーブルが存在していない場合、テーブルを作成し、すべてのデータを挿入します。削除が成功しなかった場合、操作は失敗します。
- CREATE — 新しいテーブルを、デフォルトの名前を使用して作成します。既にテーブル名が存在する場合は、接尾辞として `job_datetime` を名前につけたうえで新しいテーブルを作成し、一意性を保ちます。これで、すべてのデータが新しいテーブルに挿入されます。テーブルが存在する場合、最終的なテーブル名には接尾辞が付きます。テーブルが存在しない場合、テーブルが作成されます。いずれの場合も、テーブルが新しく作成されます。

詳細オプション

「[AWS Glue で Amazon Redshift Spark コネクタを使用する](#)」を参照してください。

AWS Glue Studio での Azure Cosmos DB に対する接続

AWS Glue は、Azure Cosmos DB のための組み込みサポートを提供します。AWS Glue Studio は、Azure Cosmos DB for NoSQL に接続してデータ統合ジョブをオーサリングし、AWS Glue

Studio サーバーレス Spark ランタイム上でそれらのジョブを実行するためのビジュアルインターフェイスを提供します。

トピック

- [Azure Cosmos DB 接続の作成](#)
- [Azure Cosmos DB ソースノードの作成](#)
- [Azure Cosmos DB ターゲットノードの作成](#)
- [詳細オプション](#)

Azure Cosmos DB 接続の作成

前提条件:

- Azure では、AWS Glue で使用する Azure Cosmos DB キー (cosmosKey) を特定または生成する必要があります。詳細については、Azure ドキュメントの「[Azure Cosmos DB のデータへのアクセスをセキュリティで保護する](#)」を参照してください。

Azure Cosmos DB に対する接続を設定するには:

1. AWS Secrets Manager で、Azure Cosmos DB キーを使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
 - [key/value ペア] を選択する際に、*cosmosKey* という値を持つキー `spark.cosmos.accountKey` のペアを作成します。
2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、将来的に AWS Glue で使用するために、接続名 *connectionName* を維持します。
 - [接続タイプ] を選択する際に、[Azure Cosmos DB] を選択します。
 - [AWS Secret] をクリックして、*secretName* を入力します。

Azure Cosmos DB ソースノードの作成

必要な前提条件

- 前のセクション [the section called “Azure Cosmos DB 接続の作成”](#) で説明したように、AWS Secrets Manager シークレットを使用して設定された AWS Glue Azure Cosmos DB 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 読み取り元とする Azure Cosmos DB for NoSQL コンテナ。コンテナの識別情報が必要になります。

Azure Cosmos for NoSQL コンテナは、データベースとコンテナによって識別されます。Azure Cosmos for NoSQL API に接続する際には、データベース *cosmosDBName* とコンテナ *cosmosContainerName* の名前を指定する必要があります。

Azure Cosmos DB データソースの追加

[データソース — Azure Cosmos DB] ノードを追加するには:

1. Azure Cosmos DB データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずです。接続を作成する必要がある場合は、[Azure Cosmos DB 接続を作成] を選択します。詳細については、前の「[the section called “Azure Cosmos DB 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. [Cosmos DB データベース名] を選択します。読み取り元のデータベースの名前である *cosmosDBName* を指定します。
3. [Azure Cosmos DB コンテナ] を選択します。読み取り元のコンテナの名前である *cosmosContainerName* を指定します。
4. 必要に応じて、[Azure Cosmos DB カスタムクエリ] を選択します。Azure Cosmos DB から特定の情報を取得するには、SQL SELECT クエリを指定します。
5. [Azure Cosmos のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

Azure Cosmos DB ターゲットノードの作成

必要な前提条件

- 前のセクション [the section called “Azure Cosmos DB 接続の作成”](#) で説明したように、AWS Secrets Manager シークレットを使用して設定された AWS Glue Azure Cosmos DB 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 書き込み先とする Azure Cosmos DB テーブル。コンテナの識別情報が必要になります。接続メソッドを呼び出す前にコンテナを作成する必要があります。

Azure Cosmos for NoSQL コンテナは、データベースとコンテナによって識別されます。Azure Cosmos for NoSQL API に接続する際には、データベース *cosmosDBName* とコンテナ *cosmosContainerName* の名前を指定する必要があります。

Azure Cosmos DB データターゲットの追加

[データターゲット — Azure Cosmos DB] ノードを追加するには:

1. Azure Cosmos DB データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずですが、接続を作成する必要がある場合は、[Azure Cosmos DB 接続を作成] を選択します。詳細については、前の「[the section called “Azure Cosmos DB 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. [Cosmos DB データベース名] を選択します。読み取り元のデータベースの名前である *cosmosDBName* を指定します。
3. [Azure Cosmos DB コンテナ] を選択します。読み取り元のコンテナの名前である *cosmosContainerName* を指定します。
4. [Azure Cosmos のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

詳細オプション

Azure Cosmos DB ノードを作成する際に、高度なオプションを指定できます。これらのオプションは Spark AWS Glue スクリプトのプログラミング時に使用できるオプションと同じです。

「[the section called “Azure Cosmos DB 接続”](#)」を参照してください。

AWS Glue Studio での Azure SQL に対する接続

AWS Glue は、Azure SQL のための組み込みサポートを提供します。AWS Glue Studio は、Azure SQL に接続してデータ統合ジョブをオーサリングし、AWS Glue Studio サーバーレス Spark ランタイム上でそれらのジョブを実行するためのビジュアルインターフェイスを提供します。

トピック

- [Azure SQL 接続の作成](#)
- [Azure SQL ソースノードの作成](#)
- [Azure SQL ターゲットノードの作成](#)
- [詳細オプション](#)

Azure SQL 接続の作成

AWS Glue から Azure SQL に接続するには、Azure SQL 認証情報を作成して AWS Secrets Manager シークレットに保存し、そのシークレットを Azure SQL AWS Glue 接続に関連付ける必要があります。

Azure SQL に対する接続を設定するには:

1. AWS Secrets Manager で、Azure SQL 認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
 - [key/value ペア] を選択する際に、*azuresqlUsername* という値を持つキー user のペアを作成します。
 - [key/value ペア] を選択する際に、*azuresqlPassword* という値を持つキー password のペアを作成します。
2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、将来的に AWS Glue で使用するために、接続名 *connectionName* を維持します。
 - [接続タイプ] を選択する際に、[Azure SQL] を選択します。
 - Azure SQL URL を指定する場合は、JDBC エンドポイント URL を入力します。

URL は、次のような形式になります:

```
jdbc:sqlserver://databaseServerName:databasePort;databaseName=azuresqlDBName
```

AWS Glue には次の URL プロパティが必要です。

- `databaseName` – 接続先の Azure SQL のデフォルトデータベース。

Azure SQL Managed Instances の JDBC URL の詳細については、[Microsoft のドキュメント](#)を参照してください。

- [AWS Secret] をクリックして、`secretName` を入力します。

Azure SQL ソースノードの作成

必要な前提条件

- 前のセクション [the section called “Azure SQL 接続の作成”](#) で説明したように、AWS Secrets Manager シークレットを使用して設定された AWS Glue Azure SQL 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 読み取り元とする Azure SQL テーブル (`tableName`)。

Azure SQL テーブルは、データベース、スキーマ、テーブル名によって識別されます。Azure SQL に接続する際には、データベース名とテーブル名を指定する必要があります。スキーマがデフォルトの「public」でない場合は、スキーマも指定する必要があります。データベースは、`connectionName` の URL プロパティ、`dbtable` を通じたスキーマおよびテーブル名を介して指定されます。

Azure SQL データソースの追加

[データソース — Azure SQL] ノードを追加するには:

1. Azure SQL データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずです。接続を作成する必要がある場合は、[Azure SQL 接続を作成] を選択します。詳細については、前の「[the section called “Azure SQL 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. [Azure SQL ソース] オプションを選択します。

- [単一のテーブルを選択] – 単一のテーブルからすべてのデータにアクセスします。
 - [カスタムクエリを入力] – カスタムクエリに基づいて、複数のテーブルからデータセットにアクセスします。
3. 単一のテーブルを選択した場合は、*tableName* を入力します。

[カスタムクエリを入力] を選択した場合は、TransactSQL SELECT クエリを入力します。
 4. [Azure SQL のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

Azure SQL ターゲットノードの作成

必要な前提条件

- 前のセクション [the section called “Azure SQL 接続の作成”](#) で説明したように、AWS Secrets Manager シークレットを使用して設定された AWS Glue Azure SQL 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 書き込み先とする Azure SQL テーブル (*tableName*)。

Azure SQL テーブルは、データベース、スキーマ、テーブル名によって識別されます。Azure SQL に接続する際には、データベース名とテーブル名を指定する必要があります。スキーマがデフォルトの「public」でない場合は、スキーマも指定する必要があります。データベースは、*connectionName* の URL プロパティ、dbtable を通じたスキーマおよびテーブル名を介して指定されます。

Azure SQL データターゲットの追加

[データターゲット — Azure SQL] ノードを追加するには:

1. Azure SQL データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずです。接続を作成する必要がある場合は、[Azure SQL 接続を作成] を選択します。詳細については、前の「[the section called “Azure SQL 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. *tableName* を指定して [テーブル名] を設定します。
3. [Azure SQL のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

詳細オプション

Azure SQL ノードを作成する際に、高度なオプションを指定できます。これらのオプションは Spark AWS Glue スクリプトのプログラミング時に使用できるオプションと同じです。

「[the section called “Azure SQL 接続”](#)」を参照してください。

BigQuery での Google への接続 AWS Glue Studio

Note

AWS Glue for Spark を使用して、4.0 以降のバージョンの Google AWS Glue BigQuery のテーブルとの間で読み書きを行うことができます。AWS Glue ジョブ BigQuery で Google をプログラムで設定するには、「[こちら](#)」を参照してください。

[BigQuery 接続](#)。

AWS Glue Studio は、[こちら](#)に接続し BigQuery、データ統合ジョブを作成し、AWS Glue Studio サーバーレス Spark ランタイムで実行するためのビジュアルインターフェイスを提供します。

トピック

- [BigQuery 接続を作成する](#)
- [BigQuery ソースノードの作成](#)
- [BigQuery ターゲットノードを作成する](#)
- [詳細オプション](#)

BigQuery 接続を作成する

AWS Glue から Google BigQuery へ接続するには、Google Cloud Platform の認証情報を作成して AWS Secrets Manager シークレットに保存し、そのシークレットを Google BigQuery AWS Glue 接続に関連付ける必要があります。

BigQuery への接続を設定するには:

1. Google Cloud Platform で、関連するリソースを作成して特定します。
 - 接続したい BigQuery テーブルを含む GCP プロジェクトを作成または特定します。

- BigQuery API を有効にします。詳細については、「[Use the BigQuery Storage Read API to read table data](#)」を参照してください。
2. Google Cloud Platform で、サービスアカウントの認証情報を作成してエクスポートします。

BigQuery 認証情報ウィザードを使用すると、「[認証情報の作成](#)」というステップを迅速に実行できます。

GCP でサービスアカウントを作成するには、「[サービス アカウントを作成する](#)」にあるチュートリアルに従ってください。

- プロジェクトを選択するときは、BigQuery テーブルを含むプロジェクトを選択します。
- サービスアカウントの GCP IAM ロールを選択するときは、BigQuery テーブルの読み取り、書き込み、作成を行う BigQuery ジョブを実行するための適切な権限を付与するロールを追加または作成します。

サービスアカウントの認証情報を作成するには、「[サービス アカウント キーを作成する](#)」にあるチュートリアルに従ってください。

- キータイプを選択するときは、[JSON] を選択します。

これで、サービスアカウントの認証情報が記載された JSON ファイルがダウンロードされたはずですが、これは次のように表示されます。

```
{
  "type": "service_account",
  "project_id": "*****",
  "private_key_id": "*****",
  "private_key": "*****",
  "client_email": "*****",
  "client_id": "*****",
  "auth_uri": "https://accounts.google.com/o/oauth2/auth",
  "token_uri": "https://oauth2.googleapis.com/token",
  "auth_provider_x509_cert_url": "https://www.googleapis.com/oauth2/v1/certs",
  "client_x509_cert_url": "*****",
  "universe_domain": "googleapis.com"
}
```

3. ダウンロードした認証情報ファイルを base64 でエンコードします。AWS CloudShell セッションなどでは、コマンドラインから `cat credentialsFile.json | base64 -w 0` コマンド

を実行してこれを実行できます。このコマンドの出力 *credentialString* を保持してください。

4. AWS Secrets Manager で、Google Cloud Platform の認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
 - キーと値のペアを選択するときは、*credentialString* という値を持つキー `credentials` のペアを作成します。
5. AWS Glue データカタログで、<https://docs.aws.amazon.com/glue/latest/dg/console-connections.html> にある手順に従って接続を作成します。接続を作成したら、次のステップのために接続名 *connectionName* を保存しておきます。
 - [接続タイプ] を選択するときは、Google BigQuery を選択してください。
 - [AWS Secret] をクリックして、*secretName* を入力します。
6. AWS Glue ジョブに関連付けられている IAM ロールに *secretName* を読み取るアクセス許可を付与します。
7. AWS Glue ジョブ設定で、[Additional network connection] として [*connectionName*] を指定します。

BigQuery ソースノードの作成

必要な前提条件

- BigQuery タイプの AWS Glue データカタログ接続
- 接続時に使用される Google BigQuery 認証情報の AWS Secrets Manager シークレット。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 読み取る対象のテーブルおよび対応する Google Cloud プロジェクトの名前とデータセット。

BigQuery データソースを追加する

データソース — BigQuery ノードを追加するには

1. BigQuery データソースの接続を選択します。すでに作成したので、ドロップダウンに表示されているはずですが。接続を作成する必要がある場合は、[Create BigQuery connection] を選択します。詳細については、「[コネクタと接続の使用に関する概要](#)」を参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. 読み取る BigQuery データを特定し、[BigQuery Source] オプションを選択します。
 - Choose a single table - テーブルからすべてのデータを取得できます。
 - Enter a custom query — クエリを指定することで、取得するデータをカスタマイズできます。
3. 読み取るデータの説明を記述します。

(必須) テーブルを含むプロジェクト、または該当する場合は、請求元の親プロジェクトに[親プロジェクト]を設定します。

1つのテーブルを選択した場合は、[テーブル] を次の形式で Google BigQuery テーブルの名前に設定します。[dataset].[table]

クエリを選択した場合は、それを [クエリ] に提供します。クエリでは、[project].[dataset].[tableName] の形式の完全修飾テーブル名の付いたテーブルを参照します。

4. BigQuery プロパティを指定します。

1つのテーブルを選択した場合は、追加のプロパティを指定する必要はありません。

クエリを選択した場合は、以下の [Custom Google BigQuery properties] を指定する必要があります。

- viewsEnabled を true に設定します。
- materializationDataset をデータセットに設定します。AWS Glue 接続を通じて提供された認証情報で認証された GCP プリンシパルは、このデータセットにテーブルを作成できる必要があります。

BigQuery ターゲットノードを作成する

必要な前提条件

- BigQuery タイプの AWS Glue データカタログ接続
- 接続時に使用される Google BigQuery 認証情報の AWS Secrets Manager シークレット。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 書き込み対象のテーブルおよび対応する Google Cloud プロジェクトの名前とデータセット。

BigQuery データターゲットを追加する

[Data target – BigQuery] ノードを追加するには:

1. BigQuery データターゲットの接続を選択します。すでに作成したので、ドロップダウンに表示されているはずです。接続を作成する必要がある場合は、[Create BigQuery connection] を選択します。詳細については、[「コネクタと接続の使用に関する概要」](#)を参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. どの BigQuery テーブルに書き込むかを特定し、[Write method] を選択します。
 - 直接 — BigQuery Storage Write API を使用して BigQuery に直接書き込みます。
 - 間接 — Google Cloud Storage に書き込んでから BigQuery にコピーします。

間接的に書き込む場合は、一時 GCS バケットを使用して送信先の GCS ロケーションを指定します。AWS Glue 接続時に追加の設定を行う必要があります。詳細については、[「Using indirect write with Google BigQuery」](#)を参照してください。

3. 読み取るデータの説明を記述します。

(必須) テーブルを含むプロジェクト、または該当する場合は、請求元の親プロジェクトに[親プロジェクト]を設定します。

1 つのテーブルを選択した場合は、[テーブル] を次の形式で Google BigQuery テーブルの名前に設定します。[dataset].[table]

詳細オプション

BigQuery ノードの作成時に高度なオプションを指定できます。これらのオプションは、Spark スクリプト AWS Glue のプログラミング時に使用できるオプションと同じです。

デ AWS Glue ペロツパーガイド [BigQuery の接続オプションリファレンス](#)を参照してください。

AWS Glue Studio での MongoDB に対する接続

AWS Glue は、MongoDB のための組み込みサポートを提供します。AWS Glue Studio は、MongoDB に接続してデータ統合ジョブをオーサリングし、AWS Glue Studio サーバーレス Spark ランタイム上でそれらのジョブを実行するためのビジュアルインターフェイスを提供します。

トピック

- [MongoDB 接続の作成](#)
- [MongoDB ソースノードの作成](#)
- [MongoDB ターゲットノードの作成](#)
- [詳細オプション](#)

MongoDB 接続の作成

前提条件:

- MongoDB インスタンスが Amazon VPC 内にある場合は、トラフィックがパブリックインターネットを経由することなく、AWS Glue ジョブが MongoDB インスタンスと通信できるように Amazon VPC を設定します。

Amazon VPC で、AWS Glue がジョブの実行中に使用する [VPC]、[サブネット]、および [セキュリティグループ] を特定または作成します。さらに、MongoDB インスタンスとこの場所の間のネットワークトラフィックを許可するように Amazon VPC が設定されているようにする必要があります。ネットワークレイアウトに基づいて、セキュリティグループルール、ネットワーク ACL、NAT ゲートウェイ、およびピアリング接続の変更が必要になる場合があります。

MongoDB に対する接続を設定するには:

1. 必要に応じて、AWS Secrets Manager で、MongoDB 認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください

い。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。

- [key/value ペア] を選択する際に、*mongodbUser* という値を持つキー `username` のペアを作成します。

[key/value ペア] を選択する際に、*mongodbPass* という値を持つキー `password` のペアを作成します。

2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、将来的に AWS Glue で使用するために、接続名 *connectionName* を維持します。

- [接続タイプ] を選択する際には、[MongoDB] または [MongoDB Atlas] を選択します。
- [MongoDB URL] または [MongoDB Atlas URL] を選択する場合は、MongoDB インスタンスのホスト名を入力します。

MongoDB URL は、`mongodb://mongoHost:mongoPort/mongoDBname` の形式で指定されます。

MongoDB Atlas URL は、`mongodb+srv://mongoHost:mongoPort/mongoDBname` の形式で指定されます。

接続にデフォルトのデータベースを指定する場合、*mongoDBname* はオプションです。

- Secrets Manager シークレットを作成することを選択した場合は、AWS Secrets Manager の [認証情報タイプ] を選択します。

その後、[AWS シークレット] で *secretName* を入力します。

- [ユーザー名とパスワード] を入力することを選択した場合は、*mongodbUser* および *mongodbPass* を入力します。

3. 次の状況では、追加の設定が必要になる場合があります。

- Amazon VPC の AWS でホストされている MongoDB インスタンスの場合
 - MongoDB セキュリティ認証情報を定義する AWS Glue 接続に、Amazon VPC 接続に関する情報を提供する必要があります。接続を作成または更新する際に、[ネットワークオプション] で [VPC]、[サブネット]、および [セキュリティグループ] を設定します。

AWS Glue MongoDB 接続を作成した後、AWS Glue ジョブを実行する前に次のステップを実行する必要があります。

- ビジュアルエディタで AWS Glue ジョブを使用する場合、ジョブが MongoDB に接続できるように Amazon VPC の接続に関する情報を入力する必要があります。Amazon VPC 内の適切な場所を特定し、それを AWS Glue MongoDB 接続に提供します。
- Secrets Manager シークレットを作成することを選択した場合は、AWS Glue ジョブに関連付けられた IAM ロールに *secretName* を読み取るための許可を付与します。

MongoDB ソースノードの作成

必要な前提条件

- 前のセクション「[the section called “MongoDB 接続の作成”](#)」で説明した AWS Glue MongoDB 接続。
- Secrets Manager シークレットを作成することを選択した場合は、接続によって使用されるシークレットを読み取るためのジョブに対する適切な許可が必要です。
- 読み取り元とする MongoDB コレクション。コレクションの識別情報が必要になります。

MongoDB コレクションは、データベース名とコレクション名である *mongodbName* および *mongodbCollection* によって識別されます。

MongoDB データソースの追加

[データソース — MongoDB] ノードを追加するには:

1. MongoDB データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずです。接続を作成する必要がある場合は、[MongoDB 接続を作成] を選択します。詳細については、前の「[the section called “MongoDB 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. [データベース] を選択します。 *mongodbName* を入力します。
3. [コレクション] を選択します。 *mongodbCollection* を入力します。
4. [パーティショナー]、[パーティションサイズ (MB)]、および [パーティションキー] を選択します。パーティションパラメータについての詳細は、「[the section called “connectionType": "mongodb" ソースとする”](#)」を参照してください。
5. [MongoDB のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

MongoDB ターゲットノードの作成

必要な前提条件

- 前のセクション [the section called “MongoDB 接続の作成”](#) で説明したように、AWS Secrets Manager シークレットを使用して設定された AWS Glue MongoDB 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 書き込み先とする MongoDB テーブルである *tableName*。

MongoDB データターゲットの追加

[データターゲット — MongoDB] ノードを追加するには:

1. MongoDB データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずですが、接続を作成する必要がある場合は、[MongoDB 接続を作成] を選択します。詳細については、前の「[the section called “MongoDB 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. [データベース] を選択します。 *mongodbName* を入力します。
3. [コレクション] を選択します。 *mongodbCollection* を入力します。
4. [パーティショナー]、[パーティションサイズ (MB)]、および [パーティションキー] を選択します。パーティションパラメータについての詳細は、「[the section called “connectionType”: "mongodb" ソースとする”](#)」を参照してください。
5. 必要に応じて、[書き込みを再試行] を選択します。
6. [MongoDB のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

詳細オプション

MongoDB ノードを作成する際に、高度なオプションを指定できます。これらのオプションは Spark AWS Glue スクリプトのプログラミング時に使用できるオプションと同じです。

「[the section called “MongoDB 接続”](#)」を参照してください。

AWS Glue Studio での OpenSearch Service に対する接続

AWS Glue は、Amazon OpenSearch Service のための組み込みサポートを提供します。AWS Glue Studio は、Amazon OpenSearch Service に接続してデータ統合ジョブをオーサリングし、AWS Glue Studio サーバーレス Spark ランタイム上でそれらのジョブを実行するためのビジュアルインターフェイスを提供します。この機能は OpenSearch サービスのサーバーレスとは互換性がありません。

トピック

- [OpenSearch Service 接続の作成](#)
- [OpenSearch Service のソースノードの作成](#)
- [OpenSearch Service ターゲットノードの作成](#)
- [詳細オプション](#)

OpenSearch Service 接続の作成

前提条件:

- 読み取り元とするドメインエンドポイント *aosEndpoint* とポート *aosPort* を特定するか、または Amazon OpenSearch Service ドキュメントの手順に従ってリソースを作成します。ドメインの作成の詳細については、Amazon OpenSearch Service ドキュメントの「[Amazon OpenSearch Service ドメインの作成と管理](#)」を参照してください。

Amazon OpenSearch Service ドメインエンドポイントのデフォルト形式は、`https://search-domainName-unstructuredIdContent.region.es.amazonaws.com` です。ドメインエンドポイントの識別の詳細については、Amazon OpenSearch Service ドキュメントの「[Amazon OpenSearch Service ドメインの作成と管理](#)」を参照してください。

ドメインの HTTP 基本認証情報、*aosUser*、および *aosPassword* を特定または生成します。

OpenSearch Service に対する接続を設定するには:

1. AWS Secrets Manager で、OpenSearch Service 認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。

- [key/value ペア] を選択する際に、*aosUser* という値を持つキー `opensearch.net.http.auth.user` のペアを作成します。
 - [key/value ペア] を選択する際に、*aosPassword* という値を持つキー `opensearch.net.http.auth.pass` のペアを作成します。
2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、将来的に AWS Glue で使用するために、接続名 *connectionName* を維持します。
- [接続タイプ] を選択する場合は、[OpenSearch Service] を選択します。
 - ドメインエンドポイントを選択する場合は、*aosEndpoint* を入力します。
 - ポートを選択する場合は、*aosPort* を入力します。
 - [AWS Secret] をクリックして、*secretName* を入力します。

OpenSearch Service のソースノードの作成

必要な前提条件

- 前のセクション [the section called “OpenSearch Service 接続の作成”](#) で説明したような、AWS Secrets Manager シークレットを使用して設定された AWS Glue OpenSearch Service 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 読み取り元とする OpenSearch Service インデックスである *aosIndex*。

OpenSearch Service データソースの追加

[データソース – OpenSearch Service] ノードを追加するには:

1. OpenSearch Service データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずです。接続を作成する必要がある場合は、[OpenSearch Service 接続を作成] を選択します。詳細については、前の「[the section called “OpenSearch Service 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. [インデックス] には、読み取り元とするインデックスを指定します。

3. オプションで、より具体的な結果を提供する OpenSearch クエリである [クエリ] を入力します。OpenSearch クエリの作成の詳細については、「[the section called “OpenSearch Service から読み取る”](#)」を参照してください。
4. [OpenSearch Service のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

OpenSearch Service ターゲットノードの作成

必要な前提条件

- 前のセクション [the section called “OpenSearch Service 接続の作成”](#) で説明したような、AWS Secrets Manager シークレットを使用して設定された AWS Glue OpenSearch Service 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 書き込み先とする OpenSearch Service インデックスである *aosIndex*。

OpenSearch Service データターゲットの追加

[データターゲット – OpenSearch Service] ノードを追加するには:

1. OpenSearch Service データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずです。接続を作成する必要がある場合は、[OpenSearch Service 接続を作成] を選択します。詳細については、前の「[the section called “OpenSearch Service 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. [インデックス] には、読み取り元とするインデックスを指定します。
3. [OpenSearch Service のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

詳細オプション

OpenSearch Service ノードを作成する際に、高度なオプションを指定できます。これらのオプションは Spark AWS Glue スクリプトのプログラミング時に使用できるオプションと同じです。

「[the section called “OpenSearch Service 接続”](#)」を参照してください。

での Salesforce への接続 AWS Glue Studio

Salesforce は、販売、カスタマーサービス、e コマースなどを支援する顧客関係管理 (CRM) ソフトウェアを提供しています。Salesforce ユーザーの場合は、Salesforce アカウント AWS Glue に接続できます。次に、ETL ジョブのデータソースまたは送信先として Salesforce を使用できます。これらのジョブを実行して、Salesforce と AWS のサービス、またはその他のサポートされているアプリケーション間でデータを転送します。

トピック

- [AWS Glue Salesforce のサポート](#)
- [接続を作成および使用するための API オペレーションを含むポリシー](#)
- [Salesforce の設定](#)
- [Salesforce 接続の設定](#)
- [Salesforce エンティティからの読み取り](#)
- [Salesforce への書き込み](#)
- [Salesforce 接続オプション](#)
- [Salesforce コネクタの制限事項](#)
- [Salesforce の JWT ベアラー OAuth フローを設定する](#)

AWS Glue Salesforce のサポート

AWS Glue は Salesforce を次のようにサポートします。

ソースとしてサポートされていますか？

はい。AWS Glue ETL ジョブを使用して Salesforce からデータをクエリできます。

ターゲットとしてサポートされていますか？

はい。AWS Glue ETL を使用して Salesforce にレコードを書き込むことができます。

サポートされている Salesforce API バージョン

次の Salesforce API バージョンがサポートされています。

- v58.0
- v59.0
- v60.0

接続を作成および使用するための API オペレーションを含むポリシー

次のサンプルポリシーでは、接続の作成と使用に必要な AWS IAM アクセス許可について説明します。新しいロールを作成する場合は、以下を含むポリシーを作成します。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "glue:ListConnectionTypes",
        "glue:DescribeConnectionType",
        "glue:RefreshOAuth2Tokens",
        "glue:ListEntities",
        "glue:DescribeEntity"
      ],
      "Resource": "*"
    }
  ]
}
```

次の IAM ポリシーを使用してアクセスを許可することもできます。

- [AWSGlueServiceRole](#) – さまざまな AWS Glue プロセスがユーザーに代わって実行するために必要なリソースへのアクセスを許可します。これらのリソースには AWS Glue、Amazon S3、IAM、CloudWatch Logs、Amazon EC2 が含まれます。このポリシーで指定されたリソースの命名規則に従うと、AWS Glue プロセスに必要なアクセス許可が付与されます。このポリシーは、通常、クローラ、ジョブ、開発エンドポイントを定義するときに指定されたロールにアタッチされます。
- [AWSGlueConsoleFullAccess](#) – ポリシーがアタッチされている ID が AWS マネジメントコンソールを使用する場合に、AWS Glue リソースへのフルアクセスを許可します。このポリシーで指定されたリソースの命名規則に従った場合、ユーザーは完全なコンソール機能を使用できます。このポリシーは通常、AWS Glue コンソールのユーザーにアタッチされます。

Salesforce の設定

を使用して Salesforce AWS Glue との間でデータを転送する前に、次の要件を満たす必要があります。

最小要件

以下に、最小要件を示します。

- Salesforce アカウントがある。
- Salesforce アカウントで API アクセスが有効になっています。API アクセスは、Enterprise、Unlimited、Developer、Performance の各エディションでデフォルトで有効になっています。
- Salesforce アカウントでは、接続されたアプリケーションをインストールできます。この機能にアクセスできない場合は、Salesforce 管理者にお問い合わせください。詳細については、Salesforce ヘルプの「[Connected Apps](#)」を参照してください。

これらの要件を満たすと、Salesforce アカウント AWS Glue に接続する準備が整います。は AWS 、マネージド接続アプリケーションを使用して残りの要件 AWS Glue を処理します。

Salesforce 用の AWS マネージド接続アプリケーション

AWS マネージド接続アプリケーションは、より少ないステップで Salesforce 接続を作成するのに役立ちます。Salesforce では、接続されたアプリケーションは、などの外部アプリケーションが Salesforce データにアクセス AWS Glue することを許可するフレームワークです。

- AWS Glue コンソールを使用して Salesforce 接続を作成します。
- 接続を設定するときは、OAuth 許可タイプを認証コード に設定します。

Salesforce 接続の設定

Salesforce 接続を設定するには：

1. AWS Secrets Manager で、次の詳細を含むシークレットを作成します。
 - a. JWT_TOKEN 許可タイプの場合 - シークレットには、JWT_TOKEN キーとその値が含まれている必要があります。
 - b. AuthorizationCode グラントタイプの場合: カスタマーマネージド接続アプリケーションの場合、シークレットには接続されたアプリケーションのコンシューマーシークレットをキー USER_MANAGED_CLIENT_APPLICATION_CLIENT_SECRET として含める必要があります。AWS マネージド接続アプリケーションの場合、空のシークレット、または一時的な値を持つシークレット。
 - c. 注: で接続ごとにシークレットを作成する必要があります AWS Glue。

2. AWS Glue データカタログで、以下の手順に従って接続を作成します。

- a. 接続タイプを選択するときは、Salesforce を選択します。
- b. 接続する Salesforce の INSTANCE_URL を指定します。
- c. Salesforce 環境を指定します。
- d. が引き受ける AWS Glue ことができ、以下のアクションに対するアクセス許可を持つ AWS IAM ロールを選択します。
- e. 接続に使用する OAuth2 許可タイプを選択します。グラントタイプは、が Salesforce と AWS Glue 通信してデータへのアクセスをリクエストする方法を決定します。選択すると、接続を作成する前に満たす必要がある要件に影響します。次のいずれかのタイプを選択できます。
 - JWT_BEARER グラントタイプ: このグラントタイプは、Salesforce インスタンス内の特定のユーザーのアクセス許可で JSON ウェブトークン (JWT) を事前に作成できるため、自動化シナリオに適しています。作成者は、JWT の有効期間を制御できます。AWS Glue は、JWT を使用して、Salesforce APIs の呼び出しに使用されるアクセストークンを取得できます。

このフローでは、ユーザーが Salesforce インスタンスで接続されたアプリケーションを作成し、ユーザーに JWT ベースのアクセストークンを発行できるようにする必要があります。

JWT ベアラー OAuth フロー用の接続アプリケーションの作成については、[server-to-server 「統合用の OAuth 2.0 JWT ベアラーフロー」](#) を参照してください。Salesforce 接続アプリケーションを使用して JWT ベアラーフローを設定するには、「」を参照してください [Salesforce の JWT ベアラー OAuth フローを設定する](#)。

- AUTHORIZATION_CODE 付与タイプ: この付与タイプは、ユーザーを認証するためにユーザーをサードパーティー認証サーバーにリダイレクトすることに依存するため、「3 レッグ OAuth」と見なされます。AWS Glue コンソール経由で接続を作成するときに使用されます。接続を作成するユーザーは、デフォルトで AWS Glue 接続されたアプリケーション (AWS Glue マネージドクライアントアプリケーション) に依存する場合があります。このアプリケーションでは、Salesforce インスタンス URL 以外の OAuth 関連情報を提供する必要はありません。AWS Glue コンソールはユーザーを Salesforce にリダイレクトし、ユーザーは Salesforce インスタンスにアクセスするためにログインし AWS Glue、リクエストされたアクセス許可を付与する必要があります。

コンソールから接続を作成する場合でも、ユーザーは Salesforce で独自の接続アプリケーションを作成し、独自のクライアント ID とクライアントシークレットを指定することができます AWS Glue。このシナリオでは、引き続き Salesforce にリダイレクトされ、ログインしてリソース AWS Glue へのアクセスを承認します。

このグラントタイプは、更新トークンとアクセストークンになります。アクセストークンは有効期間が短く、更新トークンを使用したユーザー操作なしで自動的に更新される場合があります。

認証コード OAuth フロー用の接続アプリケーションの作成については、[「認証コードと認証情報フロー用の接続アプリケーションの設定」](#)を参照してください。

- f. この接続に使用する secretName を選択して AWS Glue、トークンを配置します。
 - g. ネットワークを使用する場合は、ネットワークオプションを選択します。AWS Glue ジョブに関連付けられた IAM ロールに、 を読み取るアクセス許可を付与します secretName。
3. AWS Glue ジョブに関連付けられた IAM ロールに、 を読み取るアクセス許可を付与します secretName。
 4. AWS Glue ジョブ設定で、追加のネットワーク接続 connectionName として を指定します。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "secretsmanager:DescribeSecret",
        "secretsmanager:GetSecretValue",
        "secretsmanager:PutSecretValue",
        "ec2:CreateNetworkInterface",
        "ec2:DescribeNetworkInterface",
        "ec2>DeleteNetworkInterface",
      ],
      "Resource": "*"
    }
  ]
}
```

Salesforce エンティティからの読み取り

前提条件

読み取る Salesforce sObject。または AccountCase や などのオブジェクト名が必要になります Opportunity。

例:

```
salesforce_read = glueContext.create_dynamic_frame.from_options(  
    connection_type="salesforce",  
    connection_options={  
        "connectionName": "connectionName",  
        "ENTITY_NAME": "Account",  
        "API_VERSION": "v60.0"  
    }  
)
```

パーティションクエリ

Spark で同時実行を利用する NUM_PARTITIONS 場合は PARTITION_FIELD、追加の Spark LOWER_BOUND オプション UPPER_BOUND、、、を指定できます。これらのパラメータを使用すると、元のクエリは Spark タスクで同時に実行できるサブクエリ NUM_PARTITIONS の数に分割されます。

- PARTITION_FIELD: クエリのパーティション化に使用するフィールドの名前。
- LOWER_BOUND: 選択したパーティションフィールドの包括的な下限値。

timestamp フィールドでは、Spark SQL クエリで使用される Spark タイムスタンプ形式を使用できます。

有効な値の例：

```
"TIMESTAMP \"1707256978123\""  
"TIMESTAMP '2024-02-06 22:02:58.123 UTC'"  
"TIMESTAMP \"2018-08-08 00:00:00 Pacific/Tahiti\""  
"TIMESTAMP \"2018-08-08 00:00:00\""  
"TIMESTAMP \"-123456789\" Pacific/Tahiti"  
"TIMESTAMP \"1702600882\""
```

- UPPER_BOUND: 選択したパーティションフィールドの排他的上限値。
- NUM_PARTITIONS: パーティションの数。

例：

```
salesforce_read = glueContext.create_dynamic_frame.from_options(  
    connection_type="salesforce",  
    connection_options={  
        "connectionName": "connectionName",  
        "ENTITY_NAME": "Account",
```

```
"API_VERSION": "v60.0",
"PARTITION_FIELD": "SystemModstamp"
"LOWER_BOUND": "TIMESTAMP '2021-01-01 00:00:00 Pacific/Tahiti'"
"UPPER_BOUND": "TIMESTAMP '2023-01-10 00:00:00 Pacific/Tahiti'"
"NUM_PARTITIONS": "10"
}
```

Salesforce への書き込み

前提条件

書き込み先の Salesforce sObject。または AccountCase や などのオブジェクト名が必要になります Opportunity。

Salesforce コネクタは、次の 4 つの書き込みオペレーションをサポートしています。

- INSERT
- UPSERT
- UPDATE
- DELETE

UPSERT 書き込みオペレーションを使用する場合、レコードの外部 ID フィールドを指定するには、を指定 ID_FIELD_NAMES する必要があります。

例

```
salesforce_write = glueContext.write_dynamic_frame.from_options(
    frame=frameToWrite,
    connection_type="salesforce",
    connection_options={
        "connectionName": "connectionName",
        "ENTITY_NAME": "Account",
        "API_VERSION": "v60.0",
        "WRITE_OPERATION": "INSERT"
    }
}
```

Salesforce 接続オプション

Salesforce の接続オプションは次のとおりです。

- ENTITY_NAME (文字列) - (必須) 読み取り/書き込みに使用されます。Salesforce のオブジェクトの名前。
- API_VERSION (文字列) - (必須) 読み取り/書き込みに使用されます。使用する Salesforce Rest API バージョン。
- SELECTED_FIELDS(List<String>) - デフォルト: empty(SELECT *)。読み込みに使用されます。オブジェクトに選択する列。
- FILTER_PREDICATE (文字列) - デフォルト: 空。読み込みに使用されます。Spark SQL 形式である必要があります。
- QUERY (文字列) - デフォルト: 空。読み込みに使用されます。完全な Spark SQL クエリ。
- PARTITION_FIELD (文字列) - 読み取りに使用されます。クエリのパーティション化に使用されるフィールド。
- LOWER_BOUND (文字列) - 読み取りに使用されます。選択したパーティションフィールドの包括的な下限値。
- UPPER_BOUND (文字列) - 読み取りに使用されます。選択したパーティションフィールドの排他的上限値。
- NUM_PARTITIONS (整数) - デフォルト: 1。読み込みに使用されます。読み取り用のパーティションの数。
- IMPORT_DELETED_RECORDS (文字列) - デフォルト: FALSE。読み取りに使用されます。クエリ中に削除レコードを取得するには。
- WRITE_OPERATION (文字列) - デフォルト: INSERT。書き込みに使用されます。値は INSERT、UPDATE、UPSERT、DELETE である必要があります。
- ID_FIELD_NAMES (文字列) - デフォルト: null。UPSERT にのみ使用されます。

Salesforce コネクタの制限事項

Salesforce コネクタの制限は次のとおりです。

- Spark SQL のみがサポートされ、Salesforce SOQL はサポートされていません。
- ジョブのブックマークはサポートされません。

Salesforce の JWT ベアラー OAuth フローを設定する

OAuth 2.0 JSON ウェブトークン と server-to-server の統合を有効にする方法については、Salesforce のパブリックドキュメントを参照してください。 [OAuth](#)

PEM ファイルの証明書とキーのペアの作成

PEM ファイルの証明書とキーのペアを作成する

```
openssl req -newkey rsa:4096 -new -nodes -x509 -days 3650 -keyout key.pem -out cert.pem
```

JWT を使用した Salesforce 接続アプリケーションの作成

1. [Salesforce](#) にログインし、右上の設定歯車をクリックして、**セットアップ** を選択します。
2. 左側の App Manager に移動します。(プラットフォームツール > アプリケーション > App Manager)
3. 新しい接続アプリ を選択します。
4. アプリ名を入力し、残りは自動入力します。
5. OAuth 設定を有効にする のチェックボックスをオンにします。
6. コールバック URL を設定します。JWT には使用されないため、https://localhost. を使用できません。
7. 「デジタル署名の使用」のチェックボックスをオンにします。
8. 前に作成した cert.pem ファイルをアップロードします。
9. 必要なアクセス許可を追加します。
 - a. APIs (api) を使用してユーザーデータを管理します。
 - b. カスタムアクセス許可 (custom_permissions) にアクセスします。
 - c. ID URL サービス (ID、プロフィール、E メール、アドレス、電話) にアクセスします。
 - d. 一意のユーザー識別子 (openid) にアクセスします。
 - e. リクエストはいつでも実行します (refresh_token、オフラインアクセス)。
10. 名前付きユーザーの Issue JSON Web Token (JWT) ベースのアクセストークンのチェックボックスをオンにします。
11. [保存] を選択します。
12. [Continue] を選択します。
13. [コンシューマーの詳細を管理] を選択します。
14. コンシューマーキー (クライアント ID) をコピーします。
15. コンシューマーシークレット (クライアントシークレット) をコピーします。
16. [キャンセル] をクリックします。

JSON ウェブトークン (JWT) の生成

1. キーペアを pkcs12 に変換します (プロンプトが表示されたらエクスポートパスワードを設定します)。

```
openssl pkcs12 -export -in cert.pem -inkey key.pem -name jwtcert > jwtcert.p12
```

2. pkcs12 から Java キーストアを作成します (プロンプトが表示されたら送信先キーストアパスワードを設定し、送信元キーストアパスワードの以前のエクスポートパスワードを指定します)。

```
keytool -importkeystore -srckeystore jwtcert.p12 -destkeystore keystore.jks -srcstoretype pkcs12 -alias jwtcert
```

3. keystore.jks に jwtcert エイリアスが含まれていることを確認します (プロンプトが表示されたら、以前の送信先キーストアのパスワードを入力します)。

```
keytool -keystore keystore.jks -list
```

4. Salesforce ドキュメントに記載されている Java クラス JWTEExample を使用して、署名付きトークンを生成します。

- a. 必要に応じて claimArray の値を編集します。

- claimArray [0] = クライアント ID
- claimArray [1] = salesforce ユーザー ID
- claimArray [2] = salesforce ログイン URL
- claimArray [4] = エポックからのミリ秒単位の有効期限。3660624000000 は 2085-12-31 です。

- b. path/to/keystore を keystore.jks への正しいパスに置き換えます。

- c. keystorepassword を、入力した送信先キーストアパスワードに置き換えます。

- d. privatekeypassword を、入力したソースキーストアのパスワードに置き換えます。

- e. コードをコンパイルします。このコードは、base64 エンコーディングの [Apache Commons Codec](#) によって異なります。

```
javac -classpath "../../commons-codec-1.16.1.jar" JWTEExample.java
```

- f. コードを実行します。

```
java -classpath ".:commons-codec-1.16.1.jar" JWTEExample
```

5. 接続されたアプリと JWT が作成されても、ユーザーはアプリの認証を受ける必要があります。2つのアプローチについては、<https://mannharleen.github.io/2020-03-03-salesforce-jwt/> のステップ 3 を参照してください。

上記の手順が完了したら、Salesforce からアクセストークンを取得するために使用できる JSON ウェブトークン (JWT) を出力する必要があります。

入力例

```
export password for pkcs12: awsglue
destination keystore password for jks: awsglue
source keystore password for jks: awsglue

claimArray[0] = "client-id";
claimArray[1] = "my@email.com";
claimArray[2] = "https://login.salesforce.com";
claimArray[3] = "3660624000000";

path to keystore: ./keystore.jks
keystore password: awsglue
privatekey password: awsglue
```

出力例:

```
eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9.eyJzdWIiOiIxMjM0NTY3ODkwIiwibmFtZSI6IiwiaWF0Ij
```

便利なリンク :

- <https://www.base64encode.org/>
- <https://jwt.io/>
- https://help.salesforce.com/s/articleView?id=sf.remoteaccess_oauth_jwt_flow.htm

JWTExample.java :

```
import org.apache.commons.codec.binary.Base64;
import java.io.*;
import java.security.*;
import java.text.MessageFormat;

public class JWTExample {
```

```
public static void main(String[] args) {

    String header = "{\"alg\":\"RS256\"}";
    String claimTemplate = "'{'iss\": \"{0}\", \"sub\": \"{1}\", \"aud\": \"{2}\",
    \"exp\": \"{3}\"}'";

    try {
        StringBuffer token = new StringBuffer();

        //Encode the JWT Header and add it to our string to sign
        token.append(Base64.encodeBase64URLSafeString(header.getBytes("UTF-8")));

        //Separate with a period
        token.append(".");

        //Create the JWT Claims Object
        String[] claimArray = new String[5];
        claimArray[0] = "value";
        claimArray[1] = "my@email.com";
        claimArray[2] = "https://login.salesforce.com";
        claimArray[3] = Long.toString( ( System.currentTimeMillis()/1000 ) + 300);
        MessageFormat claims;
        claims = new MessageFormat(claimTemplate);
        String payload = claims.format(claimArray);

        //Add the encoded claims object
        token.append(Base64.encodeBase64URLSafeString(payload.getBytes("UTF-8")));

        //Load the private key from a keystore
        KeyStore keystore = KeyStore.getInstance("JKS");
        keystore.load(new FileInputStream("./keystore.jks"), "awsglue".toCharArray());
        PrivateKey privateKey = (PrivateKey) keystore.getKey("jwtcert",
        "awsglue".toCharArray());

        //Sign the JWT Header + "." + JWT Claims Object
        Signature signature = Signature.getInstance("SHA256withRSA");
        signature.initSign(privateKey);
        signature.update(token.toString().getBytes("UTF-8"));
        String signedPayload = Base64.encodeBase64URLSafeString(signature.sign());

        //Separate with a period
        token.append(".");
```

```
//Add the encoded signature
token.append(signedPayload);

System.out.println(token.toString());

} catch (Exception e) {
    e.printStackTrace();
}
}
}
```

AWS Glue Studio での SAP HANA に対する接続

AWS Glue は、SAP HANA のための組み込みサポートを提供します。AWS Glue Studio は、SAP HANA に接続してデータ統合ジョブをオーサリングし、AWS Glue Studio サーバーレス Spark ランタイム上でそれらのジョブを実行するためのビジュアルインターフェイスを提供します。

トピック

- [SAP HANA 接続の作成](#)
- [SAP HANA ソースノードの作成](#)
- [SAP HANA ターゲットノードを作成する](#)
- [詳細オプション](#)

SAP HANA 接続の作成

AWS Glue から SAP HANA に接続するには、SAP HANA 認証情報を作成して AWS Secrets Manager シークレットに保存し、そのシークレットを SAP HANA AWS Glue 接続に関連付ける必要があります。SAP HANA サービスと AWS Glue の間のネットワーク接続を設定する必要があります。

前提条件:

- SAP HANA サービスが Amazon VPC 内にある場合は、トラフィックがパブリックインターネットを経由することなく、AWS Glue ジョブが SAP HANA サービスと通信できるように Amazon VPC を設定します。

Amazon VPC で、AWS Glue がジョブの実行中に使用する [VPC]、[サブネット]、および [セキュリティグループ] を特定または作成します。さらに、SAP HANA エンドポイントとこの場所の間のネットワークトラフィックを許可するように Amazon VPC が設定されているようにする必要があります。

あります。ジョブでは、SAP HANA JDBC ポートとの TCP 接続を確立する必要があります。SAP HANA ポートの詳細については、[SAP HANA のドキュメント](#)を参照してください。ネットワークレイアウトに基づいて、セキュリティグループルール、ネットワーク ACL、NAT ゲートウェイ、およびピアリング接続の変更が必要になる場合があります。

SAP HANA に対する接続を設定するには:

1. AWS Secrets Manager で、SAP HANA 認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
 - [key/value ペア] を選択する際に、*saphanaUsername* という値を持つキー user のペアを作成します。
 - [key/value ペア] を選択する際に、*saphanaPassword* という値を持つキー password のペアを作成します。
2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、将来的に AWS Glue で使用するために、接続名 *connectionName* を維持します。
 - [接続タイプ] を選択する際に、[SAP HANA] を選択します。
 - [SAP HANA URL] を入力する場合は、インスタンスの URL を入力します。

SAP HANA JDBC URL の形式は

```
jdbc:sap://saphanaHostname:saphanaPort?databaseName=saphanaDBname,Parameter
```

です

AWS Glue には次の JDBC URL パラメータが必要です。

- *databaseName* – 接続先の SAP HANA のデフォルトデータベース。
- [AWS Secret] をクリックして、*secretName* を入力します。

AWS Glue SAP HANA 接続を作成した後、AWS Glue ジョブを実行する前に次のステップを実行する必要があります。

- AWS Glue ジョブに関連付けられている IAM ロールに *secretName* を読み取るアクセス許可を付与します。

SAP HANA ソースノードの作成

必要な前提条件

- 前のセクション [the section called “SAP HANA 接続の作成”](#) で説明したように、AWS Secrets Manager シークレットを使用して設定された AWS Glue SAP HANA 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 読み取り元とする SAP HANA テーブル、*tableName*、またはクエリ *targetQuery*。

テーブルは、SAP HANA テーブル名とスキーマ名 (形式: *schemaName.tableName*) を使用して指定できます。テーブルがデフォルトのスキーマ「public」にある場合、スキーマ名と「.」区切り文字は必要ありません。この *tableIdentifier* を呼び出します。データベースは *connectionName* の JDBC URL パラメータとして指定されることに注意してください。

SAP HANA データソースの追加

[データソース — SAP HANA] ノードを追加するには:

1. SAP HANA データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずです。接続を作成する必要がある場合は、[SAP HANA 接続を作成] を選択します。詳細については、前の「[the section called “SAP HANA 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. [SAP HANA ソース] のオプションを選択します。
 - [単一のテーブルを選択] – 単一のテーブルからすべてのデータにアクセスします。
 - [カスタムクエリを入力] — カスタムクエリに基づいて、複数のテーブルからデータセットにアクセスします。
3. 単一のテーブルを選択した場合は、*tableName* を入力します。

[カスタムクエリを入力] を選択した場合は、SQL SELECT クエリを入力します。

4. [SAP HANA のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

SAP HANA ターゲットノードを作成する

必要な前提条件

- 前のセクション [the section called “SAP HANA 接続の作成”](#) で説明したように、AWS Secrets Manager シークレットを使用して設定された AWS Glue SAP HANA 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 書き込み先とする SAP HANA テーブルである *tableName*。

テーブルは、SAP HANA テーブル名とスキーマ名 (形式: *schemaName.tableName*) を使用して指定できます。テーブルがデフォルトのスキーマ「public」にある場合、スキーマ名と「.」区切り文字は必要ありません。この *tableIdentifier* を呼び出します。データベースは *connectionName* の JDBC URL パラメータとして指定されることに注意してください。

SAP HANA データターゲットの追加

[データターゲット — SAP HANA] ノードを追加するには:

1. SAP HANA データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずです。接続を作成する必要がある場合は、[SAP HANA 接続を作成] を選択します。詳細については、前の「[the section called “SAP HANA 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. *tableName* を指定して [テーブル名] を設定します。
3. [Teradata のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

詳細オプション

SAP HANA ノードを作成する際に、高度なオプションを指定できます。これらのオプションは Spark AWS Glue スクリプトのプログラミング時に使用できるオプションと同じです。

「[the section called “SAP HANA 接続”](#)」を参照してください。

AWS Glue Studio での Snowflake への接続

Note

AWS Glue for Spark を使用して、AWS Glue 4.0 以降のバージョンの Snowflake のテーブルからの読み取りとテーブルへの書き込みを行うことができます。AWS Glue ジョブで Snowflake 接続をプログラマ的に設定する方法については、「[Redshift 接続](#)」を参照してください。

AWS Glue には、Snowflake に対するサポートが組み込まれています。AWS Glue Studio には、Snowflake に接続してデータ統合ジョブをオーサリングし、AWS Glue Studio の Spark ランタイム上でジョブをサーバーレスで実行できる、ビジュアルインターフェイスが用意されています。

トピック

- [Snowflake 接続の作成](#)
- [Snowflake ソースノードの作成](#)
- [Snowflake ターゲットノードの作成](#)
- [詳細オプション](#)

Snowflake 接続の作成

AWS Glue Studio で データソース - Snowflake ノードを追加するには、既存の AWS Glue Snowflake 接続を選択するか、新しい接続を作成します。Snowflake に接続するように設定された JDBC タイプの接続ではなく、SNOWFLAKE タイプの接続を選択する必要があります。次の手順に従って、AWS Glue Snowflake 接続を作成します。

Snowflake 接続を作成するには

1. Snowflake でユーザー *snowflakeUser* とパスワード *snowflakePassword* を作成します。
2. このユーザーが操作する Snowflake ウェアハウス *snowflakeWarehouse* を決定します。これを Snowflake の *snowflakeUser* の DEFAULT_WAREHOUSE として設定するか、次のステップのために覚えておきます。
3. AWS Secrets Manager で、Snowflake の認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。

シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。

- [キー/値のペア] をクリックして、sfUser キーを使用して *snowflakeUser* のペアを作成します。
 - [キー/値のペア] をクリックして、sfPassword キーを使用して *snowflakePassword* のペアを作成します。
 - [キー/値のペア] をクリックして、sfWarehouse キーを使用して *snowflakeWarehouse* のペアを作成します。Snowflake にデフォルト設定がある場合、これは必要ありません。
4. AWS Glue データカタログで、「[AWS Glue 接続の追加](#)」にある手順に従って接続を作成します。接続を作成したら、次のステップのために接続名 *connectionName* を保存しておきます。
- [接続タイプ] をクリックして、[Snowflake] を選択します。
 - [Snowflake URL] をクリックして、Snowflake インスタンスのホスト名を指定します。URL では、*account_identifier.snowflakecomputing.com* という形式のホスト名を使用します。
 - [AWS Secret] をクリックして、*secretName* を入力します。

Snowflake ソースノードの作成

必要となる許可

Snowflake データソースを使用する AWS Glue Studio ジョブには追加の許可が必要です。ETL ジョブに許可を追加する方法の詳細については、「[ETL ジョブに必要な IAM アクセス許可を確認する](#)」を参照してください。

SNOWFLAKE AWS Glue 接続は、AWS Secrets Manager シークレットを使用して認証情報を提供します。AWS Glue Studio のジョブおよびデータプレビューのロールには、このシークレットを読み取る許可が必要です。

Snowflake データソースの追加

前提条件:

- Snowflake 認証情報用の AWS Secrets Manager シークレット
- Snowflake タイプの AWS Glue データカタログ接続

データソース - Snowflake ノードを追加するには

1. Snowflake データソースの接続を選択します。既に接続が存在していて、それらの接続から選択できることが前提です。接続の作成が必要な場合は、[Create Snowflake connection] をクリックします。詳細については、「[コネクタと接続の使用に関する概要](#)」を参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。URL、セキュリティグループ、サブネット、アベイラビリティーゾーン、説明、作成時刻 (UTC) と最終更新時刻 (UTC) のタイムスタンプなど、接続についての情報が表示されます。

2. Snowflake ソースオプションを選択します。
 - [Choose a single table] - これは、単一の Snowflake テーブルからアクセスするデータを含むテーブルです。
 - [Enter custom query] - カスタムクエリに基づき、Snowflake の複数のテーブルからデータセットにアクセスできます。
3. 単一のテーブルを選択した場合は、Snowflake スキーマの名前を入力します。

または、[カスタムクエリを入力] を選択します。このオプションを選択すると、Snowflake の複数のテーブルからカスタムデータセットにアクセスできます。このオプションを選択した場合は、Snowflake クエリを入力します。

4. [パフォーマンスとセキュリティ] (オプション) で、次を設定します。
 - [クエリプッシュダウンを有効化] - Snowflake インスタンスに作業をオフロードするかどうかを選択します。
5. [カスタムの Snowflake プロパティ] (オプション) で、必要に応じてパラメータと値を入力します。

Snowflake ターゲットノードの作成

必要となる許可

Snowflake データソースを使用する AWS Glue Studio ジョブには追加の許可が必要です。ETL ジョブに許可を追加する方法の詳細については、「[ETL ジョブに必要な IAM アクセス許可を確認する](#)」を参照してください。

SNOWFLAKE AWS Glue 接続は、AWS Secrets Manager シークレットを使用して認証情報を提供します。AWS Glue Studio のジョブおよびデータプレビューのロールには、このシークレットを読み取る許可が必要です。

Snowflake データターゲットの追加

Snowflake ターゲットノードを作成するには

1. 既存の Snowflake テーブルをターゲットとして選択するか、新しいテーブル名を入力します。
2. データターゲット - Snowflake ターゲットノードを使用する場合、次のオプションから選択できます。
 - APPEND — テーブルがすでに存在している場合は、新しいデータをテーブルにすべて挿入してダンプします。テーブルが存在しない場合は、新規にテーブルを作成し、そこに新しいデータをすべて挿入します。
 - MERGE — 指定した条件に基づいて、AWS Glue がターゲットとなるテーブルのデータを更新するか、またはデータを追加します。

[オプション] を選択します。

- キーと簡単なアクションの選択 — ソースデータとターゲットデータセットとの、マッチングキーとして使用する列を選択します。
 - 一致した場合、次のオプションを指定します。
 - ターゲットのデータセットにあるレコードを、ソースのデータで更新します。
 - ターゲットのデータセットにあるレコードを削除します。
 - 一致しない場合、次のオプションを指定します。
 - ターゲットのデータセットに、新しい行としてソースデータを挿入します。
 - 何もしない。
- カスタム MERGE ステートメントの入力 — その後 [MERGE ステートメントの検証] を選択し、ステートメントが有効か、無効かを検証できます。
- TRUNCATE — 既にテーブルが存在している場合は、ターゲットのテーブルの内容を削除してから、ターゲットのテーブルを削除します。削除が成功してから、すべてのデータを挿入します。テーブルが存在していない場合、テーブルを作成し、すべてのデータを挿入します。削除が成功しなかった場合、操作は失敗します。
- DROP — 既にテーブルが存在している場合は、テーブルのメタデータとデータを削除します。削除が成功してから、すべてのデータを挿入します。テーブルが存在していない場合、テーブルを作成し、すべてのデータを挿入します。削除が成功しなかった場合、操作は失敗します。

詳細オプション

「AWS Glue デベロッパーガイド」の「[Snowflake connections](#)」を参照してください。

AWS Glue Studio での Teradata Vantage に対する接続

AWS Glue は、Teradata Vantage のための組み込みサポートを提供します。AWS Glue Studio は、Teradata に接続してデータ統合ジョブをオーサリングし、AWS Glue Studio サーバーレス Spark ランタイム上でそれらのジョブを実行するためのビジュアルインターフェイスを提供します。

トピック

- [Teradata Vantage 接続の作成](#)
- [Teradata ソースノードの作成](#)
- [Teradata ターゲットノードの作成](#)
- [詳細オプション](#)

Teradata Vantage 接続の作成

AWS Glue から Teradata Vantage に接続するには、Teradata 認証情報を作成して AWS Secrets Manager シークレットに保存し、そのシークレットを AWS Glue Teradata 接続に関連付ける必要があります。

前提条件:

- Amazon VPC を通じて Teradata 環境にアクセスしている場合は、AWS Glue ジョブが Teradata 環境と通信できるように Amazon VPC を設定します。パブリックインターネット経由で Teradata 環境にアクセスすることは推奨されていません。

Amazon VPC で、AWS Glue がジョブの実行中に使用する [VPC]、[サブネット]、および [セキュリティグループ] を特定または作成します。さらに、Teradata インスタンスとこの場所の間のネットワークトラフィックを許可するように Amazon VPC が設定されているようにする必要があります。ジョブでは、Teradata クライアントポートとの TCP 接続を確立する必要があります。Teradata ポートの詳細については、[Teradata のドキュメント](#)を参照してください。

ネットワークレイアウトに応じて、安全な VPC 接続を実現するには、Amazon VPC および他のネットワークサービスの変更が必要な場合があります。AWS 接続の詳細については、Teradata ドキュメントの「[AWS 接続オプション](#)」を参照してください。

AWS Glue Teradata 接続を設定するには:

1. Teradata 設定で、AWS Glue が接続するユーザーとパスワード (*teradataUser* および *teradataPassword*) を識別または作成します。詳細については、Teradata ドキュメントの「[Vantage セキュリティの概要](#)」を参照してください。
2. AWS Secrets Manager で、Teradata 認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
 - [key/value ペア] を選択する際に、*teradataUsername* という値を持つキー *user* のペアを作成します。
 - [key/value ペア] を選択する際に、*teradataPassword* という値を持つキー *password* のペアを作成します。
3. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、次のステップのために接続名 *connectionName* を保存しておきます。
 - [接続タイプ] を選択する際に、[Teradata] を選択します。
 - [JDBC URL] を入力する場合は、インスタンスの URL を入力します。JDBC URL に特定のカンマ区切りの接続パラメータをハードコーディングすることもできます。URL は次の形式に準拠する必要があります:
`jdbc:teradata://teradataHostname/ParameterName=ParameterValue,ParameterName`

サポートされる URL パラメータには次が含まれます。
 - DATABASE – デフォルトでアクセスするホスト上のデータベースの名前。
 - DBS_PORT – データベースポート。非標準ポートで実行する場合に使用されます。
 - [認証情報タイプ] を選択する場合は、[AWS Secrets Manager] を選択し、[AWS シークレット] を *secretName* に設定します。
4. 次の状況では、追加の設定が必要になる場合があります。
 - Amazon VPC の AWS でホストされている Teradata インスタンスの場合
 - Teradata セキュリティ認証情報を定義する AWS Glue 接続に、Amazon VPC 接続に関する情報を提供する必要があります。接続を作成または更新する際に、[ネットワークオプション] で [VPC]、[サブネット]、および [セキュリティグループ] を設定します。

Teradata ソースノードの作成

必要な前提条件

- 前のセクション [the section called “Teradata Vantage 接続の作成”](#) で説明したように、AWS Secrets Manager シークレットを使用して設定された AWS Glue Teradata Vantage 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。
- 読み取り元とする Teradata テーブル、*tableName*、またはクエリ *targetQuery*。

Teradata データソースの追加

[データソース — Teradata] ノードを追加するには:

1. Teradata データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずですが、接続を作成する必要がある場合は、[新しい接続を作成] を選択します。詳細については、前の「[the section called “Teradata Vantage 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. [Teradata ソース] オプションを選択します。
 - [単一のテーブルを選択] – 単一のテーブルからすべてのデータにアクセスします。
 - [カスタムクエリを入力] – カスタムクエリに基づいて、複数のテーブルからデータセットにアクセスします。
3. 単一のテーブルを選択した場合は、*tableName* を入力します。

[カスタムクエリを入力] を選択した場合は、SQL SELECT クエリを入力します。

4. [Teradata のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

Teradata ターゲットノードの作成

必要な前提条件

- 前のセクション [the section called “Teradata Vantage 接続の作成”](#) で説明したように、AWS Secrets Manager シークレットを使用して設定された AWS Glue Teradata Vantage 接続。
- 接続で使用されるシークレットを読み取るためのジョブに対する適切なアクセス許可。

- 書き込み先とする Teradata テーブルである `tableName`。

Teradata データターゲットの追加

[データターゲット — Teradata] ノードを追加するには:

1. Teradata データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずですが、接続を作成する必要がある場合は、[Teradata 接続を作成] を選択します。詳細については、「[コネクタと接続の使用に関する概要](#)」を参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. `tableName` を指定して [テーブル名] を設定します。
3. [Teradata のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

詳細オプション

Teradata ノードを作成する際に、高度なオプションを指定できます。これらのオプションは Spark AWS Glue スクリプトのプログラミング時に使用できるオプションと同じです。

「[the section called “Teradata Vantage 接続”](#)」を参照してください。

AWS Glue Studio での Vertica に対する接続

AWS Glue は、Vertica のための組み込みサポートを提供します。AWS Glue Studio は、Vertica に接続してデータ統合ジョブをオーサリングし、AWS Glue Studio サーバーレス Spark ランタイム上でそれらのジョブを実行するためのビジュアルインターフェイスを提供します。

トピック

- [Vertica 接続の作成](#)
- [Vertica ソースノードの作成](#)
- [Vertica ターゲットノードの作成](#)
- [詳細オプション](#)

Vertica 接続の作成

前提条件:

- データベースの読み取りおよび書き込み時に一時ストレージとして使用する Amazon S3 バケットまたはフォルダ。 *tempS3Path* によって参照されます。

Note

AWS Glue ジョブデータプレビューで Vertica を使用する場合、一時ファイルは *tempS3Path* から自動的に削除されない場合があります。一時ファイルを確実に削除するには、[データプレビュー] ペインで [セッションの終了] を選択して、データプレビューセッションを直接終了します。

データプレビューセッションが直接終了することを保証できない場合は、Amazon S3 ライフサイクル構成を設定して古いデータを削除することを検討してください。ジョブの最大実行時間にマージンを加えた値に基づいて、49 時間を超える時間が経過しているデータを削除することをお勧めします。Amazon S3 ライフサイクルの設定の詳細については、Amazon S3 ドキュメントの「[ストレージのライフサイクルの管理](#)」を参照してください。

- AWS Glue ジョブロールに関連付けることができる、Amazon S3 パスに対する適切な許可を持つ IAM ポリシー。
- Vertica インスタンスが Amazon VPC 内にある場合は、トラフィックがパブリックインターネットを経由することなく、AWS Glue ジョブが Vertica インスタンスと通信できるように Amazon VPC を設定します。

Amazon VPC で、AWS Glue がジョブの実行中に使用する [VPC]、[サブネット]、および [セキュリティグループ] を特定または作成します。さらに、Vertica インスタンスとこの場所の間のネットワークトラフィックを許可するように Amazon VPC が設定されているようにする必要があります。ジョブは、Vertica クライアントポート (デフォルトは 5433) との TCP 接続を確立する必要があります。ネットワークレイアウトに基づいて、セキュリティグループルール、ネットワーク ACL、NAT ゲートウェイ、およびピアリング接続の変更が必要になる場合があります。

Vertica に対する接続を設定するには:

1. AWS Secrets Manager で、Vertica 認証情報、*verticaUsername* および *verticaPassword* を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。

- [key/value ペア] を選択する際に、*verticaUsername* という値を持つキー user のペアを作成します。
 - [key/value ペア] を選択する際に、*verticaPassword* という値を持つキー password のペアを作成します。
2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、次のステップのために接続名 *connectionName* を保存しておきます。
- [接続タイプ] を選択する際に、[Vertica] を選択します。
 - [Vertica ホスト] を選択する場合は、Vertica インストールのホスト名を入力します。
 - [Vertica ポート] を選択すると、Vertica インストールを使用できるポートが選択されます。
 - [AWS Secret] をクリックして、*secretName* を入力します。
3. 次の状況では、追加の設定が必要になる場合があります。
- Amazon VPC の AWS でホストされている Vertica インスタンスの場合
 - Vertica セキュリティ認証情報を定義する AWS Glue 接続に、Amazon VPC 接続に関する情報を提供します。接続を作成または更新する際に、[ネットワークオプション] で [VPC]、[サブネット]、および [セキュリティグループ] を設定します。

AWS Glue ジョブを実行する前に、次のステップを実行する必要があります。

- AWS Glue ジョブの許可に関連付けられた IAM ロールを *tempS3Path* に付与します。
- AWS Glue ジョブに関連付けられている IAM ロールに *secretName* を読み取るアクセス許可を付与します。

Vertica ソースノードの作成

必要な前提条件

- 前のセクション「[the section called “Vertica 接続の作成”](#)」で説明したように、Vertica タイプの AWS Glue データカタログ接続である *connectionName* と、一時的な Amazon S3 の場所である *tempS3Path*。
- 読み取り元とする Vertica テーブル、*tableName*、またはクエリ *targetQuery*。

Vertica データソースの追加

[データソース — Vertica] ノードを追加するには:

1. Vertica データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずですが、接続を作成する必要がある場合は、[Vertica 接続を作成] を選択します。詳細については、前の「[the section called “Vertica 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. テーブルを含む [データベース] を選択します。
3. [Amazon S3 のステージングエリア] を選択し、S3A URI を *tempS3Path* に入力します。
4. [Vertica ソース] を選択します。
 - [単一のテーブルを選択] – 単一のテーブルからすべてのデータにアクセスします。
 - [カスタムクエリを入力] — カスタムクエリに基づいて、複数のテーブルからデータセットにアクセスします。
5. 単一のテーブルを選択した場合は、*tableName* を入力し、オプションで [スキーマ] を選択します。

[カスタムクエリを入力] を選択した場合は、SQL SELECT クエリを入力し、オプションで [スキーマ] を選択します。

6. [Vertica のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

Vertica ターゲットノードの作成

必要な前提条件

- 前のセクション「[the section called “Vertica 接続の作成”](#)」で説明したように、Vertica タイプの AWS Glue データカタログ接続である *connectionName* と、一時的な Amazon S3 の場所である *tempS3Path*。

Vertica データターゲットの追加

[データターゲット — Vertica] ノードを追加するには:

1. Vertica データソース用の接続を選択します。すでに作成したので、ドロップダウンに表示されているはずですが、接続を作成する必要がある場合は、[Vertica 接続を作成] を選択します。詳細については、前の「[the section called “Vertica 接続の作成”](#)」セクションを参照してください。

接続を選択したあとは、[プロパティを表示] をクリックすると、接続のプロパティを表示できます。

2. テーブルを含む [データベース] を選択します。
3. [Amazon S3 のステージングエリア] を選択し、S3A URI を *tempS3Path* に入力します。
4. *tableName* を入力し、オプションで [スキーマ] を選択します。
5. [Vertica のカスタムプロパティ] で、必要に応じてパラメータと値を入力します。

詳細オプション

Vertica ノードを作成する際に、高度なオプションを指定できます。これらのオプションは Spark AWS Glue スクリプトのプログラミング時に使用できるオプションと同じです。

「[the section called “Vertica 接続”](#)」を参照してください。

AWS Glue Studio でのコネクタと接続の使用

AWS Glue には、JDBC 接続を使用する最も一般的なデータストア (Amazon Redshift、Amazon Aurora、Microsoft SQL Server、MySQL、MongoDB、PostgreSQL など) に対するサポートが組み込まれています。また、AWS Glue では、抽出、変換、ロード (ETL) ジョブ用にカスタム JDBC ドライバも使用できます。SaaS アプリケーションなど、ネイティブにサポートされていないデータストアに対しては、コネクタを使用することができます。

コネクタとは AWS Glue Studio 内でデータストアに対するアクセスを支援するための、オプションのコードパッケージです。AWS Marketplace で提供されている複数のコネクタをサブスクライブすることができます。

ETL ジョブを作成する場合、ネイティブにサポートされているデータストア、からのコネクタ AWS Marketplace、または独自のカスタムコネクタを使用できます。コネクタを使用するには、最初にコネクタのための接続を作成する必要があります。接続には、特定のデータストアに接続するために必要なプロパティが含まれます。ETL ジョブでは、データソースおよびデータターゲットとの接続

を使用します。コネクタと接続は、データストアへのアクセスを容易にするために連携して動作します。

トピック

- [コネクタと接続の使用に関する概要](#)
- [AWS Glue Studio にコネクタを追加する](#)
- [使用可能な接続](#)
- [コネクタ用の接続を作成する](#)
- [カスタムコネクタを使用したジョブのオーサリング](#)
- [コネクタと接続を管理する](#)
- [カスタムコネクタの開発](#)
- [AWS Glue Studio でのコネクタおよび接続の使用に関する制約事項](#)

コネクタと接続の使用に関する概要

接続には、特定のデータストアに接続するために必要なプロパティが含まれます。作成した接続は AWS Glue Data Catalog に保存されます。コネクタを選択し、そのコネクタに基づいて接続を作成します。

でネイティブにサポートされていないデータストアのコネクタをサブスクライブし AWS Marketplace、接続の作成時にそれらのコネクタを使用できます。さらに、デベロッパーは独自のコネクタを作成し、それを接続の作成に使用することもできます。

Note

AWS Marketplace でカスタムまたはコネクタを使用して作成された接続は、AWS Glue StudioAWS Glueタイプがに設定された状態でコンソールに表示されます。UNKNOWN

以下のステップで、AWS Glue Studio でコネクタを使用するための全体的なプロセスを説明します。

1. でコネクタをサブスクライブするか AWS Marketplace、AWS Glue Studio独自のコネクタを開発してアップロードします。詳細については、「[AWS Glue Studio にコネクタを追加する](#)」を参照してください。

2. コネクタの使用法に関する情報を確認します。この情報は、コネクタ製品ページの [Usage] (使用方法) タブに表示されます。たとえば、この製品ページ「[AWS GlueConnector for Google](#)」の「使用状況」タブをクリックすると BigQuery、「その他のリソース」セクションに、このコネクタの使用に関するブログへのリンクが表示されます。他のコネクタについては、コネクタの製品ページ [Cloudwatch Logs connector for AWS Glue](#) のように、[Overview] (概要) セクションで使用手順に関するリンクが表示されます。
3. 接続を作成します。使用するコネクタを選択し、ログイン認証情報、URI 文字列、仮想プライベートクラウド (VPC) 情報など、接続に関する追加情報を提供します。詳しくは、「[コネクタ用の接続を作成する](#)」を参照してください。
4. ジョブ用に IAM ロールを作成します。ジョブは、作成時に指定する [IAM role] (IAM ロール) のアクセス許可があることを想定します。この IAM ロールには、データストアを承認し、そこからのデータ抽出、およびデータを書き込むために必要なアクセス許可を有する必要があります。
5. ETL ジョブを作成し、その ETL ジョブのためにデータソースプロパティを設定します。カスタムコネクタプロバイダーの指示に従って、接続オプションと認証情報を指定します。詳しくは、「[カスタムコネクタを使用したジョブのオーサリング](#)」を参照してください。
6. [ビジュアル ETL と AWS Glue Studio](#) での説明を参考に、変換を追加するか新しいデータストアを追加しながら、ETL ジョブをカスタマイズします。
7. データターゲットにコネクタを使用している場合は、ETL ジョブ用のデータターゲットプロパティを設定します。カスタムコネクタプロバイダーの指示に従って、接続オプションと認証情報を指定します。詳しくは、「[the section called “カスタムコネクタを使用したジョブのオーサリング”](#)」を参照してください。
8. [ジョブのプロパティを変更する](#) での説明のように、ジョブのプロパティを構成して、ジョブの実行環境をカスタマイズします。
9. ジョブを実行します。

AWS Glue Studio にコネクタを追加する

コネクタとは、データストアと AWS Glue 間の通信を容易にするための、一連のコードのことです。で提供されているコネクタを購読することも AWS Marketplace、独自のカスタムコネクタを作成することもできます。

トピック

- [AWS Marketplace コネクタを購読する](#)
- [カスタムコネクタを作成する](#)

AWS Marketplace コネクタを購読する

AWS Glue Studio AWS Marketplaceからコネクタを簡単に追加できます。

AWS Marketplace からコネクタを追加するには AWS Glue Studio

1. AWS Glue Studio コンソールのナビゲーションペインで、[Connectors] (コネクタ) をクリックします。
2. [Connectors] (コネクタ) ページで、[Go to AWS Marketplace] をクリックします。
3. で AWS Marketplace、「おすすめ商品」で、使用したいコネクタを選択します。表示されているコネクタの1つを選択するか、または検索することができます。コネクタを、名前またはタイプにより検索したり、オプションを使用して検索結果を絞り込むことができます。

表示されているコネクタのいずれかを使用するには、[View product] (製品を表示する) をクリックします。検索結果からコネクタを見つけた場合は、そのコネクタの名前を選択します。

4. コネクタの製品ページで、そのコネクタに関するタブを開いて情報を表示します。そのコネクタを購入する場合は、[Continue to Subscribe] (続行してサブスクライブする) をクリックします。
5. 支払い情報を入力し、[Continue to Configure] (設定に進む) をクリックします。
6. [Configure this software] (このソフトウェアを設定する) ページで、デプロイの方法と使用するコネクタのバージョンを選択します。[Continue to Launch] (続行して起動する) をクリックします。
7. [Launch this software] (このソフトウェアを起動する) ページでは、コネクタプロバイダから提供される [Usage Instructions] (使用手順) を確認することができます。次に進む準備ができたなら、[AWS Glue Studio で接続をアクティブにする] を選択します。

少し待機すると、コンソールに、AWS Glue Studio の [Create marketplace connection] (マーケットプレイス接続の作成) ページが表示されます。

8. [コネクタ用の接続を作成する](#) の説明を参考に、このコネクタを使用する接続を作成します。

または、[Activate connector only] (アクティブなコネクタのみ) を選択し、この時点での接続の作成をスキップすることも可能です。後にコネクタを使用する際には、先に接続を作成する必要があります。

カスタムコネクタを作成する

独自のコネクタを作成し、そのコネクタのコードを AWS Glue Studio にアップロードすることもできます。

カスタムコネクタは、AWS Glue Spark ランタイム API を介して AWS Glue Studio に組み込まれます。Spark、Athena、または JDBC インターフェイスと準拠している任意のコネクタを、AWS Glue Spark ランタイムを使用してプラグインすることができます。これにより、カスタムコネクタで利用できる任意の接続オプションを渡すことができます。

[AWS Glue 接続](#)でのすべての接続プロパティをカプセル化して、ETL ジョブにその接続名を指定することができます。Data Catalog 接続を統合することで、単一の Spark アプリケーションからの複数の呼び出しや、異なるアプリケーション間での同じ接続プロパティの使用ができるようになります。

接続には、他のオプションを指定することもできます。AWS Glue Studio が生成するジョブスクリプトには、指定された接続オプションを使用してコネクタをプラグインする接続を使用する Datasource エントリが含まれています。例:

```
Datasource = glueContext.create_dynamic_frame.from_options(connection_type =
"custom.jdbc", connection_options = {"dbTable":"Account","connectionName":"my-custom-
jdbc-
connection"}, transformation_ctx = "DataSource0")
```

AWS Glue Studio にカスタムコネクタを追加する方法

1. カスタムコネクタ用のコードを作成します。詳しくは、「[カスタムコネクタの開発](#)」を参照してください。
2. コネクタに、AWS Glue 機能に対するサポートを追加します。以下に、これらの機能について、さらにそれらが AWS Glue Studio で生成されたジョブスクリプトでどのように使用されるかに関する例をいくつか挙げてみます。
 - データ型のマッピング – コネクタは、基盤データストアから列を読み込む際に、列をタイプキャストすることができます。例えば {"INTEGER":"STRING"} の dataTypeMapping では、レコードの解析と DynamicFrame の構築時に、Integer 型のすべての列を String 型の列に変換します。これにより、ユーザーは任意のタイプに列をキャストすることができます。

```
DataSource0 = glueContext.create_dynamic_frame.from_options(connection_type
= "custom.jdbc", connection_options = {"dataTypeMapping":{"INTEGER":"STRING"}",
connectionName":"test-connection-jdbc"}, transformation_ctx = "DataSource0")
```

- 並列読み取りのパーティション化 – AWS Glue は、列にあるデータを分割することで、データストアから並列データを読み取れるようにします。パーティション列、パーティションの下限、パーティションの上限、およびパーティション数を指定する必要があります。この機

能により、データの並列処理と、Spark アプリケーションに割り当てる複数の Spark エグゼキュータの使用が可能になります。

```
DataSource0 = glueContext.create_dynamic_frame.from_options(connection_type = "custom.jdbc", connection_options = {"upperBound":"200","numPartitions":"4", "partitionColumn":"id","lowerBound":"0","connectionName":"test-connection-jdbc"}, transformation_ctx = "DataSource0")
```

- 認証情報の保存に使用 AWS Secrets Manager — Data Catalog 接続には、secretId AWS Secrets Managerに保存されているシークレット用の機能を含めることもできます。AWS シークレットには認証情報と認証情報を安全に保存し、AWS Glue 実行時に提供できます。または、以下に示すように、Spark スクリプトから secretId を指定することもできます。

```
DataSource = glueContext.create_dynamic_frame.from_options(connection_type = "custom.jdbc", connection_options = {"connectionName":"test-connection-jdbc", "secretId"-> "my-secret-id"}, transformation_ctx = "DataSource0")
```

- 行述語と列射影によるソースデータのフィルタリング – AWS Glue Spark ランタイムでは、行述語と列射影を使用しながら、SQL クエリをプッシュダウンすることでソースにあるデータをフィルタリングすることもできます。これにより ETL ジョブは、フィルタリングされたデータを、プッシュダウンをサポートするデータストアからより迅速にロードできます。SELECT id, name, department FROM department WHERE id < 200. は、JDBC データ・ソースにプッシュダウンされた SQL クエリの例です。

```
DataSource = glueContext.create_dynamic_frame.from_options(connection_type = "custom.jdbc", connection_options = {"query":"SELECT id, name, department FROM department WHERE id < 200","connectionName":"test-connection-jdbc"}, transformation_ctx = "DataSource0")
```

- ジョブのブックマーク – AWS Glue は、JDBC ソースからのデータの増分ロードをサポートしています。AWS Glue は、最後に処理された (データストアからの) レコードを追跡し、後続の ETL ジョブ実行で新しいデータレコードを処理します。ジョブのブックマークは、その列が順番に増減しているのであれば、ブックマークキーのデフォルト列としてプライマリキーを使用します。ジョブのブックマークの詳細については、AWS Glue デベロッパーガイドの「[ジョブ ブックマーク](#)」を参照してください。

```
DataSource0 = glueContext.create_dynamic_frame.from_options(connection_type = "custom.jdbc", connection_options = {"jobBookmarkKeys":["empno"], "jobBookmarkKeysSortOrder"
```

```
:"asc", "connectionName":"test-connection-jdbc"}, transformation_ctx =  
"DataSource0")
```

3. カスタムコネクタを JAR ファイルとしてパッケージ化し、そのファイルを Amazon S3 にアップロードします。
4. カスタムコネクタをテストします。詳細については、「[Glue カスタムコネクタ:ローカル検証テストガイド](#)」の説明を参照してください。GitHub
5. AWS Glue Studio コンソールのナビゲーションペインで、[Connectors] (コネクタ) をクリックします。
6. [Connectors] (コネクタ) ページで、[Create custom connector] (カスタムコネクタを作成) をクリックします。
7. [Create custom connector] (カスタムコネクタの作成) ページで、次の情報を入力します。
 - カスタムコード JAR ファイルの Amazon S3 内ロケーションへのパス。
 - AWS Glue Studio によって使用されるコネクタの名前。
 - コネクタのタイプ (JDBC、Spark、または Athena)。
 - コネクタを使用するために AWS Glue Studio が呼び出す、カスタムコード内のエントリポイントの名前。
 - JDBC コネクタの場合、このフィールドは JDBC ドライバのクラス名です。
 - Spark コネクタの場合、このフィールドは、データソースの完全修飾クラス名、またはそのエイリアスである必要があります。これは、format 演算子を使用して Spark データソースをロードする際に使用します。
 - (JDBC のみ) データストアの JDBC 接続で使用されるベース URL。
 - (オプション) カスタムコネクタの説明。
8. [Create connector] (コネクタを作成) をクリックします。
9. [コネクタ用の接続を作成する](#) の説明を参考に、[Connectors] (コネクタ) ページで、そのコネクタを使用する接続を作成します。

使用可能な接続

コネクタの接続を作成する際に、次の接続を使用できます。

- Amazon Aurora – セキュリティ、バックアップと復元、メモリ内アクセラレーションを組み込んだ、スケーラブルで高性能なリレーショナルデータベースエンジン。

- Amazon DocumentDB – MongoDB と SQL API をサポートする、スケーラブルで可用性の高いフルマネージドのドキュメントデータベースサービス。
- Amazon Redshift – MongoDB および SQL API をサポートする、スケーラブルで可用性の高いフルマネージドのドキュメントデータベースサービス。
- Azure SQL – スケーラブルで信頼性が高く安全なデータストレージと管理機能を提供する、Microsoft Azure のクラウドベースのリレーショナルデータベースサービス。
- Cosmos DB – スケーラブルで高性能なデータストレージとクエリ機能を提供する、Microsoft Azure のグローバルに分散されたクラウドデータベースサービス。
- Google BigQuery – 大規模なデータセットで高速 SQL クエリを実行するためのサーバーレスクラウドデータウェアハウス。
- JDBC – 接続と、データ接続とのインタラクションに Java API を使用するリレーショナルデータベース管理システム (RDBMS)。
- Kafka – リアルタイムのデータストリーミングとメッセージングに使用されるオープンソースのストリーム処理プラットフォーム。
- MariaDB – 強化されたパフォーマンス、スケーラビリティ、および機能を提供する、コミュニティによって開発された MySQL のフォーク。
- MongoDB – 高いスケーラビリティ、柔軟性、パフォーマンスを提供する、クロスプラットフォームのドキュメント指向データベース。
- MongoDB Atlas – MongoDB のデプロイの管理とスケーリングを簡素化する、MongoDB のクラウドベースの Database as a Service (DBaaS) オファリング。
- Microsoft SQL Server – 堅牢なデータストレージ、分析、レポート機能を提供する、Microsoft のリレーショナルデータベース管理システム (RDBMS)。
- MySQL – ウェブアプリケーションで広く使用されており、その信頼性とスケーラビリティで知られる、オープンソースのリレーショナルデータベース管理システム (RDBMS)。
- ネットワーク- ネットワークデータソースは、データ統合プラットフォームによってアクセスできる、ネットワークからアクセス可能なリソースまたはサービスを表します。
- OpenSearch – OpenSearch データソースは、 に接続してデータを取り込む OpenSearch ことのできるアプリケーションです。
- Oracle – 堅牢なデータストレージ、分析、レポート機能を提供する、Oracle Corporation のリレーショナルデータベース管理システム (RDBMS)。
- PostgreSQL – 堅牢なデータストレージ、分析、レポート機能を提供する、オープンソースのリレーショナルデータベース管理システム (RDBMS)。

- Salesforce – Salesforce は、販売、カスタマーサービス、e コマースなどを支援する顧客関係管理 (CRM) ソフトウェアを提供しています。Salesforce ユーザーの場合は、Salesforce アカウント AWS Glue に接続できます。その後、ETL ジョブのデータソースまたは送信先として Salesforce を使用できます。これらのジョブを実行して、Salesforce と AWS のサービス、またはその他のサポートされているアプリケーション間でデータを転送します。
- SAP HANA – 高速データ処理、高度な分析、リアルタイムのデータ統合を提供する、インメモリデータベースおよび分析プラットフォーム。
- Snowflake – スケーラブルで高性能なデータストレージと分析サービスを提供するクラウドベースのデータウェアハウス。
- Teradata – 高性能のデータストレージ、分析、レポート機能を提供する、リレーショナルデータベース管理システム (RDBMS)。
- Vertica – 高速なクエリパフォーマンス、高度な分析、およびスケーラビリティを提供する、ビッグデータ分析用に設計されたカラム指向の分析データウェアハウス。

コネクタ用の接続を作成する

AWS Glue 接続は、特定のデータストアの接続情報を保存する Data Catalog オブジェクトです。接続には、ログイン認証情報、URI 文字列、Virtual Private Cloud (VPC) 情報などが含まれます。Data Catalog で接続を作成すると、ジョブを作成するたびにすべての接続の詳細を指定する必要がなくなります。

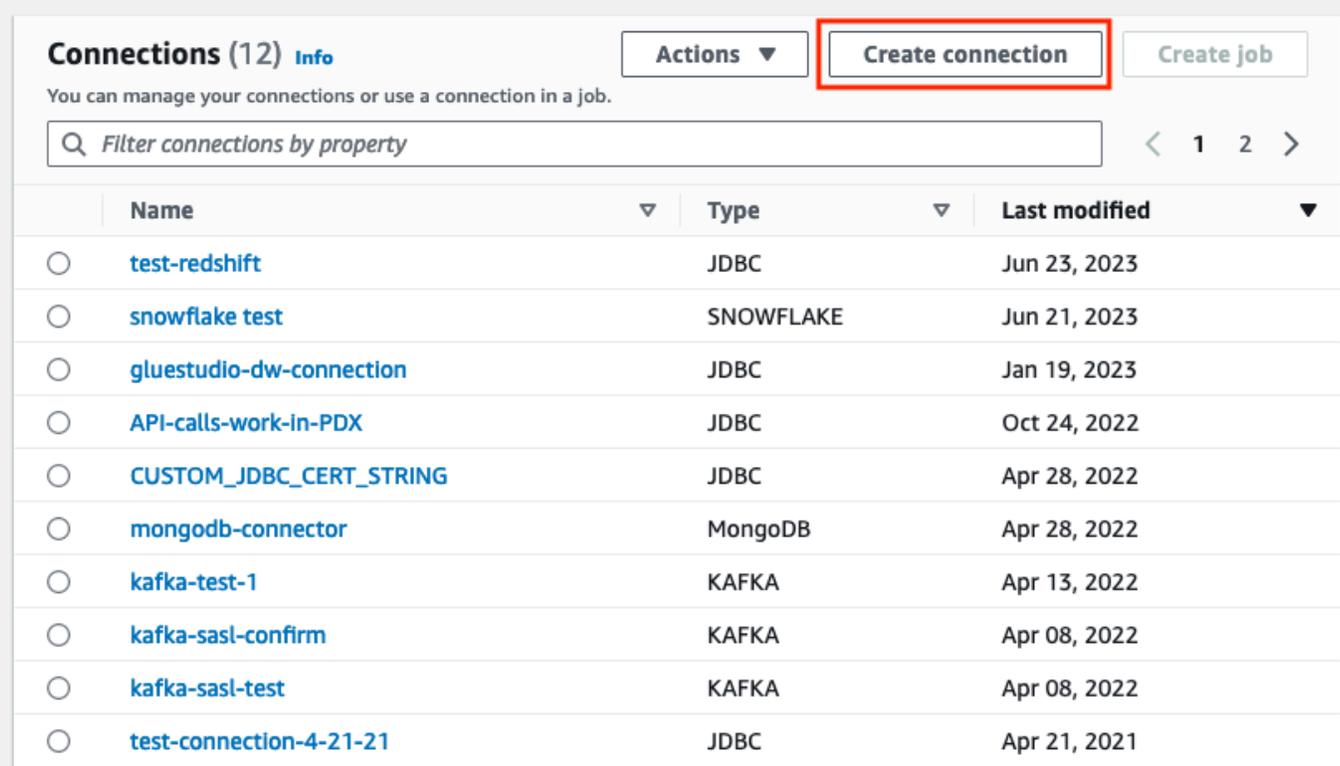
コネクタのために接続を作成するには

1. AWS Glue Studio コンソールのナビゲーションペインで、[Connectors] (コネクタ) をクリックします。[接続] セクションで、[接続を作成] を選択します。
2. [データ接続を作成] ウィザードのステップ 1 で、接続を作成するデータソースを選択します。使用可能なデータソースを表示するには、以下のような複数の方法があります。
 - タブを選択して、使用可能なデータソースをフィルタリングします。デフォルトでは、[すべてのコネクタ] が選択されています。
 - [リスト] に切り替えてデータソースをリストとして表示するか、または [グリッド] に切り替えてグリッドレイアウトで使用可能なコネクタを表示します。
 - データソースのリストを絞り込むには、検索バーを使用します。入力すると、検索に一致するソースが表示され、一致しないソースはビューから削除されます。

データソースを選択したら、[次へ] を選択します。

3. ウィザードのステップ 2 で接続を設定します。

接続の詳細を入力します。選択したコネクタのタイプに応じて、追加情報の入力を促されます。



The screenshot shows the AWS Glue Connections console. At the top, there is a header with 'Connections (12) Info', an 'Actions' dropdown menu, and a 'Create connection' button highlighted with a red box. Below the header is a search bar with the placeholder text 'Filter connections by property'. The main content is a table listing connections with columns for Name, Type, and Last modified. The table contains 12 rows of connection data.

	Name	Type	Last modified
<input type="radio"/>	test-redshift	JDBC	Jun 23, 2023
<input type="radio"/>	snowflake test	SNOWFLAKE	Jun 21, 2023
<input type="radio"/>	gluestudio-dw-connection	JDBC	Jan 19, 2023
<input type="radio"/>	API-calls-work-in-PDX	JDBC	Oct 24, 2022
<input type="radio"/>	CUSTOM_JDBC_CERT_STRING	JDBC	Apr 28, 2022
<input type="radio"/>	mongodb-connector	MongoDB	Apr 28, 2022
<input type="radio"/>	kafka-test-1	KAFKA	Apr 13, 2022
<input type="radio"/>	kafka-sasl-confirm	KAFKA	Apr 08, 2022
<input type="radio"/>	kafka-sasl-test	KAFKA	Apr 08, 2022
<input type="radio"/>	test-connection-4-21-21	JDBC	Apr 21, 2021

4. [データ接続を作成] ウィザードのステップ 1 で、接続を作成するデータソースを選択します。使用可能なデータソースを表示するには、いくつかの方法があります。デフォルトでは、使用可能なすべてのデータソースがグリッドレイアウトで表示されます。以下の操作も可能です。

- [リスト] に切り替えてデータソースをリストとして表示するか、または [グリッド] に切り替えてグリッドレイアウトで使用可能なコネクタを表示します。
- データソースのリストを絞り込むには、検索バーを使用します。入力すると、検索に一致するソースが表示され、一致しないソースはビューから削除されます。

Choose data source



The screenshot shows the 'Choose data source' dialog box. It has a title 'Choose data source' and a sub-header 'Data sources (21)'. Below the sub-header is a search bar with the placeholder text 'Find data sources'. To the right of the search bar are two buttons: 'Grid' (highlighted in blue) and 'List'.

データソースを選択したら、[次へ] を選択します。

5. ウィザードのステップ 2 で接続を設定します。

接続の詳細を入力します。選択したコネクタの種類によっては、追加の接続情報の入力が必要になる場合があります。これには次が含まれる場合があります。

- [接続の詳細] – これらのフィールドは、接続先のデータソースに応じて変わります。例えば、Amazon DocumentDB データベースに接続しようとしている場合は、Amazon DocumentDB URL を入力します。Amazon Aurora に接続しようとしている場合は、データベースインスタンスを選択し、データベース名を入力します。Amazon Aurora に必要な [接続の詳細] は次のとおりです。

The screenshot shows the 'Configure connection' wizard in the AWS Glue console. The breadcrumb navigation is 'AWS Glue > Connectors > Create connection'. The left sidebar shows the progress: Step 1 'Choose data source', Step 2 'Configure connection' (active), Step 3 'Set properties', and Step 4 'Review and create'. The main content area is titled 'Configure connection' and contains a 'Connection details' section. This section has a dropdown for 'Database instances' with the text 'Provisioned Amazon Relational Database Service instances.' and 'Choose one JDBC URL'. Below it is a 'Database name' text input field. The 'Credential type' section has two radio buttons: 'Username and password' (selected) and 'AWS Secrets Manager'. Below that are 'Username' and 'Password' text input fields. At the bottom right of the form are 'Cancel', 'Previous', and 'Next' buttons.

- 認証情報タイプ – [ユーザー名とパスワード]、または [AWS Secrets Manager] のいずれかを選択します。要求された認証情報を入力します。
 - JDBC を使用するコネクタの場合は、データストアに JDBC URL を作成するために必要な情報を入力します。
 - Virtual Private Cloud (VPC) を使用する場合は、VPC のネットワーク情報を入力します。
6. ウィザードのステップ 3 で接続のプロパティを設定します。このステップのオプション部分として、説明とタグを追加できます。名前は必須であり、デフォルト値が事前に入力されています。[Next] を選択します。
 7. 接続ソース、詳細、プロパティを確認します。変更を加える必要がある場合は、ウィザードのステップで [編集] を選択します。準備ができたら、[接続を作成] を選択します。

[Create connection] (接続の作成) を選択します。

再び [Connectors] (コネクタ) ページが表示され、作成された接続が情報バナーに表示されます。AWS Glue Studio ジョブで接続を使用することができるようになりました。

Kafka 接続を作成する

Kafka 接続を作成するときには、ドロップダウンメニューから Kafka を選択すると追加設定が表示され、構成できるようになります。

- Kafka クラスターの詳細
- 認証
- 暗号化
- ネットワークオプション

Kafka クラスターの詳細を設定する

1. クラスターの場所を選択します。Amazon Managed Streaming for Apache Kafka (MSK) クラスターまたは Customer managed Apache Kafka クラスターから選択できます。Amazon Managed Streaming for Apache Kafka の詳細については、[Amazon Managed Streaming for Apache Kafka \(MSK\)](#) を参照してください。

Note

Amazon Managed Streaming for Apache Kafka は TLS および SASL/SCRAM-SHA-512 認証方法のみをサポートします。

Kafka cluster details [Info](#)

Cluster location

Amazon managed streaming for Apache Kafka (MSK)

Customer managed Apache Kafka

Kafka bootstrap server URLs [Info](#)

A comma-separated list of bootstrap server URLs. Include the port number.

Enter list of URLs, separated by commas

Example: b-1.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094, b-2.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094, b-3.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094

2. Kafka ブートストラップサーバーの URL を入力します。各サーバーをカンマで区切ることで、複数のサーバーを入力できます。URL の末尾に `:<port number>` を足してポート番号を入力します。

例: `b-1.vpc-test-2.034a88o.kafka-us-east-1.amazonaws.com:9094`

認証方法の選択

Authentication [Info](#)

Authentication method

Choose authentication method

AWS Glue は、認証用の Simple Authentication and Security Layer (SASL) フレームワークをサポートしています。SASL フレームワークはさまざまな認証メカニズムをサポートしており、AWS Glue は SCRAM プロトコル (ユーザー名およびパスワード)、GSSAPI プロトコル (Kerberos プロトコル)、PLAIN プロトコル (ユーザー名およびパスワード) を提供します。

ドロップダウンメニューから認証方法を選択するときには、次のクライアント認証方法を選択できません。

- None - 認証なし。これは、テスト目的で接続する場合に便利です。
- SASL/SCRAM-SHA-512 - これを選択すると、この認証方法のための認証情報が指定されます。2つのオプションがあります。
- AWS Secrets Manager を使用 (推奨) - このオプションを選択すると、認証情報が AWS Secrets Manager に保存され、必要になった際に AWS Glue からその情報へアクセスすることを許可できます。SSL または SASL 認証の認証情報を格納するシークレットを指定します。

Authentication Info

Authentication method
SASL/SCRAM-SHA-512

Authentication credentials

Use AWS Secrets Manager (recommended)
Store your token in AWS Secrets Manager, and let AWS Glue access it when needed.

Provide username and password directly
Provide your username and password directly to AWS Glue.

Secret from [AWS Secrets Manager](#)

Search secret by name or type ARN

- ユーザー名とパスワードを直接指定します。
- SASL/GSSAPI (Kerberos) - このオプションを選択すると、キータブファイル、krb5.conf ファイルの場所を選択して、Kerberos プリンシパル名と Kerberos サービス名を入力することができます。キータブファイルと krb5.conf ファイルの場所は、Simple Storage Service (Amazon S3) がある場所の中にする必要があります。MSK は SASL/GSSAPI をまだサポートしていないため、このオプションは Customer Managed Apache Kafka クラスターでのみ使用できます。詳細については、[MIT Kerberos ドキュメント: キータブ](#)を参照してください。
- SASL/PLAIN - 認証用の認証情報を指定するためにこの認証方法を選択します。2つのオプションがあります。
 - AWS Secrets Manager を使用 (推奨) - このオプションを選択すると、認証情報が AWS Secrets Manager に保存され、必要になった際に AWS Glue からその情報へアクセスすることを許可できます。SSL または SASL 認証の認証情報を格納するシークレットを指定します。
 - ユーザー名とパスワードを直接指定します。

- SSL クライアント認証 - このオプションを選択すると、Simple Storage Service (Amazon S3) を参照することで Kafka クライアントキーストアの場所を選択できます。オプションで、Kafka クライアントキーストアのパスワードと Kafka クライアントキーのパスワードを入力できます。

Authentication [Info](#)

Authentication method
SSL client authentication ▼

Kafka client keystore location
s3://bucket/prefix/object [View](#) [Browse S3](#)

Path must be in the form s3://bucket/prefix/path/. It must end with the file name and .jks extension.

Kafka client keystore password - optional
Enter password

Kafka client key password - optional
Enter password

暗号化設定の構成

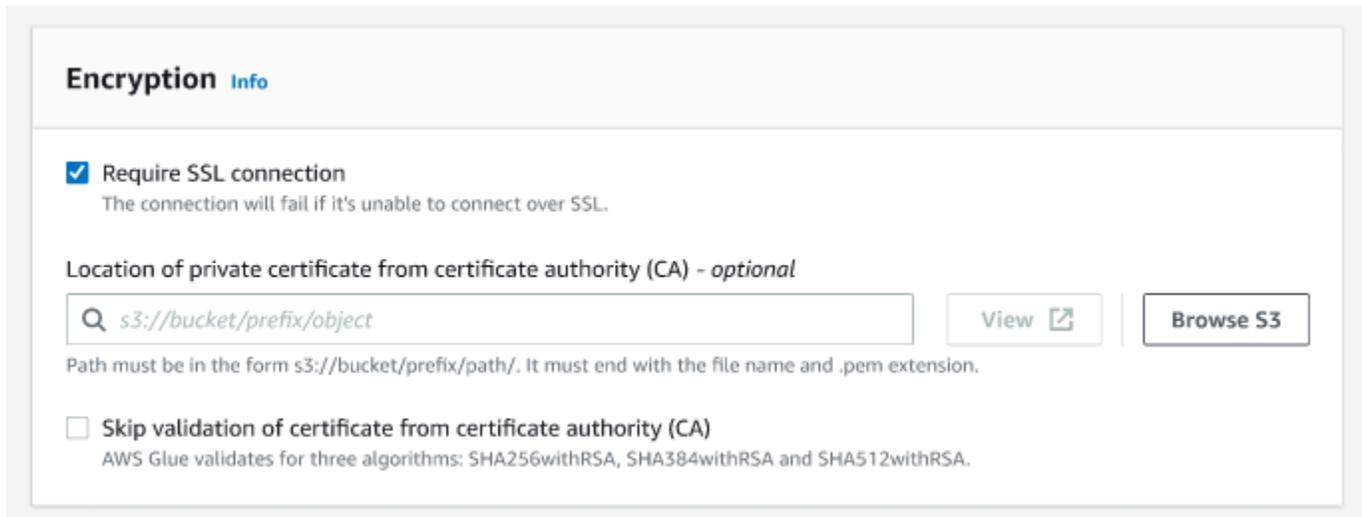
1. Kafka 接続に SSL 接続が必要な場合は、[Require SSL connection] (SSL 接続が必要) チェックボックスを選択します。SSL 経由で接続できない場合、接続は失敗することに留意してください。暗号化用の SSL は、任意の認証方法 (SASL/SCRAM-SHA-512、SASL/GSSAPI、SASL/PLAIN、SSL クライアント認証) と使用できるオプションです。

認証方式が [SSL client authentication] (SSL クライアント認証) に設定されている場合、このオプションは自動的に選択され、変更を防ぐために無効化されます。

2. (オプション)。認証機関 (CA) からプライベート証明書の場所を選択します。証明書の場所は S3 が存在する場所の中である必要があることに留意してください。[Browse] (参照) を選択して、接続された S3 バケットからファイルを選択します。パスは s3://bucket/prefix/filename.pem の形式で指定する必要があります。ファイル名と .pem 拡張子で終わる必要があります。
3. 認証機関 (CA) からの証明書の検証はスキップできます。[Skip validation of certificate from certificate authority (CA)] (認証機関 (CA) からの証明書の検証をスキップする) チェックボックス

を選択してください。このチェックボックスにチェックが入っていない場合、AWS Glue は、3つのアルゴリズムの証明書を検証します。

- SHA256withRSA
- SHA384withRSA
- SHA512withRSA



Encryption [Info](#)

Require SSL connection
The connection will fail if it's unable to connect over SSL.

Location of private certificate from certificate authority (CA) - *optional*

[View](#) [Browse S3](#)

Path must be in the form s3://bucket/prefix/path/. It must end with the file name and .pem extension.

Skip validation of certificate from certificate authority (CA)
AWS Glue validates for three algorithms: SHA256withRSA, SHA384withRSA and SHA512withRSA.

(オプション) ネットワークオプション

以下は、VPC、サブネット、およびセキュリティグループを設定するためのオプションのステップです。AWS Glue ジョブを仮想プライベートクラウド (VPC) サブネット内の Amazon EC2 インスタンスで実行する必要がある場合は、追加の VPC に固有の設定情報を提供する必要があります。

1. データストアが含まれる VPC (仮想プライベートクラウド) を選択します。
2. VPC を持つサブネットを選択します。
3. VPC サブネット内のデータストアへのアクセスを許可するセキュリティグループを 1 つ以上選択します。セキュリティグループは、サブネットに接続されている ENI に関連付けられています。すべての TCP ポートに対して、自己参照のインバウンドルールを持つセキュリティグループを少なくとも 1 つ選択する必要があります。

▼ Network options - *optional*

If your AWS Glue job needs to run on [Amazon Elastic Compute Cloud](#) (EC2) instances in a virtual private cloud (VPC) subnet, you must provide additional VPC-specific configuration information.

VPC [Info](#)

Choose the virtual private cloud that contains your data source.

 ▼

Subnet [Info](#)

Choose the subnet within your VPC.

 ▼

Security groups [Info](#)

Choose one or more security groups to allow access to the data store in your VPC subnet. Security groups are associated to the ENI attached to your subnet. You must choose at least one security group with a self-referencing inbound rule for all TCP ports.

 ▼

カスタムコネクタを使用したジョブのオーサリング

AWS Glue Studio では、データソースノードとデータターゲットノードの両方に対し、コネクタと接続を使用できます

トピック

- [データソースに対しコネクタを使用するジョブを作成する](#)
- [コネクタを使用するノードのソースプロパティを設定する](#)
- [コネクタを使用するノードのターゲットプロパティの設定](#)

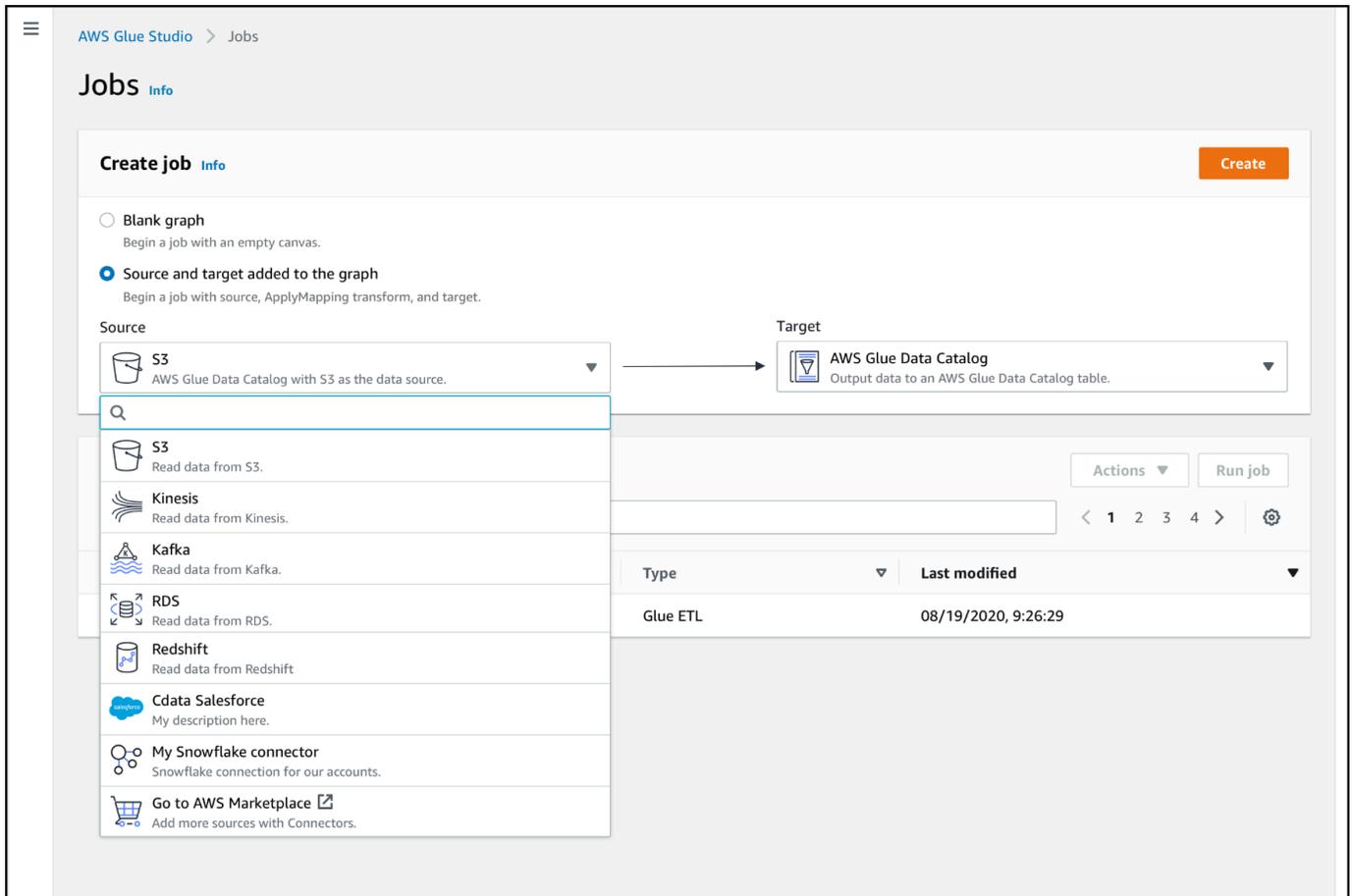
データソースに対しコネクタを使用するジョブを作成する

新しいジョブの作成時に、データソースとデータターゲットで使用するコネクタを選択できます。

データソースまたはデータターゲットに対しコネクタを使用するジョブを作成するには

1. AWS Management Console にサインインし、<https://console.aws.amazon.com/gluestudio/> [AWS Glue Studio](#) にあるコンソールを開きます。
2. [Connectors] (コネクタ) ページの [Your connections] (使用する接続) リソースリストで、ジョブで使用する接続を選択してから、[Create job] (ジョブを作成) をクリックします。

または、AWS Glue Studio の [Jobs] (ジョブ) ページの [Create job] (ジョブを作成) で、[Source and target added to the graph] (グラフに追加されたソースとターゲット) を選択します。[Source] (ソース) ドロップダウンリストから、ジョブで使用するカスタムコネクタを選択します。また、[Target] (ターゲット) のためにコネクタを選択することもできます。



3. [Create] (作成) をクリックして、ビジュアルジョブエディタを開きます。
4. 「[コネクタを使用するノードのソースプロパティを設定する](#)」にある説明に従い、データソースノードを設定します。
5. 「[ビジュアル ETL と AWS Glue Studio](#)」を参考に、変換、追加のデータストア、データターゲットを追加しながら、ETL ジョブを作成します。
6. 「[ジョブのプロパティを変更する](#)」のように、ジョブのプロパティを設定して、ジョブの実行環境をカスタマイズします。
7. このジョブを保存して、実行します。

コネクタを使用するノードのソースプロパティを設定する

データソース用にコネクタを使用するジョブを作成すると、ビジュアルジョブエディタには、コネクタ用に設定されたデータソースのノードを含むジョブグラフが表示されます。このノードには、データソースのプロパティを設定する必要があります。

コネクタを使用するデータソースノードのプロパティを設定するには

1. ジョブグラフでコネクタのデータソースノードを選択するか、新しいノードを追加して [Node type] (ノードタイプ) でコネクタを選択します。次に、右側のノードの詳細パネルで、[Data source properties] (データソースのプロパティ) タブを選択します (まだ選択されていない場合)。

The screenshot displays the AWS Glue console interface for a job titled "Combine legislator data". At the top right, there are buttons for "Save" and "Run", with a red notification "Job has not been saved". The main area shows a job graph with several nodes. On the right, the "Node properties" panel is open to the "Data source properties - Connector" tab. This panel includes sections for "Output schema", "Connection Info" (with a dropdown menu showing "MyEsConn"), and "Schema Info" (with an "Add schema" button). Below these is a section for "Connection options".

2. [Data source properties] (データソースのプロパティ) タブで、このジョブに使用する接続を選択します。

各接続タイプに必要な追加情報を入力します。

JDBC

- [Data source input type] (データソースの入力タイプ): テーブル名または SQL クエリのどちらを、データソースとして指定するか選択します。選択した選択に応じて、以下の追加情報を設定する必要があります。
- [Table name] (テーブル名): データソース内のテーブルの名前。データソースでテーブルという用語が使用されていない場合は、カスタムコネクタの使用情報 (を参照 AWS Marketplace) に示されているように、適切なデータ構造の名前を指定します。
- [Filter predicate] (フィルター述語): データソースを読み取る際に使用する条件句。これは、データのサブセットを取得するために使用される、WHERE 句と類似しています。
- [Query code] (クエリコード): データソースから特定のデータセットを取得するために使用する SQL クエリを入力します。基本 SQL クエリの例は以下のとおりです。

```
SELECT column_list FROM  
           table_name WHERE where_clause
```

- [Schema] (スキーマ): AWS Glue Studio がデータソースにアクセスする際には、Data Catalog テーブルからメタデータ情報を取得する代わりに、接続内に格納された情報を使用します。このため、データソースのスキーマメタデータを指定する必要があります。[Add schema] (スキーマを追加) をクリックして、スキーマエディタを開きます。

スキーマエディタの使用方法については、「[カスタム変換ノードでスキーマを編集する](#)」を参照してください。

- [Partition column] (パーティション列): (オプション) データの読み取りをパーティション化するには、[Partition column] (パーティション列)、[Lower bound] (下限)、[Upper bound] (上限)、および [Number of partitions] (パーティションの数) を、それぞれ指定します。

lowerBound および upperBound 値は、パーティションのストライドを決定するために使用されます (テーブル内の行のフィルタリングには使用しません)。返されるテーブル内のすべての行は、パーティション化されています。

Note

列のパーティショニングは、データの読み取りに使用されるクエリに対し、さらにパーティショニング条件を追加します。テーブル名の代わりにクエリを使用する場

合は、指定されたパーティショニング条件でクエリが動作することを確認する必要があります。例:

- "SELECT col1 FROM table1" の形式のクエリでパーティション列を使用する場合、末尾に WHERE 句を追加してそのクエリをテストします。
- クエリを "SELECT col1 FROM table1 WHERE col2=val" 形式で行っている場合、AND とパーティション列を使用する式を使用して、WHERE 句を拡張することでそのクエリをテストします。

- データ型のキャスト: データソースが JDBC で使用できないデータ型を使用している場合、このセクションを使用して、データソースのデータ型を JDBC データ型に変換する方法を指定します。データ型の変換方法には、最大 50 種類までを指定できます。データソース内で同じデータ型を使用しているすべての列は、同じ方法で変換されます。

例えば、データソース内に Float データ型を使用する 3 つの列があり、Float データ型に対し JDBC の String データ型への変換を指定している場合には、Float データ型を使用する 3 つの列がすべて String データ型に変換されます。

- ジョブのブックマークキー: ジョブのブックマークは、AWS Glue が状態情報を保持することと、古いデータの再処理を防ぐことを助けます。ブックマークキーとして 1 つ以上の列を指定します。AWS Glue Studio はブックマークキーを使用して ETL ジョブの以前の実行中にすでに処理されたデータを追跡します。カスタムブックマークキーに使用する列は、厳密に単調に増加または減少する必要がありますが、そこにギャップを含むことは可能です。

ブックマークキーを複数入力した場合は、それらは結合され単一の複合キーを形成します。複合ジョブブックマークキーには、重複する列を含めることはできません。ブックマークキーを指定しない場合、AWS Glue Studio はデフォルトでプライマリキーをブックマークキーとして使用します。ただし、そのプライマリキーが連続して (ギャップなく) 増減することが条件です。ジョブブックマークのプロパティでは有効になっているものの、テーブルにプライマリキーがない場合には、カスタムジョブブックマークキーを指定する必要があります。この指定を行わないと、デフォルトとして使用するプライマリキーの検索が失敗し、ジョブの実行も失敗します。

- [Job bookmark keys sorting order] (Job ブックマークキーの並べ替え順序): キー値を、連続的に増加させるか減少させるかを選択します。

Spark

- [Schema] (スキーマ): AWS Glue Studio がデータソースにアクセスする際には、Data Catalog テーブルからメタデータ情報を取得する代わりに、接続内に格納された情報を使用します。このため、データソースのスキーマメタデータを指定する必要があります。[Add schema] (スキーマを追加) をクリックして、スキーマエディタを開きます。

スキーマエディタの使用方法については、「[カスタム変換ノードでスキーマを編集する](#)」を参照してください。

- [Connection options:] (接続オプション): 追加の接続情報やオプションを提供するために、必要に応じて追加のキーと値のペアを入力します。例えば、データベース名、テーブル名、ユーザー名、パスワードを入力します。

たとえば、で説明されているように OpenSearch、次のキーと値のペアを入力します。[the section called “チュートリアル: Elasticsearch 向けの AWS Glue コネクタを使用する”](#)

- es.net.http.auth.user : *username*
- es.net.http.auth.pass : *password*
- es.nodes : https://<*Elasticsearch endpoint*>
- es.port : 443
- path: <*Elasticsearch resource*>
- es.nodes.wan.only : true

最低限使用する接続オプションの例については、[MinimalSparkConnectorTest.scala](#) on のサンプルテストスクリプトを参照してください。このサンプルには GitHub、接続で通常指定する接続オプションが示されています。

Athena

- [Table name] (テーブル名): データソース内のテーブルの名前。Athena-CloudWatch ログからの読み取りにコネクタを使用している場合は、all_log_streams テーブル名を入力します。
- [Athena schema name:] (Athena のスキーマ名): テーブルを含むデータベースに対応する、Athena データソース内のスキーマを選択します。Athena-CloudWatch ログからの読み取りにコネクタを使用している場合は、/aws/glue/*name* と同様のスキーマ名を入力します。

- [Schema] (スキーマ): AWS Glue Studio がデータソースにアクセスする際には、Data Catalog テーブルからメタデータ情報を取得する代わりに、接続内に格納された情報を使用します。このため、データソースのスキーマメタデータを指定する必要があります。[Add schema] (スキーマを追加) をクリックして、スキーマエディタを開きます。

スキーマエディタの使用方法については、「[カスタム変換ノードでスキーマを編集する](#)」を参照してください。

- [Additional connection options] (その他の接続オプション): 他の接続情報やオプションを提供する際には、必要に応じて追加のキーと値のペアを入力します。

例については、<https://github.com/aws-samples/aws-glue-samples/tree/master/GlueCustomConnectors/Development/Athena README.md>にあるファイルを参照してください。このドキュメント内に示す手順では、サンプルコードにより、必要最小限の接続オプション (tableName、schemaName、および className) を使用しています。コード例では、これらのオプションを optionsMap 変数の一部として指定しており、それらを実際に指定することで接続を使用できるようになります。

3. (オプション) 必要な情報を指定した後は、[Output schema] (出力スキーマ) タブを選択することで、出力されたデータソース用のデータスキーマを、ノードの詳細パネルに表示できるようになります。このタブに表示されるスキーマは、ジョブグラフに追加される任意の子ノードによって使用されます。
4. (オプション) ノードおよびデータソースのプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、データソースからのデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

コネクタを使用するノードのターゲットプロパティの設定

データターゲットタイプにコネクタを使用する場合は、データターゲットノードのプロパティを設定する必要があります。

コネクタを使用するデータターゲットノードでプロパティを設定するには

1. ジョブグラフで、コネクタのデータターゲットノードを選択します。次に、右側のノードの詳細パネルで、[Data target properties] (データターゲットのプロパティ) タブを選択します (選択されていない場合)。

2. [Data target properties] (データターゲットのプロパティ) タブで、ターゲットへの書き込みに使用する接続を選択します。

各接続タイプに必要な追加情報を入力します。

JDBC

- [Connection] (接続): コネクタで使用する接続を選択します。接続の作成方法については、[「コネクタ用の接続を作成する」](#)を参照してください。
- [Table name] (テーブル名): データターゲット内のテーブルの名前。データターゲットがテーブルという用語を使用していない場合は、カスタムコネクタの使用情報 (を参照) に示されているように、適切なデータ構造の名前を指定します。AWS Marketplace
- [Batch size] (バッチサイズ) (オプション): 1 回のオペレーションでターゲットテーブルに挿入する、行数またはレコード数を入力します。デフォルト値は 1000 行です。

Spark

- [Connection] (接続): コネクタで使用する接続を選択します。以前に接続を作成していない場合は、[Create connection] (接続を作成する) をクリックして作成します。接続の作成方法については、[「コネクタ用の接続を作成する」](#)を参照してください。
- [Connection options:] (接続オプション): 追加の接続情報やオプションを提供するために、必要に応じて追加のキーと値のペアを入力します。データベース名、テーブル名、ユーザー名、パスワードを入力することもできます。

たとえば、で説明されているように OpenSearch、次のキーと値のペアを入力します。[the section called “チュートリアル: Elasticsearch 向けの AWS Glue コネクタを使用する”](#)

- es.net.http.auth.user : *username*
- es.net.http.auth.pass : *password*
- es.nodes : `https://<Elasticsearch endpoint>`
- es.port : 443
- path: `<Elasticsearch resource>`
- es.nodes.wan.only : true

最低限使用する接続オプションの例については、[MinimalSparkConnectorTest.scala](#) on のサンプルテストスクリプトを参照してください。このサンプルには GitHub、接続で通常指定する接続オプションが示されています。

3. 必要な情報を指定した後は、[Output schema] (出力スキーマ) タブを選択することで、出力されたデータソース用のデータスキーマを、ノードの詳細パネルに表示できるようになります。

コネクタと接続を管理する

コネクタと接続を管理するには、AWS Glue の [コネクタ] ページを使用します。

トピック

- [コネクタと接続の詳細を表示する](#)
- [コネクタと接続を編集する](#)
- [コネクタおよび接続を削除する](#)
- [コネクタのサブスクリプションをキャンセルする](#)

コネクタと接続の詳細を表示する

[Connectors] (コネクタ) ページの [Your connectors] (使用中のコネクタ)、および [Your connections] (使用中の接続) リソーステーブルで、コネクタと接続に関する概要情報を表示できます。詳細な情報を表示するには、以下の手順を実行します。

コネクタまたは接続の詳細を表示するには

1. AWS Glue Studio コンソールのナビゲーションペインで、[Connectors] (コネクタ) をクリックします。
2. 詳細情報を表示するコネクタまたは接続を選択します。
3. [Actions] (アクション)、[View details] (詳細を表示) の順にクリックして、選択したコネクタまたは接続の詳細ページを開きます。
4. 詳細ページでは、コネクタまたは接続に対し、[Edit] (編集) または [Delete] (削除) を選択できます。
 - コネクタの場合は、[Create connection] (接続を作成する) をクリックすることで、そのコネクタを使用する新しい接続を作成できます。
 - 接続の場合は、[Create job] (ジョブを作成する) をクリックすると、その接続を使用するジョブを作成できます。

コネクタと接続を編集する

[Connectors] (コネクタ) ページを使用して、コネクタと接続に保存されている情報を変更します。

コネクタまたは接続を変更するには

1. AWS Glue Studio コンソールのナビゲーションペインで、[Connectors] (コネクタ) をクリックします。
2. 変更するコネクタまたは接続を選択します。
3. [Actions] (アクション) を選択して、[Edit] (編集) を選択します。

また、[View details] (詳細を表示する) をクリックし、コネクタまたは接続の詳細ページを開いて、[Edit] (編集) を実行することもできます。

4. [Edit connector] (コネクタを編集する) または [Edit connection] (接続を編集する) ページで、情報を更新した上で、[Save] (保存) をクリックします。

コネクタおよび接続を削除する

コネクタと接続を削除するには、[Connectors] (コネクタ) ページを使用します。コネクタを削除すると、そのコネクタのために作成された接続もすべて削除されます。

AWS Glue Studio からコネクタを削除する方法

1. AWS Glue Studio コンソールのナビゲーションペインで、[Connectors] (コネクタ) をクリックします。
2. 削除するコネクタまたは接続を選択します。
3. [Actions] (アクション) を選択してから、[Delete] (削除) をクリックします。

または、[View details] (詳細を表示する) をクリックし、コネクタまたは接続の詳細ページで、[Delete] (削除) を実行することもできます。

4. 「Delete」と入力し、コネクタまたは接続を削除することを確認した上で、[Delete] (削除) をクリックします。

コネクタを削除すると、そのコネクタのために作成された接続もすべて削除されます。

ここで削除される接続を使用していたジョブも、すべて機能しなくなります。これらのジョブは、別のデータストアを使用するように編集することも、あるいは削除することもできます。ジョブを削除する方法については、「[ジョブの削除](#)」を参照してください。

コネクタを削除しても、そのコネクタの AWS Marketplace に対するサブスクリプションはキャンセルされません。削除されたコネクタのサブスクリプションを削除するには、「[コネクタのサブスクリプションをキャンセルする](#)」の手順に従います。

コネクタのサブスクリプションをキャンセルする

AWS Glue Studio 接続とコネクタをから削除した後で、AWS Marketplace コネクタが不要になった場合はサブスクリプションをキャンセルできます。

Note

コネクタのサブスクリプションをキャンセルしても、コネクタや接続はアカウントから削除されません。コネクタと、それに関連する接続を使用するジョブは、以後コネクタを使用できなくなり失敗します。

からコネクタの購読を解除または再購読する前に AWS Marketplace、AWS Marketplace その製品に関連付けられている既存の接続とコネクタを削除する必要があります。

でコネクタの購読を解除するには AWS Marketplace

1. <https://console.aws.amazon.com/marketplace> [AWS Marketplace](#) でコンソールにサインインします。
2. [Manage subscriptions] (サブスクリプションを管理する) をクリックします。
3. [Manage subscriptions] (サブスクリプションの管理) ページで、キャンセルするコネクタのサブスクリプションの横にある [Manage] (管理) をクリックします。
4. [Actions] (アクション)、[Delete Application] (アプリケーションの削除) の順にクリックします。
5. 実行中のインスタンスがアカウントで課金されることを受け入れるチェックボックスをオンにして、[Yes, cancel subscription] (はい、サブスクリプションをキャンセルします) をクリックします。

カスタムコネクタの開発

コードを記述すると、データストアとの間でデータの読み取りや書き込みを実行したり、AWS Glue Studio ジョブで使用するためにデータのフォーマットを行ったりできます。Spark、Athena、JDBC データストアのためのコネクタを作成できます。に掲載されているサンプルコードには、GitHub 実装する必要のある基本インターフェースの概要が記載されています。

コネクタコードを作成するには、ローカルの開発環境が必要です。コネクタを書き込むには、任意の IDE やコマンドラインエディタを使用することもできます。開発環境の例には以下があります。

- ローカルの AWS Glue ETL Maven ライブラリを使用したローカル Scala 環境 (AWS Glue デベロッパーガイドの「[Developing Locally with Scala](#)」を参照)。
- <https://www.jetbrains.com/idea/> からダウンロードして使用する、IntelliJ IDE。

トピック

- [Spark コネクタの開発](#)
- [Athena コネクタの開発](#)
- [JDBC コネクタの開発](#)
- [AWS Glue Studio でのカスタムコネクタ使用例](#)
- [以下のコネクタの開発 AWS Glue AWS Marketplace](#)

Spark コネクタの開発

Spark DataSource API V2 (Spark 2.4) を使用して Spark コネクタを作成して、データを読み取ることができます。

カスタム Spark コネクタを作成するには

<https://github.com/aws-samples/tree/master/development/spark/README.md> にある [Spark AWS Glue GitHub コネクタを開発するためのサンプルライブラリの手順に従ってください](#)。aws-glue-samples GlueCustomConnectors

Athena コネクタの開発

Athena コネクタを作成し、AWS Glue および AWS Glue Studio でカスタムのデータソースをクエリするために使用できます。

カスタムの Athena コネクタを作成するには

<https://github.com/aws-samples/aws-glue-samples/tree/master/GlueCustomConnectors/development/Athena> にある [Athena AWS Glue GitHub コネクタを開発するためのサンプルライブラリの手順に従ってください](#)。

JDBC コネクタの開発

データストアにアクセスするために、JDBC を使用するコネクタを作成できます。

カスタムの JDBC コネクタを作成するには

1. ローカルの開発環境に、AWS Glue Spark ランタイムライブラリをインストールします。[https://github.com/aws-samples/ /tree/master/ /development/ /README.md](https://github.com/aws-samples/tree/master/development/README.md)にあるサンプルライブラリの説明を参照してください。AWS Glue GitHub [aws-glue-samples](#) [GlueCustomConnectors](#) [GlueSparkRuntime](#)
2. データソースからデータを取得するための JDBC ドライバーを実装します。Java SE 8 向けの [Java ドキュメント](#)を参照してください。

作成しているコード内に、AWS Glue Studio がコネクタの位置を特定するために使用するエントリポイントを記述します。[Class name] (クラス名) フィールドは、JDBC ドライバーへの完全なパスを指定する必要があります。

3. GlueContext API を使用して、コネクタによりデータを読み取ります。必要に応じてユーザーは、AWS Glue Studio コンソールを使用して、データソースへの接続を設定するための他の入力オプションを追加できます。カスタム JDBC コネクタを使用して JDBC データベースから読み書きする方法を示すコード例については、「[カスタム値と connectionType 値](#)」を参照してください。AWS Marketplace

AWS Glue Studio でのカスタムコネクタ使用例

カスタムコネクタの使用例については、次のブログを参照してください。

- [AWS Glue を使用したデータストア用のカスタムコネクタの開発、テスト、およびデプロイ](#)
- Apache Hudi: [Writing to Apache Hudi tables using AWS Glue Custom Connector](#)
- グーグル BigQuery: [AWS Glue カスタムコネクタを使用してグーグルから Amazon S3 BigQuery にデータを移行する](#)
- Snowflake (JDBC): [Performing data transformations using Snowflake and AWS Glue](#)
- SingleStore: [とを使用した高速 ETL の構築 SingleStore AWS Glue](#)
- Salesforce: [Ingest Salesforce data into Amazon S3 using the CData JDBC custom connector with AWS Glue](#)
- MongoDB: [Building AWS Glue Spark ETL jobs using Amazon DocumentDB \(with MongoDB compatibility\) and MongoDB](#)
- Amazon Relational Database Service (Amazon RDS): Amazon RDS [用の独自の JDBC ドライバーを持ち込んで AWS Glue Spark ETL ジョブを構築する](#)
- MySQL (JDBC): [https://github.com/aws-samples/ aws-glue-samples /blob/master/ /Development/ Spark/ SQL.scala](https://github.com/aws-samples/aws-glue-samples/blob/master/Development/Spark/SQL.scala) [GlueCustomConnectors](#) [SparkConnectorMy](#)

以下のコネクタの開発 AWS GlueAWS Marketplace

AWS パートナーは、カスタムコネクタを作成してアップロードし、AWS Marketplace AWS Glue 顧客に販売することができます。

コネクタコードの開発プロセスは、カスタムコネクタの場合と同様です。ただし、コネクタコードのアップロードと検証のプロセスには、より詳細な要素が含まれます。GitHub Web サイトの「[コネクタの作成](#)」 [AWS Marketplaceの手順を参照してください](#)。

AWS Glue Studio でのコネクタおよび接続の使用に関する制約事項

カスタムコネクタまたはからのコネクタを使用する場合は AWS Marketplace、以下の制限に注意してください。

- カスタムコネクタ用に作成された接続では、TestConnection API は使用できません。
- Data Catalog 接続でのパスワードの暗号化は、カスタムコネクタではサポートされていません。
- JDBC コネクタを使用するデータソースノードでフィルター述語を指定した場合は、ジョブブックマークを使用できなくなります。
- Marketplace 接続の作成は、AWS Glue Studio ユーザーインターフェイス以外ではサポートされていません。

Visual ETL ジョブを使用してデータソースに接続する

新しいジョブの作成中、AWS Glue でのビジュアル ETL ジョブ作成時に接続を使用して、データに接続できます。これには、コネクタを使用してデータを読み込むソースノードと、データを書き出す場所を指定するターゲットノードを追加します。

トピック

- [データソースノードのプロパティを変更する](#)
- [データソースにデータカタログテーブルを使用する](#)
- [データソースにコネクタを使用する](#)
- [データソースに Amazon S3 内のファイルを使用する](#)
- [ストリーミングデータソースの使用](#)
- [リファレンス](#)

データソースノードのプロパティを変更する

データソースのプロパティを指定するには、初めにジョブ図でデータソースノードを選択します。次に、ノードの詳細パネルの右側で、ノードのプロパティを設定します。

データソースノードのプロパティを変更するには

1. 新規または保存済みのジョブのビジュアルエディタに移動します。
2. ジョブ図でデータソースノードを選択します。
3. ノードの詳細パネルの [Node properties] (ノードのプロパティ) タブを選択して、次の情報を入力します。
 - Name (名前): (オプション) ジョブ図でノードに関連付ける名前を入力します。この名前は、このジョブのすべてのノードにおいて一意である必要があります。
 - Node type (ノードタイプ): ノードタイプによって、ノードで実行されるアクションが決定します。[Node type] (ノードタイプ) のオプションのリストで、見出し [Data source] (データソース) の下にリストされている値のいずれかを選択します。
4. データソースのプロパティの情報を設定します。詳細については、次のセクションを参照してください。
 - [データソースにデータカタログテーブルを使用する](#)
 - [データソースにコネクタを使用する](#)
 - [データソースに Amazon S3 内のファイルを使用する](#)
 - [ストリーミングデータソースの使用](#)
5. (オプション) ノードおよびデータソースのプロパティを設定した後、ノードのパネルの [Output schema] (出力スキーマ) タブをクリックすると、データソースのスキーマを表示できます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。
6. (オプション) ノードおよびデータソースのプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、データソースからのデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

データソースにデータカタログテーブルを使用する

Amazon S3 とコネクタを除くすべてのデータソースでは、選択するソースタイプの AWS Glue Data Catalog にテーブルが存在する必要があります。AWS Glue はデータカタログテーブルを生成しません。

データカタログテーブルに基づいてデータソースノードを設定するには

1. 新規または保存済みのジョブのビジュアルエディタに移動します。
2. ジョブ図でデータソースノードを選択します。
3. [Data source properties] (データソースのプロパティ) タブを選択して、次の情報を入力します。
 - S3 source type (S3 ソースタイプ): (Amazon S3 データソースのみ) [Select a Catalog table] (Catalog テーブルを選択) オプションを選択して、既存の AWS Glue Data Catalog テーブルを使用します。
 - Database (データベース): このジョブに使用するソーステーブルを含むデータカタログのデータベースを選択します。検索フィールドを使用して、名前でデータベースを検索できます。
 - Table (テーブル): ソースデータに関連付けられたテーブルをリストから選択します。このテーブルは、既に AWS Glue Data Catalog に存在している必要があります。検索フィールドを使用して、名前でテーブルを検索できます。
 - Partition predicate (パーティション述語): (Amazon S3 データソースのみ) パーティション列のみを含む Spark SQL に基づいてブール式を入力します。例: "(year=='2020' and month=='04')"
 - Temporary directory (一時ディレクトリ): (Amazon Redshift データソースのみ) ETL ジョブが一時的な中間結果を書き込める Amazon S3 の作業ディレクトリの場所のパスを入力します。
 - Role associated with the cluster (クラスターに関連付けられたロール): (Amazon Redshift データソースのみ) Amazon Redshift クラスターのアクセス許可を含む、使用する ETL ジョブのロールを入力します。詳細については、「[the section called “データソースとデータターゲットのアクセス許可”](#)」を参照してください。

データソースにコネクタを使用する

ノードタイプにコネクタを選択する場合、[カスタムコネクタを使用したジョブのオーサリング](#)の説明に従って、データソースのプロパティの設定を完了します。

データソースに Amazon S3 内のファイルを使用する

データソースとして Amazon S3 を選択する場合、次のいずれかを選択できます。

- データカタログデータベースおよびテーブル。
- Amazon S3 内のバケット、フォルダ、またはファイル。

Amazon S3 バケットをデータソースとして使用する場合、AWS Glue は、ファイルの一つ、あるいはサンプルファイルとして指定したファイルを使うことによって、指定した場所にあるデータのスキーマを検出します。[Infer schema] (スキーマを推測) ボタンをクリックすると、スキーマの検出が行われます。Amazon S3 の場所またはサンプルファイルを変更する場合は、[Infer schema] (スキーマを推測) を再度クリックして、新しい情報を使用しスキーマの検出を実行します。

Amazon S3 内のファイルから直接読み取るデータソースノードを設定するには

1. 新規または保存済みのジョブのビジュアルエディタに移動します。
2. Amazon S3 ソース用に、ジョブ図でデータソースノードを選択します。
3. [Data source properties] (データソースのプロパティ) タブを選択して、次の情報を入力します。
 - S3 source type (S3 ソースタイプ): (Amazon S3 データソースのみ) [S3 location] (S3 の場所) オプションを選択します。
 - S3 URL (S3 の URL): ジョブのデータを含む Amazon S3 バケット、フォルダ、またはファイルへのパスを入力します。[Browse S3] (S3 をブラウズ) を選択すると、アカウントで利用可能な場所からパスを選択できます。
 - Recursive : S3 の場所にある子フォルダ内のファイルからデータを AWS Glue に読ませたい場合、このオプションを選択します。

子フォルダにパーティション化されたデータが含まれている場合、AWS Glue は、フォルダ名中で指定されているパーティション情報をデータカタログに追加しません。例えば、Amazon S3 で、次のフォルダについて考えてみます。

```
S3://sales/year=2019/month=Jan/day=1
S3://sales/year=2019/month=Jan/day=2
```

[Recursive] を選択し、sales フォルダを S3 の場所として選択した場合、AWS Glue は、すべての子フォルダーのデータを読みますが、年、月、日のパーティションは作成しません。

- **Data format (データ形式):** データを保存する形式を選択します。JSON、CSV、または Parquet を選択できます。選択した値により、ソースファイルからデータを読み取る方法が AWS Glue ジョブに伝えられます。

 Note

データのために正しい形式を選択していない場合、AWS Glue はスキーマは正しく推測する可能性はありますが、ジョブはソースファイルからのデータを正しく解析することはできないかもしれません。

選択した形式に応じて、追加の設定オプションを入力できます。

- **JSON (JavaScript Object Notation)**
 - **JsonPath:** テーブルスキーマの定義に使用されるオブジェクトを指す JSON パス。JSON パス式では、XPath 式が XML ドキュメントと組み合わせて使用されるのと同じように、常に JSON 構造を参照します。JSON パスの「ルートメンバーオブジェクト」は、オブジェクトや配列であっても、常に \$ として参照されます。JSON パスは、ドット表記またはブラケット表記で記述できます。

JSON パスの詳細については、GitHub ウェブサイトの「[JsonPath](#)」を参照してください。

- **Records in source files can span multiple lines (ソースファイル内のレコードが複数行にまたがることのできる):** CSV ファイル内の単一のレコードが複数行にまたがることのできる場合は、このオプションを選択します。
- **CSV (カンマ区切り値)**
 - **Delimiter (区切り記号):** 行内の各列エントリの区切りを示す文字 (;、, など) を入力します。
 - **Escape character (エスケープ文字):** エスケープ文字として使用する文字を入力します。この文字は、その直後の文字は文字通りに解釈され、区切り文字として解釈しないことを示します。
 - **Quote character (引用符文字):** 区切られた文字列を単一の値にグループ化するのに使用する文字を入力します。例えば、CSV ファイルに "This is a single value" などの値がある場合、ダブルクォート (") を選択します。

- Records in source files can span multiple lines (ソースファイル内のレコードが複数行にまたがることのできる): CSV ファイル内の単一のレコードが複数行にまたがることのできる場合は、このオプションを選択します。
- First line of source file contains column headers (ソースファイルの最初の行に列ヘッダーが含まれる): CSV ファイルの最初の行にデータではなく列ヘッダーが含まれている場合は、このオプションを選択します。
- Parquet (Apache Parquet 列指向ストレージ)

Parquet 形式で保存されたデータに対する追加の設定はありません。

- Partition predicate (パーティション述語): データソースから読み取られたデータをパーティション分割するには、パーティション列のみを含む Spark SQL に基づいてブール式を入力します。例: "(year=='2020' and month=='04')"
- Advanced options : 特定のファイルに基づくデータのスキーマの検出を AWS Glue にさせたい場合は、このセクションを展開します。
- Schema inference : AWS Glue にファイルを選択させる代わりに、特定のファイルを使用したい場合は、オプションの「Choose a sample file from S3」を選択します。
- Auto-sampled file (自動サンプリングファイル): スキーマの推測に使用する Amazon S3 内のファイルへのパスを入力します。

データソースノードを編集集中に、選択したサンプルファイルを変更する場合は、[Reload schema] (スキーマの再ロード) を選択して、新しいサンプルファイルを使用してスキーマを検出します。

4. [Infer schema] (スキーマを推測) ボタンをクリックして、Amazon S3 内のソースファイルからスキーマを検出します。Amazon S3 の場所またはサンプルファイルを変更する場合は、[Infer schema] (スキーマを推測) を再度選択して、新しい情報を使用してスキーマを推測します。

ストリーミングデータソースの使用

Amazon Kinesis Data Streams、Apache Kafka、および Amazon Managed Streaming for Apache Kafka (Amazon MSK) で、ストリーミングからのデータを消費し、連続的に実行されるストリーミング抽出、変換、ロード (ETL) ジョブを作成できます。

ストリーミングデータソースのプロパティを設定するには

1. 新規または保存済みのジョブのビジュアルグラフエディタに移動します。
2. Kafka または Kinesis Data Streams のグラフでデータソースノードを選択します。

3. [Data source properties] (データソースのプロパティ) タブを選択して、次の情報を入力します。

Kinesis

- Kinesis source type: ストリームの詳細への直接のアクセスを使うためにオプションの [Stream details] を選択するか、そこに格納されている情報を使用するために [Data Catalog table] (データカタログテーブル) を選択します。

[Stream details] を選択すると、以下の追加情報を指定します。

- Location of data stream : ストリームを現在のユーザーに関連付けるかどうか、または別のユーザーに関連付けるかどうかを選択します。
- Region: ストリームが存在する [AWS リージョン] を選択します。この情報は、データストリームにアクセスするための ARN を構築するために使用されます。
- Stream ARN: Kinesis Data Streams の Amazon リソースネーム (ARN) を入力します。ストリームが現在のアカウント内にある場合は、ドロップダウンリストからストリーム名を選択します。検索フィールドを使用して、名前、あるいは ARN でデータベースを検索できます。
- Data format: リストから、データストリームで使用される形式を選択します。

AWS Glue は、ストリーミングデータからスキーマを自動的に検出します。

[Data Catalog table] (データカタログテーブル) を選択すると、以下の追加情報を指定します。

- Database : (オプション) ストリーミングデータソースに関連付けられたテーブルを含む AWS Glue データカタログ内のデータベースを選択します。検索フィールドを使用して、名前でデータベースを検索できます。
- Table (テーブル): (オプション) ソースデータに関連付けられたテーブルをリストから選択します。このテーブルは、AWS Glue のデータカタログに既に存在している必要があります。検索フィールドを使用して、名前でテーブルを検索できます。
- Detect schema: このオプションを選択すると、スキーマの情報がデータカタログテーブルに格納されるのではなく、AWS Glue にストリーミングデータからスキーマを検出させます。[Stream details] オプションを選択すると、このオプションは自動的に有効になります。
- Starting position: デフォルトでは、ETL ジョブは [Earliest] オプションを使います。これは、ストリーム内の使用可能な最も古いレコードから始まるデータを読み取ることを意味します。代わりに [Latest] を選択することができます。これは、ETL ジョブがストリーム内の最新のレコードの直後から読み出しを開始する必要があることを示します。

- Window size (ウィンドウサイズ): デフォルトでは、ETL ジョブによるデータ処理と書き出しは 100 秒ウィンドウ単位で行われます。これにより、データを効率的に処理しつつ、想定より遅く到着するデータに対して集計を実行できます。ウィンドウサイズを変更することで、適時性や集計の精度を高めることができます。

AWS Glue ストリーミングジョブでは、読み込んだデータの追跡には、ジョブのブックマークではなくチェックポイントが使用されます。

- Connection options : 追加の接続オプションを指定するために、このセクションを展開してキーバリュペアを追加します。ここで指定できるオプションの詳細については、[AWS Glue Developer Guide]中の「"[connectionType": "kinesis"](#)」を参照してください。

Kafka

- Apache Kafka source: オプションの[Stream details]を選択し、ストリームの詳細ストリーミングソースへの直接アクセスを使用するか、[Data Catalog table] を選択して代わりに、そこに格納されている情報を使用します。

[Data Catalog table] を選択すると、以下の追加情報を指定します。

- Database : (オプション) ストリーミングデータソースに関連付けられたテーブルを含む AWS Glue データカタログ内のデータベースを選択します。検索フィールドを使用して、名前データベースを検索できます。
- Table (テーブル): (オプション) ソースデータに関連付けられたテーブルをリストから選択します。このテーブルは、AWS Glue のデータカタログに既に存在している必要があります。検索フィールドを使用して、名前テーブルを検索できます。
- Detect schema : このオプションを選択し、スキーマの情報がデータカタログテーブルに格納されるのではなく、ストリーミングデータからスキーマを AWS Glue に検出させます。Stream details オプションを選択すると、このオプションは自動的に有効になります。

Stream details を選択すると、以下の追加情報を指定します。

- Connection name: Kafka データストリームのためのアクセス情報と認証情報を含む AWS Glue の接続を選択します。Kafka ストリーミングデータソースとの接続を使用する必要があります。接続が存在しない場合は、AWS Glue のコンソールを使用して Kafka データストリームのための接続を作成します。
- Topic name: 読み取り元のトピック名を入力します。

- **Data format** : Kafka イベントストリームからデータを読み取る際に使用する形式を選択します。
- **Starting position**: デフォルトでは、ETL ジョブは[Earliest] オプションを使用します。これは、ストリーム内の使用可能な最も古いレコードから始まるデータを読み取ることを意味します。代わりに[Latest]選択することができます。これは、ETL ジョブがストリーム内の最新のレコードの直後から読み出しを開始する必要があることを示します。
- **Window size (ウィンドウサイズ)**: デフォルトでは、ETL ジョブによるデータ処理と書き出しは 100 秒ウィンドウ単位で行われます。これにより、データを効率的に処理しつつ、想定より遅く到着するデータに対して集計を実行できます。ウィンドウサイズを変更することで、適時性や集計の精度を高めることができます。

AWS Glue ストリーミングジョブでは、読み込んだデータの追跡には、ジョブのブックマークではなくチェックポイントが使用されます。

- **Connection options** : 追加の接続オプションを指定するために、このセクションを展開してキーバリュペアを追加します。ここで指定できるオプションの詳細については、[AWS Glue Developer Guide]内の「["connectionType": "kafka"](#)」を参照してください。

Note

現在、ストリーミングデータソースでのデータのプレビューはサポートされていません。

リファレンス

ベストプラクティス

- [AWS Glue を使用して Amazon S3 から Amazon Redshift にデータを段階的にロードする ETL サービスパイプラインを構築します。](#)

ETL プログラミング

- [AWS Glue での ETL の接続タイプとオプション](#)
- [JDBC connectionType の値](#)
- [Amazon Redshift との間でのデータ移動に関する詳細オプション](#)

独自の JDBC ドライバーを使用した JDBC 接続の追加

JDBC 接続を使用する際に、独自の JDBC ドライバーを使用できます。AWS Glue クローラーが使用するデフォルトドライバーがデータベースに接続できない場合、独自の JDBC ドライバーを使用できます。例えば、Postgres データベースで SHA-256 を使用したいが、古い Postgres ドライバーがこれをサポートしていない場合は、独自の JDBC ドライバーを使用できます。

サポート対象データソース

サポート対象データソース	サポートされていないデータソース
MySQL	Snowflake
Postgres	
Oracle	
Redshift	
SQL Server	
Aurora*	

* ネイティブの JDBC ドライバーが使用されている場合にサポートされます。すべてのドライバー機能を使用できるわけではありません。

JDBC 接続への JDBC ドライバーの追加

Note

独自の JDBC ドライバーバージョンを導入する場合、AWS Glue クローラーは AWS Glue ジョブと Amazon S3 バケットのリソースを使用して、用意したドライバーが自分の環境で実行されるようにします。リソースの追加使用量はアカウントに反映されます。AWS Glue クローラーとジョブのコストは、請求の AWS Glue カテゴリに含まれます。さらに、独自の JDBC ドライバーを用意しても、クローラーがドライバーの機能をすべて活用できるわけではありません。

独自の JDBC ドライバーを JDBC 接続に追加するには:

1. JDBC ドライバーファイルを Amazon S3 の場所に追加します。バケットおよび/またはフォルダを作成または使用できます。
2. AWS Glue コンソールで、左側のメニューの [データカタログ] の [接続] を選択し、新しい接続を作成します。
3. [接続プロパティ] のフィールドに入力し、[接続タイプ] で JDBC を選択します。
4. [接続アクセス] に、[JDBC URL] と [JDBC ドライバーのクラス名] を入力します (オプション)。ドライバークラス名は、AWS Glue クローラーがサポートするデータソースのものである必要があります。

Connection access

JDBC URL
Use the JDBC protocol to access Amazon Redshift, Amazon RDS, and publicly accessible databases.

JDBC syntax for most database engines is jdbc:protocol://host:port/databasename.

JDBC Driver Class name - optional

Type a custom JDBC driver class name for the crawler to connect to the data source.

JDBC Driver S3 Path - optional

Browse for or enter an existing S3 path to a .jar file.

Please note that if you choose to bring in your own JDBC driver versions to be used with Glue Crawlers, the Glue Crawlers will consume resources in Glue Jobs and S3 to ensure your provided driver are run in your environment. The additional usage of resources will be reflected in your account.

Credential type

Username and password
 Secret

Username

Password

5. JDBC ドライバーが置かれている Amazon S3 パスを [JDBC ドライバーの Amazon S3 パス] で選択します (オプションフィールド)。
6. ユーザー名とパスワード、またはシークレットを入力する場合は、認証情報タイプのフィールドに入力します。完了したら、[接続を作成] を選択します。

 Note

テスト接続は現在サポートされていません。用意した JDBC ドライバーを使用してデータソースをクローリングする場合、クローラーはこのステップをスキップします。

7. 新しく作成した接続をクローラーに追加します。AWS Glue コンソールで、左側のメニューの [データカタログ] で [クローラー] を選択し、新しいクローラーを作成します。
8. [クローラーの追加] ウィザードの [ステップ 2] で [データソースの追加] を選択します。

Add data source ✕

Data source
Choose the source of data to be crawled.

JDBC ▼

Connection
Select a connection to access the data sources below.

mysql-connection068fd134-c2f1-4234-ad6b-345968e73be8 ▼ ↻

Clear selection Add new connection [↗](#)

Include path

public/%

You can substitute the percent (%) character for a schema or table. For databases that support schemas, enter MyDatabase/MySchema/% to match all tables in MySchema within MyDatabase. Oracle Database and MySQL don't support schema in the path; instead, enter MyDatabase/%. For Oracle database without SSL, MyDatabase can be either the system identifier (SID) or the service name (SERVICE_NAME). For Oracle database with SSL, MyDatabase must be the service name (SERVICE_NAME).

Additional metadata - optional

▼

Select additional metadata properties for the crawler to crawl.

Exclude tables matching pattern

Cancel Add a JDBC data source

9. データソースとして [JDBC] を選択し、これまでの手順で作成した接続を選択します。完了
10. クローラーで独自の JDBC ドライバーを使用するには、AWS Glue クローラーが使用するロールに次のアクセス許可を追加します。
 - ジョブアクション CreateJob、DeleteJob、GetJob、GetJobRun、StartJobRun にアクセス権限を付与します。
 - IAM アクション iam:PassRole にアクセス権限を付与します。

- Amazon S3 アクション
s3:DeleteObjects、s3:GetObject、s3:ListBucket、s3:PutObject にアクセス権限を付与します。
- サービスプリンシパルに IAM ポリシーのバケット/フォルダへのアクセス権限を付与します。

IAM ポリシーの例:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "VisualEditor0",
      "Effect": "Allow",
      "Action": [
        "s3:PutObject",
        "s3:GetObject",
        "s3:ListBucket",
        "s3:DeleteObject"
      ],
      "Resource": [
        "arn:aws:s3:::bucket-name/driver-parent-folder/driver.jar",
        "arn:aws:s3:::bucket-name"
      ]
    }
  ]
}
```

11. VPC を使用している場合は、インターフェイスエンドポイントを作成し、これをルートテーブルに追加して AWS Glue へのアクセスを許可する必要があります。詳細については、「[AWS Glue 用のインターフェイス VPC エンドポイントの作成](#)」を参照してください。
12. Data Catalog で暗号化を使用している場合は、AWS KMS インターフェイスエンドポイントを作成し、ルートテーブルに追加します。詳細については、「[AWS KMS用の VPC エンドポイントの作成](#)」を参照してください。

AWS Glue 接続のテスト

ベストプラクティスとして、ETL ジョブで AWS Glue 接続を使用する前に、AWS Glue コンソールを使用して接続をテストします。AWS Glue は、接続内のパラメータを使用してデータストアにアクセスできることを確認し、エラーを報告します。AWS Glue 接続の詳細については、「[データへの接続](#)」を参照してください。

AWS Glue 接続をテストするには

1. AWS Management Console にサインインし、<https://console.aws.amazon.com/glue/> AWS Glue でコンソールを開きます。
2. ナビゲーションペインの [データカタログ] で [接続] を選択します。ナビゲーションウィンドウの [データ カタログ] 上で、[データ接続] を選択することもできます。
3. [接続] で、目的の接続の横にあるチェックボックスをオンにし、[アクション] を選択します。ドロップダウンメニューで、[接続のテスト] を選択します。
4. テスト接続ダイアログボックスでロールを選択するか、[IAM ロールの作成] を選択して AWS Identity and Access Management (IAM) コンソールに移動して新しいロールを作成します。ロールには、データストアに対する権限が必要です。
5. [確認] を選択します。

テストが開始されます。完了までに数分かかることがあります。テストが失敗した場合、[Troubleshoot] をクリックし、問題解決のステップを表示します。

6. [Logs] を選択してログインしたログを表示します。CloudWatchログを表示するには、必要な IAM アクセス許可が付与されている必要があります。詳細については、Amazon CloudWatch Logs ユーザーガイドの「[AWS CloudWatch ログの管理 \(定義済み\) ポリシー](#)」を参照してください。

すべての AWS コールを VPC を経由するように設定する

特別なジョブパラメータ `disable-proxy-v2` を使用すると、VPC 経由で Amazon S3、CloudWatch、AWS Glue などのサービスに呼び出しをルーティングできます。デフォルトでは、AWS Glue はローカルプロキシを使用して AWS Glue VPC 経由でトラフィックを送信します。これにより、Amazon S3 からスクリプトとライブラリをダウンロードし、ログとメトリクスの公開のために CloudWatch にリクエストを送信し、データカタログにアクセスするためのリクエストを AWS Glue に送信します。このプロキシは、VPC が Amazon S3、CloudWatch、AWS Glue など、他の AWS のサービスに向かう適切なルートを設定していなくても、ジョブを正常に機能させること

ができます。AWS Glue は、この動作をオフにするためのパラメータを提供しています。詳細については、「[AWS Glue によって使用される Job パラメータ](#)」を参照してください。AWS Glue は、引き続きローカルプロキシを使用して、AWS Glue ジョブの CloudWatch ログを公開します。

Note

- この機能は、AWS Glue バージョン 2.0 以降の AWS Glue ジョブでサポートされています。この特徴を利用する場合、VPC が NAT または サービス VPC エンドポイントを経由して Amazon S3 へのルートを設定していることを確認する必要があります。
- 非推奨のジョブパラメータ `disable-proxy` は、VPC 経由でスクリプトとライブラリをダウンロードするため、Amazon S3 にのみ呼び出しをルーティングします。代わりに、新しいパラメータ `disable-proxy-v2` を使用することをお勧めします。

使用例

AWS Glue で `disable-proxy-v2` のジョブを作成します。

```
aws glue create-job \  
  --name no-proxy-job \  
  --role GlueDefaultRole \  
  --command "Name=glueetl,ScriptLocation=s3://my-bucket/glue-script.py" \  
  --connections Connections="traffic-monitored-connection" \  
  --default-arguments '{"--disable-proxy-v2" : "true"}'
```

VPC の JDBC データストアに接続する

通常、これらのリソースは、パブリックインターネット経由でアクセスできないように、Amazon Virtual Private Cloud (Amazon VPC) 内に作成します。デフォルトでは、AWS Glue は VPC 内のリソースにアクセスできません。AWS Glue が VPC 内のリソースにアクセスできるようにするには、VPC サブネット ID やセキュリティグループ ID など、追加の VPC 固有設定情報を指定する必要があります。AWS Glue はこの情報を、関数がプライベート VPC 内の他のリソースに安全に接続できる [Elastic Network Interface](#) のセットアップに使用します。

VPC エンドポイントを使用する場合は、そのエンドポイントをルートテーブルに追加します。詳細については、「[AWS Glue 用のインターフェイス VPC エンドポイントの作成](#)」と「[前提条件](#)」を参照してください。

データカタログで暗号化を使用する場合は、KMS インターフェイスエンドポイントを作成してルートテーブルに追加します。詳細については、「[AWS KMS 用の VPC エンドポイントの作成](#)」を参照してください。

Elastic Network Interface を使用して VPC データにアクセスする

AWS Glue が VPC 内の JDBC データストアに接続する場合、AWS Glue は VPC データにアクセスするために、アカウントに Elastic Network Interface (プレフィックス Glue_) を作成します。AWS Glue にアタッチされている限り、このネットワークインターフェイスを削除することはできません。Elastic Network Interface 作成の一部として、AWS Glue はこれに 1 つ以上のセキュリティグループを関連付けます。AWS Glue がネットワークインターフェイスを作成できるようにするには、リソースに関連付けられているセキュリティグループがソースルールを使用したインバウンドアクセスを許可する必要があります。このルールには、リソースに関連付けられたセキュリティグループが含まれています。これにより、Elastic Network Interface は同じセキュリティグループを持つデータストアにアクセスできるようになります。

AWS Glue がコンポーネントと通信できるようにするには、すべての TCP ポートに対して自己参照のインバウンドルールを持つセキュリティグループを指定します。自己参照ルールを作成することで、ソースをすべてのネットワークではなく VPC 内の同じセキュリティグループに制限することができます。VPC のデフォルトのセキュリティグループには、すでに ALL Traffic の自己参照インバウンドルールがある場合があります。

Amazon VPC コンソールでルールを作成します。AWS Management Console を介してルールの設定を更新するには、VPC コンソール (<https://console.aws.amazon.com/vpc/>) に移動し、適切なセキュリティグループを選択します。ALL TCP のインバウンドルールを指定して、同じセキュリティグループ名をソースとして指定します。セキュリティグループルールの詳細については、「[VPC のセキュリティグループ](#)」を参照してください。

それぞれの Elastic Network Interface には、指定したサブネット内の IP アドレス範囲からプライベート IP アドレスが割り当てられます。ネットワークインターフェイスにパブリック IP アドレスが割り当てられることはありません。AWS Glue にはインターネットアクセスが必要です (例えば、VPC エンドポイントのない AWS のサービスにアクセスする場合など)。ネットワークアドレス変換 (NAT) インスタンスを VPC 内で設定するか、Amazon VPC NAT ゲートウェイを使用することができます。詳細については、Amazon VPC ユーザーガイドの「[NAT ゲートウェイ](#)」を参照してください。ネットワークインターフェイスにはパブリック IP アドレスが必要なため、VPC にアタッチされたインターネットゲートウェイをサブネットルートテーブルのルートとして直接使用することはできません。

VPC ネットワーク属性の `enableDnsHostnames` および `enableDnsSupport` を `true` に設定する必要があります。詳細については、「[Using DNS with Your VPC](#)」を参照してください。

Important

インターネットアクセスのないパブリックサブネットまたはプライベートサブネットにデータストアを配置しないでください。代わりに、NAT インスタンスまたは Amazon VPC NAT ゲートウェイを介して、インターネットアクセスのあるプライベートサブネットにのみアタッチしてください。

Elastic Network Interface プロパティ

Elastic Network Interface を作成するには、次のプロパティを指定する必要があります。

VPC

データストアを含む VPC 名。

サブネット

データストアを含む VPC 内のサブネット。

セキュリティグループ

データストアに関連付けられているセキュリティグループ。AWS Glue は、VPC サブネットにアタッチされている Elastic Network Interface にこれらのセキュリティグループを関連付けます。AWS Glue コンポーネントの通信を可能にし、他のネットワークからのアクセスを禁止するには、少なくとも1つの選択されたセキュリティグループにおいて、すべての TCP ポートの自己参照のインバウンドルールを指定する必要があります。

Amazon Redshift を使用した VPC の管理については、「[Amazon Virtual Private Cloud \(VPC\)でクラスターを管理する](#)」を参照してください。

Amazon Relational Database Service (Amazon RDS) を使用した VPC の管理については、「[VPC 内の Amazon RDS DB インスタンスの使用](#)」を参照してください。

MongoDB または MongoDB Atlas 接続を使用する

MongoDB または MongoDB Atlas の接続を作成すると、ETL ジョブで接続を使用することができます。AWS Glue Data Catalog にテーブルを作成し、テーブルの connection 属性に MongoDB または MongoDB Atlas 接続を指定します。

AWS Glue では、接続 url と MongoDB 接続の認証情報が保存されます。接続 URI の形式は次のとおりです。

- MongoDB の場合: `mongodb://host:port/database` ホストは、ホスト名、IP アドレス、UNIX ドメインソケットのいずれかにすることができます。接続文字列でポートが指定されていない場合は、デフォルトの MongoDB ポート、27017 が使用されます。
- MongoDB Atlas: `mongodb+srv://server.example.com/database` ホストは、DNS SRV レコードに対応するホスト名にすることができます。SRV 形式ではポートは不要であり、デフォルトの MongoDB ポート、27017 が使用されます。

さらに、ジョブスクリプトでオプションを指定できます。詳細については、「[the section called “MongoDB 接続”](#)」を参照してください。

VPC エンドポイントを使用して Amazon S3 データストアをクローラーする

セキュリティ、監査、またはコントロールの目的で、Amazon Virtual Private Cloud 環境 (Amazon VPC) を介してのみ Amazon S3 データストアまたは Amazon S3 backed データカタログテーブルにアクセスできるようにしたいことがあります。このトピックでは、Network 接続タイプを使用して VPC エンドポイントで Amazon S3 データストアまたは Amazon S3 backed データカタログテーブルへの接続を作成およびテストする方法について説明します。

データストアでクローラーを実行するには、次のタスクを実行します。

- [the section called “前提条件”](#)
- [the section called “Amazon S3 への接続を作成する”](#)
- [the section called “Amazon S3 への接続をテストする”](#)
- [the section called “Amazon S3 データストアのクローラーを作成する”](#)
- [the section called “クローラーの実行”](#)

前提条件

Amazon Virtual Private Cloud 環境 (Amazon VPC) を介してアクセスされるように、Amazon S3 データストアまたは Amazon S3 backed データカタログテーブルを設定するための前提条件を満たしていることを確認します。

- 設定済みの VPC。例: vpc-01685961063b0d84b。詳細については、Amazon VPC ユーザーガイドの「[Amazon VPC の使用を開始する](#)」を参照してください。
- VPC にアタッチされた Amazon S3 エンドポイント。例: vpc-01685961063b0d84b。詳細については、Amazon VPC ユーザーガイドの「[Amazon S3 におけるエンドポイント](#)」を参照してください。

The screenshot shows the AWS Management Console interface for a VPC. At the top, there is a search bar with the text "Name: privateVPC" and a filter icon. Below the search bar is a table with columns: Name, VPC ID, State, IPv4 CIDR, IPv6, DHCP options set, Main Route table, Main Network ACL, Tenancy, and Default VPC. The table contains one row for the VPC "privateVPC" with ID "vpc-01685961063b0d84b", state "available", IPv4 CIDR "192.168.1.0/24", and other details. Below the table, there is a section titled "VPC: vpc-01685961063b0d84b" with tabs for "Description", "CIDR Blocks", "Flow Logs", and "Tags". The "Description" tab is selected, showing a list of attributes and their values.

Attribute	Value
VPC ID	vpc-01685961063b0d84b
State	available
IPv4 CIDR	192.168.1.0/24
IPv6 Pool	-
Network ACL	acl-02d197f2c9fbc46be
DHCP options set	dopt-a79e5acc
Owner	261353713322
Tenancy	default
Default VPC	No
IPv6 CIDR	-
DNS resolution	Enabled
DNS hostnames	Disabled
Route table	rtb-0750198567d5b5202

- VPC エンドポイントを指すルートエントリ。例えば、VPC エンドポイント (vpce-0ec5da4d265227786) で使用されるルートテーブルの vpce-0ec5da4d265227786。

Name	Route Table ID	Explicit subnet association	Edge associations	Main	VPC ID
	rtb-0750198567d5b5202	-	-	Yes	vpc-01685961063b0d84b ...

Route Table: rtb-0750198567d5b5202

Summary

Routes

Subnet Associations

Edge Associations

Route Propagation

Tags

Edit routes

View All routes

Destination	Target	Status	Propagate
192.168.1.0/24	local	active	No
pl-7ba54012 (com.amazonaws.us-east-2.s3, 52.219.80.0/20, 3.5.128.0/22, 3.5.132.0/23, 52.219.96.0/20, 52.92.76.0/22)	vpce-0ec5da4d265227786	active	No

- VPC にアタッチされたネットワーク ACL により、トラフィックが可能になります。
- VPC にアタッチされたセキュリティグループにより、トラフィックが可能になります。

Amazon S3 への接続を作成する

通常、これらのリソースは、パブリックインターネット経由でアクセスできないように、Amazon Virtual Private Cloud (Amazon VPC) 内に作成します。デフォルトでは、AWS Glue は VPC 内のリソースにアクセスできません。AWS Glue が VPC 内のリソースにアクセスできるようにするには、VPC サブネット ID やセキュリティグループ ID など、追加の VPC 固有設定情報を指定する必要があります。Network 接続を作成するには、次の情報を指定することが必要です。

- VPC ID
- VPC 内のサブネット
- セキュリティグループ

Network 接続をセットアップするには

1. AWS Glue コンソールのナビゲーションペインで、[Add connection] (接続の追加) をクリックします。
2. 接続名を入力し、接続タイプとして、[Network] (ネットワーク) を選択します。[Next] を選択します。

Add connection



Connection properties

Connection access

Review all steps

Set up your connection's properties.

For more information, see [Working with Connections](#).

Connection name

Connection type

Description (optional)

3. VPC、サブネット、およびセキュリティグループの情報を設定します。

- VPC: データストアを含む VPC 名を選択します。
- サブネット: VPC 内のサブネットを選択します。
- セキュリティグループ: VPC 内のデータストアへのアクセスを許可する 1 つ以上のセキュリティグループを選択します。

Add connection ✕

- ✔ **Connection properties**
TestNetworkConnecti
on
Type: Network
- **Connection access**
- Review all steps

Set up access to your data store.

For more information, see [Working with Connections](#).

VPC
Choose the VPC name that contains your data store.

vpc-01685961063b0d84b | privateVPC

Subnet
Choose the subnet within your VPC.

subnet-0b350d86953aa6d60 | Range192

Security groups
Choose one or more security groups that allow access to the data store in your VPC. AWS Glue associates these security groups to the ENI attached to your subnet. To allow AWS Glue components to communicate and also prevent access from other networks, at least one chosen security group must specify a self-referencing inbound rule for all TCP ports.

<input checked="" type="checkbox"/> Group ID	Group name
<input checked="" type="checkbox"/> sg-0ce8b36fb6206c56e	default

[Back](#) [Next](#)

4. [Next] を選択します。
5. 接続情報を確認し、[Finish] (完了) をクリックします。

Add connection ✕

- Connection properties**
TestNetworkConnection
Type: Network
- Connection access**
VPC Id:
vpc-01685961063b0d84b
- Review all steps**

Connection properties

Name	TestNetworkConnection
Type	Network
Description (optional)	This is a demo Network Connection

Connection access

VPC Id	vpc-01685961063b0d84b
Subnet	subnet-0b350d86953aa6d60
Security groups	sg-0ce8b36fb6206c56e

[Back](#) [Finish](#)

Amazon S3 への接続をテストする

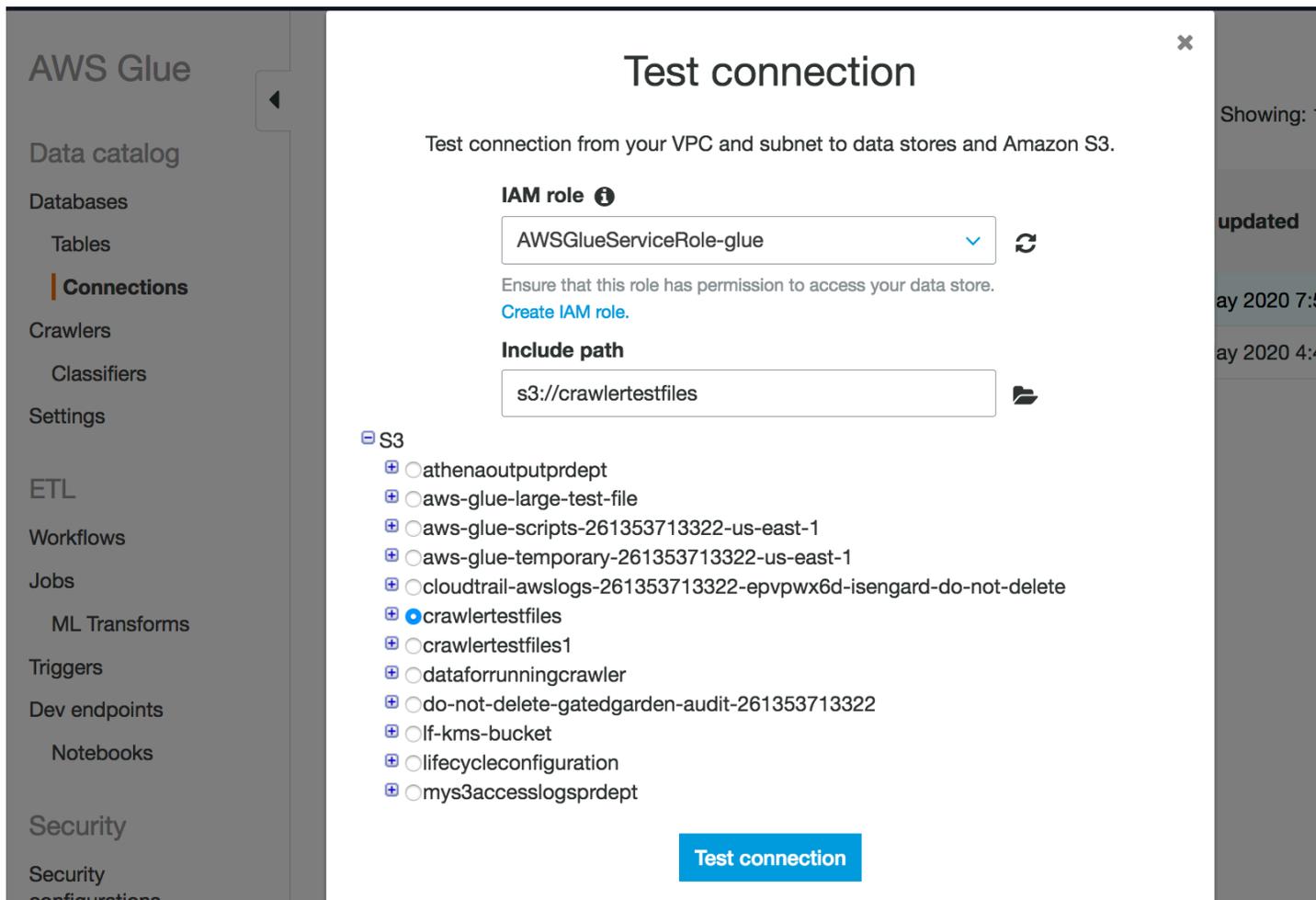
Network 接続を作成したら、VPC エンドポイントの Amazon S3 データストアへの接続をテストできます。

接続のテスト時に、次のエラーが発生することがあります。

- **INTERNET CONNECTION ERROR (インターネット接続エラー):** インターネット接続に問題があることを示します。
- **INVALID BUCKET ERROR (無効なバケットエラー):** Amazon S3 バケットに問題があることを示します。
- **S3 CONNECTION ERROR (S3 接続エラー):** Amazon S3 への接続に失敗したことを示します。
- **INVALID CONNECTION TYPE (無効な接続タイプ):** 接続タイプが想定される NETWORK という値でないことを示します。
- **INVALID CONNECTION TEST TYPE (無効な接続テストタイプ):** ネットワーク接続テストのタイプに問題があることを示します。
- **INVALID TARGET (無効なターゲット):** Amazon S3 バケットが正しく指定されていないことを示します。

Network 接続をテストするには

1. AWS Glue コンソールで [Network] (ネットワーク) 接続を選択します。
2. [Test connection] を選択します。
3. 前のステップで作成した IAM ロールを選択し、Amazon S3 バケットを指定します。
4. テストを開始するには、[Test connection] (接続のテスト) をクリックします。結果を表示するには少し時間がかかることがあります。



エラーが発生した場合は、次のチェックを行います。

- 選択したロールに正しい権限が与えられている。
- 正しい Amazon S3 バケットが指定されている。
- セキュリティグループとネットワーク ACL により、必要な着信トラフィックと発信トラフィックが許可されている。
- 指定した VPC が、Amazon S3 VPC エンドポイントに接続されている。

接続のテストに成功したら、クローラーを作成できます。

Amazon S3 データストアのクローラーを作成する

これで、作成した Network 接続を指定したクローラーを作成できます。クローラーの作成についての詳細については、「[クローラーの設定](#)」を参照してください。

1. AWS Glue コンソールのナビゲーションペインで [Crawlers] (クローラー) を選択します。
2. [Add crawler (クローラーの追加)] を選択します。
3. クローラーの名前を指定して、[Next] (次へ) をクリックします。
4. データソースを要求されたら、[S3] (S3) を選択し、Amazon S3 バケットプレフィックスと前に作成した接続を指定します。

The screenshot shows the 'Add crawler' wizard in the AWS Glue console. The wizard is titled 'Add crawler' and has a close button (X) in the top right corner. On the left side, there is a navigation pane with the following steps: 'Crawler info' (checked), 'TestNetworkConnection', 'Crawler source type' (checked), 'Data stores', 'Data store' (selected), 'IAM Role', 'Schedule', 'Output', and 'Review all steps'. The main area is titled 'Add a data store' and contains the following fields and options:

- Choose a data store:** A dropdown menu with 'S3' selected.
- Connection:** A dropdown menu with 'AddNetworkConnection' selected.
- Connection note:** 'Optionally include a Network connection to use with this S3 target. Note that each crawler is limited to one Network connection so any future S3 targets will also use the same connection (or none, if left blank).' Below this is an 'Add connection' button.
- Crawl data in:** A radio button labeled 'Specified path' is selected.
- Include path:** A text input field containing 's3://crawlerstestfiles' with a file selection icon to its right.
- Include path note:** 'All folders and files contained in the include path are crawled. For example, type s3://MyBucket/MyFolder/ to crawl all objects in MyFolder within MyBucket.'
- Exclude patterns (optional):** A section with a right-pointing arrow.
- Navigation:** 'Back' and 'Next' buttons at the bottom.

5. 必要に応じて、同じネットワーク接続に別のデータストアを追加します。
6. IAM ロールを選択します。IAM ロールでは、AWS Glue サービスと Amazon S3 バケットへのアクセスが許可されている必要があります。詳細については、「[the section called “クローラーの設定”](#)」を参照してください。

Add crawler

✓ Crawler info

TestNetworkConnecti
on

✓ Crawler source type

Data stores

✓ Data store

S3: s3://crawlertestf...

○ IAM Role

○ Schedule

○ Output

○ Review all steps

Choose an IAM role

The IAM role allows the crawler to run and access your Amazon S3 data stores. [Learn more](#)

- Update a policy in an IAM role
- Choose an existing IAM role
- Create an IAM role

IAM role ⓘ

AWSGlueServiceRole-glue



This role must provide permissions similar to the AWS managed policy, **AWSGlueServiceRole**, plus access to your data stores.

- s3://crawlertestfiles

You can also create an IAM role on the [IAM console](#).

Back

Next

7. クローラーのスケジュールを定義します。

8. データカタログの既存のデータベースを選択するか、新しいデータベースエントリを作成します。

Add crawler



✓ Crawler info

TestNetworkConnecti
on

✓ Crawler source type

Data stores

✓ Data store

S3: s3://crawlertestf...

✓ IAM Role

arn:aws:iam::2613537
13322:role/service-
role/AWSGlueService
Role-glue

✓ Schedule

Run on demand

○ Output

○ Review all steps

Configure the crawler's output

Database ⓘ

testnetworkconnectiondb

Add database

Prefix added to tables (optional) ⓘ

Type a prefix added to table names

▶ Grouping behavior for S3 data (optional)

▶ Configuration options (optional)

Back

Next

9. 残りのセットアップを完了します。

Amazon S3 backed データカタログテーブルのクローラーを作成する

これで、作成した Network 接続とカタログソースタイプを指定するクローラーを作成できます。クローラーの作成についての詳細については、「[クローラーの設定](#)」を参照してください。

1. AWS Glue コンソールのナビゲーションペインで [Crawlers] (クローラー) を選択します。
2. [Add crawler (クローラーの追加)] を選択します。
3. クローラーの名前を指定して、[Next] (次へ) をクリックします。
4. クローラーソースタイプを要求されたら、[既存のカタログテーブル] を選択し、使用可能なテーブルのリストからクロールする既存のカタログテーブルを指定します。

The screenshot shows the 'Add crawler' wizard in the AWS Glue console, specifically the 'Choose catalog tables' step. On the left, a navigation pane shows the progress: 'Crawler info' (test) is complete, 'Crawler source type' (Existing catalog tables) is selected, and 'Catalog tables', 'IAM Role', 'Schedule', 'Output', and 'Review all steps' are pending. The main area is titled 'Choose catalog tables' and contains two tables: 'Selected tables' (currently empty) and 'Available tables' (listing three tables). At the bottom, there is a 'Connection' dropdown menu set to 'Select a connection' and an 'Add connection' button.

Selected tables				Showing: 0 - 0 <>
Name	Database	Location	Classification	
No items selected				

Available tables				Showing: 1 - 3 <>
Name	Database	Location	Classification	
Add s3_event_crawl_demo	test-sampling-db	s3://s3-event-crawl-demo/	json	
Add test_int5100_idf_20210310094002_0800_obfusca...	test-large-xml	s3://crawlertickets/TEST_INT5100_IDF_20210310...	Unknown	
Add test_int5100_idf_20210310094002_0800_obfusca...	test-cx-whitelist	s3://crawlertickets/TEST_INT5100_IDF_20210310...	xml	

Connection: Select a connection [v]
Add connection

5. IAM ロールを選択します。IAM ロールでは、AWS Glue サービスと Amazon S3 バケットへのアクセスが許可されている必要があります。詳細については、「[the section called “クローラーの設定”](#)」を参照してください。
6. クローラーのスケジュールを定義します。
7. データカタログの既存のデータベースを選択するか、新しいデータベースエントリを作成します。
8. 残りのセットアップを完了し、ステップを確認します。

Add crawler
✕

- Crawler info**
test
- Crawler source type**
Existing catalog tables
- Catalog tables**
test_int5100_idf_20...
- IAM Role**
arn:aws:iam::804918391416:role/service-role/AWSGlueServiceRole-crawler-test
- Schedule**
Run on demand
- Output**
- Review all steps**

Use Lake Formation Data Catalog

Catalog tables

Database	test-large-xml
Table name	test_int5100_idf_20210310094002_0800_obfuscated_xml
Connection	test

IAM role

IAM role: arn:aws:iam::804918391416:role/service-role/AWSGlueServiceRole-crawler-test

Schedule

Schedule: Run on demand

Output

Database	
Prefix added to tables (optional)	
Create a single schema for each S3 path	true
Table level (optional)	
Data Lineage (optional)	DISABLE

▶ Configuration options

Back Finish

クローラーの実行

クローラーを実行します。

AWS Glue

Data catalog

Databases

Tables

Connections

Crawlers

A crawler connects to a data store, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

Crawler **TestNetworkConnection** was created to run on demand. [Run it now?](#) ✕

[User preferences](#)

トラブルシューティング

VPC ゲートウェイを使用した Amazon S3 バケットに関するトラブルシューティングについては、「[ゲートウェイ VPC エンドポイントを使用して S3 バケットに接続できないのはなぜですか](#)」を参照してください。

AWS Glue での接続の問題のトラブルシューティング

AWS Glue クローラーまたはジョブが、データストアにアクセスするため接続プロパティを使用する場合、接続しようとするエラーが発生する場合があります。指定したVirtual Private Cloud (VPC) とサブネットに Elastic Network Interface を作成する際に、AWS Glue はサブネット内のプライベー

ト IP アドレスを使用しています。接続で指定されたセキュリティグループは、各 Elastic Network Interface に適用されます。セキュリティグループがアウトバウンドアクセスを許可しているかどうか、またデータベースクラスターへの接続を許可しているかどうかを確認します。

さらに、Apache Spark ではドライバーノードとエグゼキューターノード間の双方向接続が必要です。セキュリティグループのいずれかが、すべての TCP ポートの進入ルールを許可する必要があります。自己参照のセキュリティグループでセキュリティグループのソースを自身に制限することによって、世界に向けて開かれないようにすることができます。

接続の問題をトラブルシューティングするために行える一般的なアクションをいくつか示します。

- 接続のポートアドレスを確認します。
- 接続もしくはシークレットの、ユーザー名とパスワードの文字列を確認します。
- JDBC データストアの場合、着信接続を許可していることを検証します。
- データストアが VPC 内でアクセスできることを検証します。
- AWS Secrets Manager を使用して接続の認証情報を保存する場合、AWS Glue のための IAM ロールには、対象のシークレットにアクセスするための許可が必要です。詳細については、「[AWS Secrets Manager ユーザーガイド](#)」の「[例: シークレット値を取得するアクセス許可](#)」を参照してください。ネットワークの設定によっては、VPC エンドポイントを作成して VPC と Secrets Manager 間のプライベート接続を確立する必要がある場合もあります。詳細については、「[AWS Secrets Manager VPC エンドポイントの使用](#)」を参照してください。

チュートリアル: Elasticsearch 向けの AWS Glue コネクタを使用する

Elasticsearch は、ログ分析、リアルタイムアプリケーションモニタリング、クリックストリーム分析などのユースケースで人気のあるオープンソースの検索および分析エンジンです。

AWS Glue Studio で Elasticsearch 向けの AWS Glue コネクタを設定することで、OpenSearch を抽出、変換、ロード (ETL) ジョブ用のデータストアとして使用できます。このコネクタは [AWS Marketplace](#) から無料で入手できます。

Note

[AWS Marketplace Elasticsearch Spark コネクタ](#) は廃止されました。代わりに、[Elasticsearch 向けの AWS Glue コネクタ](#) を使用してください。

このチュートリアルでは、最小限の手順で Amazon OpenSearch Service ノードに接続する方法を示します。

トピック

- [前提条件](#)
- [ステップ 1: \(オプション\) OpenSearch クラスター情報の AWS シークレットを作成する](#)
- [ステップ 2: コネクタをサブスクライブする](#)
- [ステップ 3: AWS Glue Studio でコネクタをアクティブにして、接続を作成します。](#)
- [ステップ 4: ETL ジョブの IAM ロールを設定する](#)
- [ステップ 5: OpenSearch 接続を使用するジョブを作成する](#)
- [ステップ 6: ジョブを実行する](#)

前提条件

このチュートリアルを使用するには、以下のものがが必要です。

- AWS Glue Studio へのアクセス
- AWS Cloud の OpenSearch クラスターにアクセスする
- (オプション) AWS Secrets Manager にアクセスします。

ステップ 1: (オプション) OpenSearch クラスター情報の AWS シークレットを作成する

接続の認証情報を安全に保存して使用するには、認証情報を AWS Secrets Manager に保存します。作成するシークレットは、チュートリアルの後半で接続で使用されます。認証情報のキーと値のペアは、通常の接続オプションとして Elasticsearch 向けの AWS Glue コネクタに供給されます。

シークレットの作成についての詳細は、AWS Secrets Manager ユーザーガイドの「[AWS Secrets Manager を使用したシークレットの作成と管理](#)」を参照してください。

AWS シークレットを作成するには

1. [AWS Secrets Manager コンソール](#) にサインインします。
2. サービスの概要ページまたは [シークレット] リストページで、[新しいシークレットの保存] を選択します。

3. [新しいシークレットの保存] ページで、[他の種類のシークレット] を選択します。このオプションの意味は、シークレットの構造と詳細を指定する必要があるということです。
4. の追加キーおよび値のペアを OpenSearch クラスターユーザー名として使用します。例:

```
es.net.http.auth.user: username
```
5. [+ 行を追加] を選択して、パスワードとして別のキーと値のペアを入力します。例:

```
es.net.http.auth.pass: password
```
6. [Next] を選択します。
7. シークレット名を入力します 例えば、my-es-secret です。オプションで説明を含めることができます。

このチュートリアルの後半で使用するシークレット名を記録し、[Next] (次へ) を選択します。

8. [Next] (次へ) をもう一度選択して、[Store] (保存) を選択してシークレットを作成します。

次のステップ

[ステップ 2: コネクタをサブスクライブする](#)

ステップ 2: コネクタをサブスクライブする

Elasticsearch 向けの AWS Glue コネクタは、[AWS Marketplace](#) で無料で入手できます。

AWS Marketplace で Elasticsearch 向けの AWS Glue コネクタにサブスクライブするには

1. まだ License Manager を使用するための AWS アカウント設定をしていない場合、次の操作を行います。
 - a. AWS License Manager コンソール、<https://console.aws.amazon.com/license-manager> を開きます。
 - b. Create customer managed license (カスタマー管理ライセンスの作成) を選択します。
 - c. [IAM permissions (one-time setup)] 画面で [I grant AWS License Manager the required permissions] を選択し、次に [Grant permissions] を選択します。

このウィンドウが表示されない場合は、必要なアクセス許可がすでに設定されています。

2. <https://console.aws.amazon.com/gluestudio/> で AWS Glue Studio コンソールを開きます。

3. AWS Glue Studio コンソールのナビゲーションペインで、メニューアイコン (☰) を展開して、[Connectors (コネクタ)] を選択します。
4. [Connectors] (コネクタ) ページで、[Go to AWS Marketplace] をクリックします。
5. AWS Marketplace の [AWS Glue Studio の製品を検索] セクションで、検索フィールドに「Elasticsearch 向けの AWS Glue コネクタ」と入力し、Enter キーを押します。
6. [Elasticsearch 向けの AWS Glue コネクタ] というコネクタの名前を選択します。
7. コネクタの製品ページで、そのコネクタに関するタブを開いて情報を表示します。続行する準備が整ったら、[Continue to Subscribe] (スクライブを続行する) を選択します。
8. ソフトウェアの利用規約を確認します。[Accept Terms] (規約に同意する) をクリックします。
9. サブスクリプションプロセスが完了すると、「Thank you for subscribing to this product! You can now configure your software.」(製品にサブスクライブしていただきありがとうございます。ソフトウェアの設定が可能になりました。) という通知が表示されます。バナーの上に [Continue to Configuration] (設定に進む) ボタンがあります。[Continue to Configuration (設定に進む)] を選択します。
10. [Configure this software] (このソフトウェアを設定) ページで、フルフィルメントオプションを選択します。AWS Glue 1.0/2.0 または AWS Glue 3.0 を選択できます。[Continue to Launch] (続行して起動する) を選択します。

次のステップ

[ステップ 3: AWS Glue Studio でコネクタをアクティブにして、接続を作成します。](#)

ステップ 3: AWS Glue Studio でコネクタをアクティブにして、接続を作成します。

[Continue to Launch (起動を続ける)] を選択すると、AWS Marketplace の [Launch this software] (ソフトウェアの起動) ページが表示されます。AWS Glue Studio でリンクを使用してコネクタをアクティブ化した後、接続を作成します。

コネクタを展開し、AWS Glue Studio で接続を作成する方法

1. AWS Marketplace コンソールの [Launch this software] (ソフトウェアの起動) ページで、Usage Instructions (使用方法) を選択し、表示されるウィンドウでリンクを選択します。

ブラウザは AWS Glue Studio コンソールの 「Create marketplace connection」 ページにリダイレクトされます。

2. 接続の名前を入力します。例えば、my-es-connection です。
3. [Connection access] (接続アクセス) セクションの [Connection credential type] (接続の認証情報タイプ) で、[User name and password] (ユーザー名とパスワード) を選択します。
4. [AWS シークレット] で、シークレット名を入力します。例えば、my-es-secret です。
5. [Network options] (ネットワークオプション) セクションで、VPC 情報を入力して OpenSearch クラスターに接続します。
6. [Create connection and activate connector] (接続の作成とコネクタのアクティブ化) を選択します。

次のステップ

[ステップ 4: ETL ジョブの IAM ロールを設定する](#)

ステップ 4: ETL ジョブの IAM ロールを設定する

AWS Glue ETL ジョブを作成するとき、使用するジョブの AWS Identity and Access Management (IAM) ロールを指定します。ロールは、Amazon S3 (ソース、ターゲット、スクリプト、ドライバーファイル、一時ディレクトリ) および AWS Glue Data Catalog オブジェクトを含む、ジョブで使用するすべてのリソースへのアクセスを許可する必要があります。

AWS Glue ETL ジョブで想定される IAM ロールは、前のセクションで作成したシークレットにもアクセスできる必要があります。デフォルトでは、AWS が管理するロール AWSGlueServiceRole はシークレットにアクセスできません。シークレットのアクセスコントロールを設定するには、「[AWS Secrets Manager の認証とアクセスコントロール](#)」および「[特定のシークレットへのアクセスの制限](#)」を参照してください。

ETL ジョブの IAM ロールを設定するには

1. [the section called “ETL ジョブに必要な IAM アクセス許可を確認する”](#) で説明しているように、アクセス許可を設定します。
2. AWS Glue Studio で説明しているように、[the section called “コネクタの使用に必要なアクセス許可”](#) でコネクタを使用するときに必要な追加のアクセス許可を設定します。

次のステップ

[ステップ 5: OpenSearch 接続を使用するジョブを作成する](#)

ステップ 5: OpenSearch 接続を使用するジョブを作成する

ETL ジョブのロールを作成した後、AWS Glue Studio で Open Spark Elasticsearch の接続とコネクタを使用するジョブを作成できます。

ジョブが Amazon Virtual Private Cloud (Amazon VPC) で実行されている場合は、VPC が正しく設定されていることを確認してください。詳細については、「[the section called “ETL ジョブの VPC を設定します”](#)」(ETL ジョブの VPC を設定します) を参照してください。

Elasticsearch Spark コネクタを使用するジョブを作成するには

1. AWS Glue Studio で コネクタ を選択します。
2. [Your connections] (接続) リストで、作成した接続を選択し、[Create job] を選択します。
3. ビジュアルジョブエディタで、[Data source node] (データソースノード) を選択します。右側の [Data source properties - Connector] (データソースのプロパティ-コネクタ) タブで、コネクタの追加情報を設定します。
 - a. [Add schema] (スキーマの追加) を選択し、データソースにデータセットのスキーマを入力します。接続には、データカタログに格納されたテーブルは使用されません。つまり、AWS Glue Studio はデータのスキーマを認識しません。このスキーマ情報は手動で指定する必要があります。スキーマエディタの使用方法については、「[the section called “カラム変換ノードでスキーマを編集する”](#)」を参照してください。
 - b. Connection options (接続オプション) を展開します。
 - c. [Add new option] (新しいオプションの追加) を選択し、AWS シークレットに入力されなかったコネクタに必要な情報を入力します。
 - es.nodes: https://<OpenSearch ドメインエンドポイント>
 - es.port: 443
 - path: test
 - es.nodes.wan.only.: true

これらの接続オプションの詳細については、「<https://www.elastic.co/guide/en/elasticsearch/hadoop/current/configuration.html>」を参照してください。

4. グラフにターゲットノードを追加します。

データターゲットは Amazon S3 にすることも、AWS Glue Data Catalog またはコネクタからの情報を使用して別の場所にデータを書き込むこともできます。例えば、データカタログテーブルを使用して Amazon RDS のデータベースに書き込むことも、コネクタをデータターゲットとして使用して、AWS Glue でネイティブにサポートされていないデータストアに書き込むこともできます。

データターゲットのコネクタを選択する場合は、そのコネクタ用に作成された接続を選択する必要があります。また、コネクタプロバイダが必要な場合は、コネクタに追加情報を提供するオプションを追加する必要があります。AWS シークレットの情報を含む接続を使用する場合は、接続オプションでユーザー名とパスワード認証を指定する必要はありません。

5. オプションで、[the section called “AWS Glue マネージドデータ変換ノードの編集”](#) の説明の通り、追加のデータソースと 1 つ以上の変換ノードを追加します。

6. ステップ 3 から [the section called “ジョブのプロパティを変更する”](#) の説明の通りジョブプロパティを設定し、ジョブを保存します。

次のステップ

[ステップ 6: ジョブを実行する](#)

ステップ 6: ジョブを実行する

ジョブを保存した後、ジョブを実行し ETL 操作が行えるようになります。

Elasticsearch 向けの AWS Glue コネクタ用に作成したジョブを実行するには

1. AWS Glue Studio コンソールを使用して、ビジュアルエディタページで [Run](実行) を選択します。
2. 成功バナーで、[Run Details] (実行の詳細) を選択するか、[Runs] (実行) タブをクリックして、ジョブの実行に関する情報を表示します。

インタラクティブセッションで AWS Glue ジョブを構築する

データエンジニアは AWS Glue のインタラクティブセッションを使用して、以前よりも迅速かつ簡単に AWS Glue のジョブをオーサリングできます。

トピック

- [AWS Glue インタラクティブセッションの概要](#)
- [AWS Glue のインタラクティブセッションの開始方法](#)
- [AWS Glue Jupyter および AWS Glue Studio ノートブックのインタラクティブセッションの設定](#)
- [AWS Glue for Ray インタラクティブセッションの開始方法\(プレビュー\)](#)
- [IAM を使用したインタラクティブセッション](#)
- [スクリプトまたはノートブックの AWS Glue ジョブジョブへの変換](#)
- [AWS Glue ストリーミングのためのインタラクティブセッション](#)
- [AWS Glueスクリプトのローカルでの開発およびテスト](#)
- [開発エンドポイント](#)

AWS Glue インタラクティブセッションの概要

AWS Glue のインタラクティブセッションでは、データ準備および分析アプリケーションの構築、テスト、実行を迅速に行うことができます。インタラクティブセッションは、データ準備のための抽出、変換、ロード (ETL) スクリプトを構築してテストするためのプログラムおよびビジュアルインターフェイスを提供します。インタラクティブセッションでは Apache Spark 分析アプリケーションが実行され、リモートの Spark ランタイム環境にオンデマンドでアクセスできます。AWS Glue はこれらのインタラクティブセッションのサーバーレス Spark を透過的に管理します。

インタラクティブセッションには柔軟性が備わっているため、お好みの環境からアプリケーションを構築し、テストできます。インタラクティブセッションは、AWS Command Line Interface とAPIを介して作成、操作することができます。Jupyter と互換性のあるノートブックを使用して、ノートブックスクリプトを視覚的に作成し、テストできます。インタラクティブセッションは、PyCharm、IntelliJ、VS Code などの IDE との統合を含め、Jupyter が統合するほぼすべてと統合できるオープンソースの Jupyter カーネルを提供します。これにより、ローカル環境でコードをオーサリングし、インタラクティブセッションのバックエンドでシームレスに実行できます。

インタラクティブセッション API を使用すると、Spark インフラストラクチャを管理する必要なく Apache Spark 分析を使用するアプリケーションをプログラムで実行できます。単一のインタラクティブセッション内で 1 つ以上の Spark ステートメントを実行できます。

したがって、インタラクティブセッションは、より速く、安価で、柔軟にデータ準備および分析アプリケーションを構築して実行することができます。インタラクティブセッションの使用方法については、このセクションのドキュメントを参照してください。[AWS Glue がサポートするマジック](#)

制限事項

- ジョブのブックマークは、インタラクティブセッションではサポートされていません。
- AWS Command Line Interface を使用したノートブックの作成はサポートされていません。

AWS Glue のインタラクティブセッションの開始方法

これらのセクションでは、AWS Glue の対話型セッションをローカルで実行する方法について説明します。

ローカルでインタラクティブセッションを設定するための前提条件

インタラクティブセッションをインストールするための前提条件は次のとおりです。

- サポートされている Python バージョンは 3.6~3.10 以降です。
- MacOS/Linux および Windows の手順については、以下のセクションを参照してください。

Jupyter AWS Glue とインタラクティブセッション Jupyter カーネルのインストール

次を使用してカーネルをローカルにインストールします。

`install-glue-kernels` コマンドは、`pyspark` カーネルと `Spark` カーネルの両方に `jupyter kernelspec` をインストールし、また適切なディレクトリにロゴをインストールします。

```
pip3 install --upgrade jupyter boto3 aws-glue-sessions
```

```
install-glue-kernels
```

Jupyter の実行

Jupyter Notebook を実行するには、以下の手順に従います。

1. 次のコマンドを実行して Jupyter Notebook をインストールします。

```
jupyter notebook
```

2. [New] (新規) を選択して、次の AWS Glue カーネルのいずれかを選択して AWS Glue に対するコーディングを開始します。

セッション認証情報とリージョンの設定

MacOS/Linux の手順

AWS Glue インタラクティブセッションには、AWS Glue ジョブおよび開発エンドポイントと同じ IAM アクセス許可が必要です。インタラクティブセッションで使用するロールは、次のいずれかの方法で指定します。

1. `%iam_role` と `%region` マジックを使用
2. `~/.aws/config` の追加行を使用

マジックを使用したセッションロールの設定

最初のセルで、実行された最初のセルに「`%iam_role <YourGlueServiceRole>`」と入力します。

`~/.aws/config` を使用したセッションロールの設定

AWS Glue インタラクティブセッションのサービスロールは、ノートブック自体に指定することも、設定と一緒に保存することもできます。AWS CLI AWS Glue ジョブで通常使用するロールがある場合、これがロールになります。AWS Glue ジョブに使用するロールがない場合は、「[AWS Glue 用の IAM アクセス許可の設定](#)」のガイドに従ってセットアップしてください。

このロールをインタラクティブセッションのデフォルトのロールとして設定するには、以下の手順を実施します。

1. テキストエディタで `~/.aws/config` を開きます。

2. AWS Glue に使用するプロファイルを探します。いずれのプロファイルも使わない場合は、[Default] プロファイルを使用します。
3. `glue_role_arn=<AWSGlueServiceRole>` に示すように、使用するロールのプロファイルに 1 行を追加します。
4. [オプション]: プロファイルにデフォルトのリージョンセットがない場合は、`region=us-east-1` でリージョンを追加し、`us-east-1` を対象リージョンに置き換えることをお勧めします。
5. 設定を保存します。

詳細については、「[IAM を使用したインタラクティブセッション](#)」を参照してください。

Windows の手順

AWS Glue インタラクティブセッションには、AWS Glue ジョブおよび開発エンドポイントと同じ IAM アクセス許可が必要です。インタラクティブセッションで使用するロールは、次のいずれかの方法で指定します。

1. `%iam_role` と `%region` マジックを使用
2. `~/.aws/config` の追加行を使用

マジックを使用したセッションロールの設定

最初のセルで、実行された最初のセルに「`%iam_role <YourGlueServiceRole>`」と入力します。

`~/.aws/config` を使用したセッションロールの設定

AWS Glue インタラクティブセッションのサービスロールは、ノートブック自体に指定することも、AWS CLI 設定と一緒に保存することもできます。AWS Glue ジョブで通常使用するロールがある場合、これがロールになります。AWS Glue ジョブに使用するロールがない場合は、「[AWS Glue での IAM アクセス許可のセットアップ](#)」のガイドに従って、設定してください。

このロールをインタラクティブセッションのデフォルトのロールとして設定するには、以下の手順を実施します。

1. テキストエディタで `~/.aws/config` を開きます。
2. AWS Glue に使用するプロファイルを探します。いずれのプロファイルも使わない場合は、[Default] プロファイルを使用します。

3. `glue_role_arn=<AWSGlueServiceRole>` に示すように、使用するロールのプロファイルに 1 行を追加します。
4. [オプション]: プロファイルにデフォルトのリージョンセットがない場合は、`region=us-east-1` でリージョンを追加し、`us-east-1` を対象リージョンに置き換えることをお勧めします。
5. 設定を保存します。

詳細については、「[IAM を使用したインタラクティブセッション](#)」を参照してください。

インタラクティブセッションプレビューからのアップグレード

カーネルがバージョン 0.27 でリリースされたときに、新しい名前でもアップグレードされました。カーネルのプレビューバージョンをクリーンアップするには、PowerShellターミナルまたはから以下を実行します。

Note

カスタムサービスモデルを必要とする他の AWS Glue プレビューに参加している場合は、このカーネルを削除することで、対象のカスタムサービスモデルも削除されます。

```
# Remove Old Glue Kernels
jupyter kernelspec remove glue_python_kernel
jupyter kernelspec remove glue_scala_kernel

# Remove Custom Model
cd ~/.aws/models
rm -rf glue/
```

SageMaker Studio でインタラクティブセッションを使用する

AWS Glue インタラクティブセッションは、データサイエンティストおよびエンジニアがデータ準備および分析アプリケーションを迅速に構築、テスト、実行するために使用できる、オンデマンドのサーバーレス Apache Spark ランタイム環境です。Amazon SageMaker Studio Classic ノートブックを起動することにより、AWS Glue インタラクティブセッションを開始できます。

詳細については、「[AWS Glue インタラクティブセッションを使用してデータの準備](#)」を参照してください。

Microsoft Visual Studio Code によるインタラクティブセッションの使用

前提条件

- AWS Glue インタラクティブセッションをインストールし、Jupyter Notebook で動作することを確認します。
- Jupyter により Visual Studio Code をダウンロードし、インストールします。詳細については、「[Jupyter Notebooks in VS Code](#)」(VS コードでの Jupyter Notebook) を参照してください。

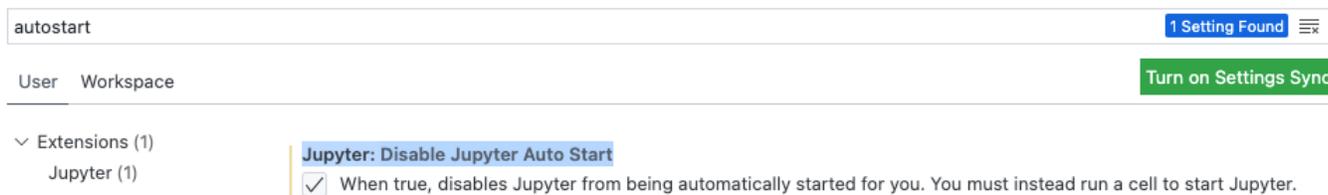
VSCode でインタラクティブセッションを始めるには

1. VS Code で Jupyter AutoStart を無効にします。

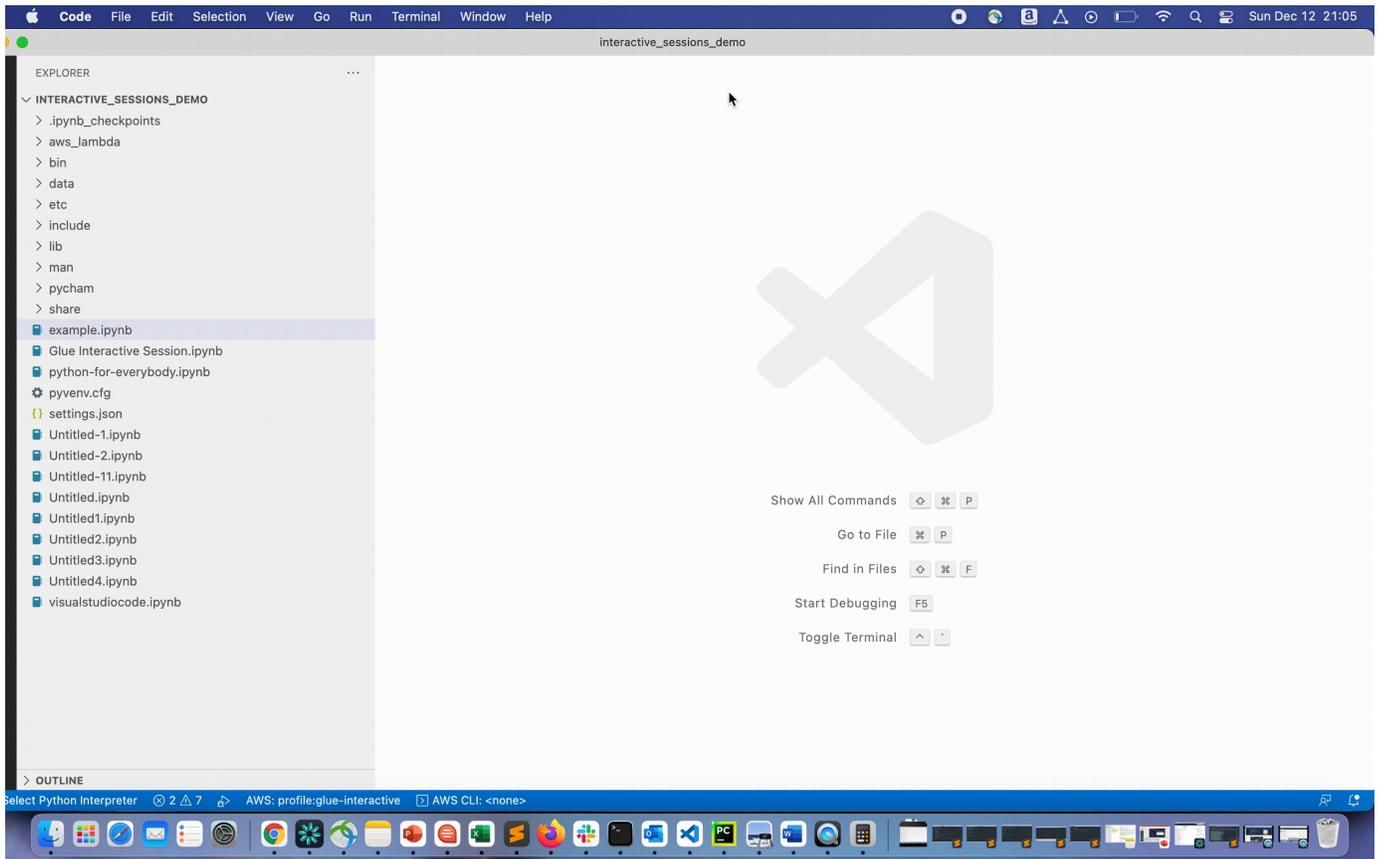
Visual Studio Code では、Jupyter カーネルが自動起動します。セッションが開始済みであるため、マジックは有効になりません。Windows で [自動起動] を無効にするには、[ファイル] > [設定] > [拡張機能] > [Jupyter] に移動し、Jupyter を右クリックして [拡張機能の設定] を選択します。

MacOS では、[コード] > [設定] > [拡張機能] > [Jupyter] に移動し、Jupyter を右クリックして [拡張機能の設定] を選択します。

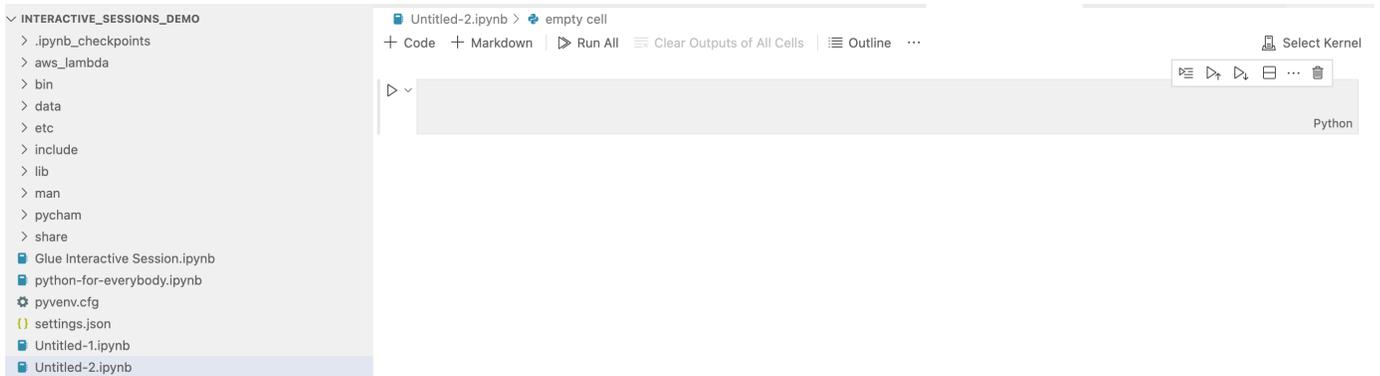
[Jupyter: 自動起動を無効にする] が表示されるまで下にスクロールします。[true の場合、Jupyter が自動的に起動されないようにする] のチェックボックスをオンにします。You must instead run a cell to start Jupyter.] (Jupyter を開始するには、セルを実行してください)



2. [File] (ファイル) > [New File] (新規ファイル) > [Save] (保存) の順に選択して、選択した名前でのこのファイルを .ipynb 拡張子として保存するか、または [Select a language] (言語の選択) の直下で [jupyter] を選択してファイルを保存します。



3. ファイルをダブルクリックします。Jupyter シェルが表示され、ノートブックが開きます。



4. Windows では、最初にファイルを作成するとき、デフォルトではカーネルが選択されていません。[Select Kernel] (カーネルの選択) をクリックすると、使用可能なカーネルのリストが表示されます。[Glue PySpark] を選択します。

MacOS で [Glue PySpark] カーネルが表示されない場合は、次の手順を試してください。

1. URL を取得するには、ローカルの Jupyter セッションを実行します。

たとえば、次のコマンドを実行して Jupyter Notebook をインストールします。

```
jupyter notebook
```

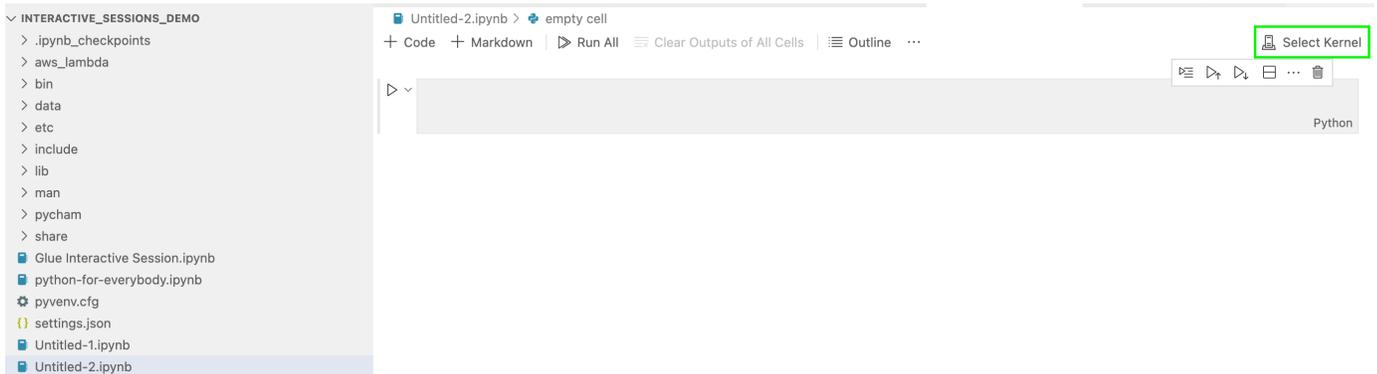
ノートブックを初めて実行すると、`http://localhost:8888/?token=3398XXXXXXXXXXXXXXXXXX` のような URL が表示されます。

URL をコピーします。

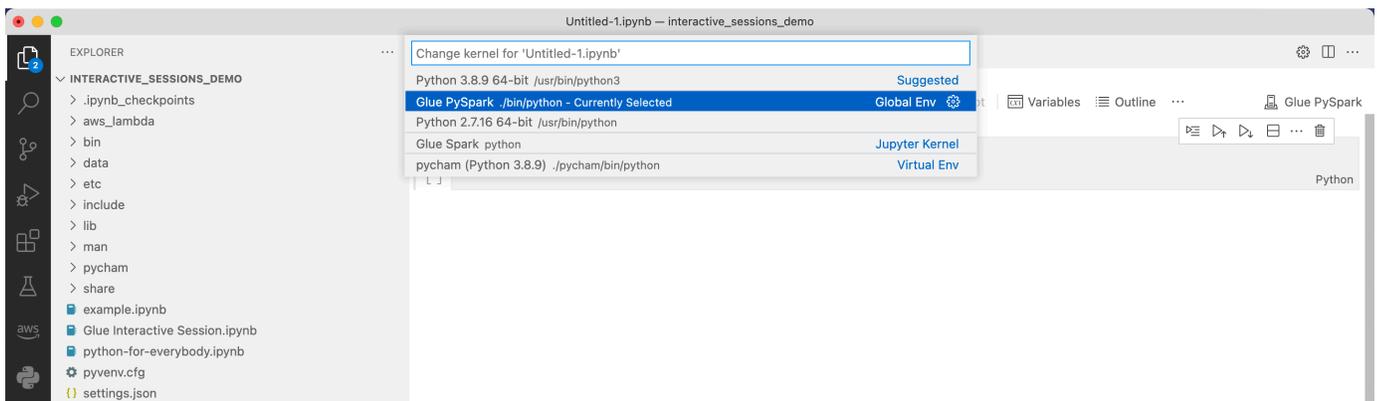
2. VS Code で、現在のカーネルをクリックし、[別のカーネルを選択...]、[既存の Jupyter Server...] を選択します。上記の手順でコピーした URL を貼り付けます。

エラーメッセージが表示される場合は、「[VS コードジュピターウィキ](#)」を参照してください。

3. 成功すると、カーネルが [Glue PySpark] に設定されます。



「Glue PySpark」または「Glue Spark」カーネルを (それぞれ Python と Scala 用に) 選択します。



ドロップダウンリストに [AWS Glue PySpark] および [AWS Glue Spark] カーネルが表示されない場合、上記のステップで AWS Glue カーネルがインストール済みであるか、Visual Studio

Code で `python.defaultInterpreterPath` 設定が正しいことを確認してください。詳細については、「[python.defaultInterpreterPath 設定の説明](#)」を参照してください。

5. AWS Glue のインタラクティブセッションを作成します。セッションの作成に進み、Jupyter Notebook と同じ方法を使用します。最初のセルの先頭に必要なマジックを指定し、コードのステートメントを実行します。

AWS Glue Jupyter および AWS Glue Studio ノートブックのインタラクティブセッションの設定

Jupyter Magics の概要

Jupyter Magics は、セル先頭またはセル本文全体として実行できるコマンドです。マジックの先頭は、Line Magics の場合は `%`、Cell Magics の場合は `%%` です。`%region` や `%connections` などの Line Magics は、セル内の複数のマジック、または次の例のようにセル本文に含まれるコードを使用して実行できます。

```
%region us-east-2
%connections my_rds_connection
dy_f = glue_context.create_dynamic_frame.from_catalog(database='rds_tables',
table_name='sales_table')
```

Cell Magics はセル全体を使用する必要があり、コマンドを複数行に渡すことができます。`%%sql` の例を以下に示します。

```
%%sql
select * from rds_tables.sales_table
```

Jupyter 向け AWS Glue インタラクティブセッションでサポートされているマジック

以下に、Jupyter Notebook の AWS Glue インタラクティブセッションで使用できるマジックを示します。

Sessions Magics

名前	型	説明
<code>%help</code>	該当なし	すべてのマジックコマンドのために説明と入力タイプのリストを返します。
<code>%profile</code>	文字列	認証情報プロバイダーとして使用するプロファイルを AWS 設定で指定します。
<code>%region</code>	文字列	セッションを初期化 AWS リージョンするを指定します。~/ <code>.aws/configure</code> からのデフォルト設定 例: <code>%region us-west-1</code>
<code>%idle_timeout</code>	Int	セルが実行されてからセッションがタイムアウトするまでの非アクティブ時間 (分)。Spark ETL セッションのデフォルトのアイドルタイムアウト値は、デフォルトのタイムアウトである 2,880分 (48時間) です。他のセッションタイプについては、該当するセッションタイプのマニュアルを参照してください。 例: <code>%idle_timeout 3000</code>
<code>%session_id</code>	該当なし	実行中のセッションのセッション ID をリターンします。
<code>%session_id_prefix</code>	文字列	<code>[session_id_prefix]-[session_id]</code> の形式ですべてのセッション ID に先行する文字列を定義します。セッション ID を指定しないと、ランダムな UUID が生成されます。AWS Glue Studioで Jupyter Notebook を実行する場合に、このマジックはサポートされていません。 例: <code>%session_id_prefix 001</code>

名前	型	説明
<code>%status</code>		現在の AWS Glue セッションのステータス (期間、構成、実行中のユーザー/ロールなど) を返します。
<code>%stop_session</code>		現在のセッションを停止します。
<code>%list_sessions</code>		名前と ID で、現在実行中のすべてのセッションをリストします。
<code>%session_type</code>	文字列	セッションタイプをストリーミング、ETL、または Ray のいずれかに設定します。 例: <code>%session_type Streaming</code>
<code>%glue_version</code>	文字列	このセッションで使用される AWS Glue のバージョン。 例: <code>%glue_version 3.0</code>

ジョブタイプを選択するためのマジック

名前	型	説明
<code>%streaming</code>	文字列	セッションタイプを AWS Glue ストリーミングに変更します。
<code>%etl</code>	文字列	セッションタイプを AWS Glue ETL に変更します。
<code>%glue_ray</code>	文字列	セッションタイプを AWS Glue for Ray に変更します。 AWS Glue 「Ray インタラクティブセッションでサポートされているマジック」 を参照してください。

AWS Glue for Spark 設定マジック

`%%configure`マジックは、セッションのすべての構成パラメータで構成される JSON 形式のディクショナリです。各パラメータは、ここで指定することも、個々のマジックを使って指定することもできます。

名前	型	説明
<code>%%configure</code>	辞書	<p>セッションのすべての構成パラメータで構成される JSON 形式のディクショナリを指定します。各パラメータは、ここで指定することも、個々のマジックを使って指定することもできます。</p> <p><code>%%configure</code> を使用する際のパラメータのリストと例については、表「<code>%%configure</code> の使用」を参照してください。</p>
<code>%iam_role</code>	文字列	<p>セッションを実行する IAM ロール ARN を指定します。<code>~/.aws/configure</code> からのデフォルト設定。</p> <p>例: <code>%iam_role AWSGlueServiceRole</code></p>
<code>%number_of_workers</code>	Int	<p>ジョブの実行時に割り当てられる、定義された <code>worker_type</code> のワーカーの数。<code>worker_type</code> も設定する必要があります。デフォルト <code>number_of_workers</code> は 5 です。</p> <p>例: <code>%number_of_workers 2</code></p>
<code>%additional_python_modules</code>	リスト	<p>クラスターに含める追加の Python モジュールのカンマ区切りリスト (PyPI または S3 からでも可)。</p> <p>例えば、<code>%additional_python_modules pandas, numpy</code> などです。</p>

名前	型	説明
%tags	文字列	<p>セッションにタグを追加します。タグを中かっこ {} で囲んで指定します。各タグ名のペアは二重引用符 (") で囲まれ、カンマ (,) で区切られます。</p> <pre data-bbox="894 443 1507 640">%tags {"billing":"Data-Platform", "team":"analytics"}</pre> <p>%status マジックを使用して、セッションに関連付けられたタグを表示します。</p> <pre data-bbox="894 800 1507 877">%status</pre> <pre data-bbox="894 911 1507 1822">Session ID: <sessionId> Status: READY Role: <example-role> CreatedOn: 2023-05-26 11:12:17.056000-07:00 GlueVersion: 3.0 Job Type: glueetl Tags: {'owner':'example-owner', 'team':'analytics', 'billing':'Data-Platform'} Worker Type: G.4X Number of Workers: 5 Region: us-west-2 Applying the following default arguments: --glue_kernel_version 0.38.0 --enable-glue-datacatalog true Arguments Passed: ['--glue_kernel_version: 0.38.0', '--enable-glue-datacatalog: true']</pre>

名前	型	説明
%%assume_role	辞書	<p>JSON 形式のディクショナリまたは IAM ロールの ARN 文字列を指定して、クロスアカウントアクセス用のセッションを作成します。</p> <p>ARN を使用した例:</p> <pre>%%assume_role { 'arn:aws:iam::XXXXXXXXXXXX: role/AWSGlueServiceRole'</pre> <p>認証情報の例:</p> <pre>%%assume_role {{ "aws_access_key_id" = "XXXXXXXXXXXX", "aws_secret_access_key" = "XXXXXXXXXXXX", "aws_session_token" = "XXXXXXXXXXXX" }}</pre>

%%configure セルマジック引数

%%configure マジックは、セッションのすべての構成パラメータで構成される JSON 形式のディクショナリです。各パラメータは、ここで指定することも、個々のマジックを使って指定することもできます。%%configure セルマジックがサポートする引数の例については、以下を参照してください。ジョブに指定された実行引数に -- プレフィックスを使用します。例：

```
%%configure
{
  "--user-jars-first": "true",
  "--enable-glue-datacatalog": "false"
```

```
}

```

ジョブパラメータの詳細については、「[ジョブのパラメータ](#)」を参照してください。

セッション設定

パラメータ	Type	説明
max_retries	Int	失敗した場合にこのジョブを再試行する最大回数。 <pre>%%configure { "max_retries": "0" }</pre>
max_concurrent_runs	Int	ジョブで許可される同時実行の最大数。 例： <pre>%%configure { "max_concurrent_runs": "3" }</pre>

セッションパラメータ

パラメータ	Type	説明
--enable-spark-ui	ブール値	Spark UI を有効にして AWS Glue ETL ジョブのモニタリングとデバッグを行います。 <pre>%%configure { "--enable-spark-ui": "true" }</pre>

パラメータ	Type	説明
		<pre>} </pre>
<code>--spark-event-logs-path</code>	文字列	<p>Amazon S3 パスを指定します。Spark UI のモニタリング機能を使用する場合。</p> <p>例 :</p> <pre>%%configure { "--spark-event-logs-path": "s3://path/to/event/logs/" }</pre>
<code>--script_location</code>	文字列	<p>ジョブを実行するスクリプトへの S3 パスを指定します。</p> <p>例 :</p> <pre>%%configure { "script_location": "s3://new- folder-here" }</pre>
<code>--SECURITY_CONFIGURATI</code> <code>ON</code>	文字列	<p>AWS Glue セキュリティ設定の名前</p> <p>例 :</p> <pre>%%configure { "--SECURITY_CONFIGURATI ON": <i>security-configura tion-name</i> , }</pre>

パラメータ	Type	説明
<code>--job-language</code>	文字列	<p>スクリプトプログラミング言語。'scala' または 'python' の値を使用できます。デフォルトは 'python' です。</p> <p>例 :</p> <pre>%%configure { "--job-language": "scala" }</pre>
<code>--class</code>	文字列	<p>Scala スクリプトのエントリポイントとなる Scala クラス。デフォルトは null です。</p> <p>例 :</p> <pre>%%configure { "--class": "className" }</pre>
<code>--user-jars-first</code>	ブール値	<p>クラスパスにあるお客様のその他の JAR ファイルに優先順位を付けます。デフォルトは null です。</p> <p>例 :</p> <pre>%%configure { "--user-jars-first": "true" }</pre>

パラメータ	Type	説明
<code>--use-postgres-driver</code>	ブール値	<p>JDBC ドライバーとの競合を避けるため、クラスパスで Postgres Amazon Redshift JDBC ドライバーを優先します。デフォルトは null です。</p> <p>例 :</p> <pre>%%configure { "--use-postgres-driver": "true" }</pre>
<code>--extra-files</code>	List(string)	<p>スクリプトを実行する前に、AWS Glue がスクリプトの作業ディレクトリにコピーする設定ファイルなどの追加ファイルを指す Amazon S3 のパス。</p> <p>例 :</p> <pre>%%configure { "--extra-files": "s3://path/to/ additional/files/" }</pre>

パラメータ	Type	説明
<code>--job-bookmark-option</code>	文字列	<p>ジョブブックマークの動作を制御します。 "、'job-bookmark-enable'、または 'job-bookmark-disable' の値を受け入れjob-bookmark-pauseます。デフォルトはjob-bookmark-disable「」です。</p> <p>例：</p> <pre>%%configure { "--job-bookmark-option": "job-bookmark-enable" }</pre>
<code>--TempDir</code>	文字列	<p>ジョブの一時ディレクトリとして使用できるバケットへの Amazon S3 のパスを指定します。デフォルトは null です。</p> <p>例：</p> <pre>%%configure { "--TempDir": "s3://path/to/temp/dir" }</pre>

パラメータ	Type	説明
<code>--enable-s3-parquet-optimized-commmitter</code>	ブール値	<p>Parquet データを Amazon S3 に書き込むために、EMRFS Amazon S3 最適化コミッターを有効にします。デフォルトは 'true' です。</p> <p>例 :</p> <pre>%%configure { "--enable-s3-parquet-optimized-commmitter": "false" }</pre>
<code>--enable-rename-algorithm-v2</code>	ブール値	<p>EMRFS 名前変更アルゴリズムのバージョンをバージョン 2 に設定します。デフォルトは 'true' です。</p> <p>例 :</p> <pre>%%configure { "--enable-rename-algorithm-v2": "true" }</pre>

パラメータ	Type	説明
<code>--enable-glue-data-catalog</code>	ブール値	<p>AWS Glue データカタログの Apache Spark Hive メタストアとしての使用を有効にします。</p> <p>例 :</p> <pre>%%configure { --"enable-glue-datacatalog": "true" }</pre>
<code>--enable-metrics</code>	ブール値	<p>ジョブの実行のジョブプロファイリングに関するメトリクスの収集を有効にします。デフォルトは 'false' です。</p> <p>例 :</p> <pre>%%configure { "--enable-metrics": "true" }</pre>
<code>--enable-continuous-cloudwatch-log</code>	ブール値	<p>AWS Glue ジョブのリアルタイムの連続ログ記録を有効にします。デフォルトは 'false' です。</p> <p>例 :</p> <pre>%%configure { "--enable-continuous-cloudw atch-log": "true" }</pre>

パラメータ	Type	説明
<code>--enable-continuous-log-filter</code>	ブール値	<p>連続ログ記録が有効であるジョブを作成または編集するときに、標準フィルタまたはフィルタなしを指定します。デフォルトは 'true' です。</p> <p>例 :</p> <pre>%%configure { "--enable-continuous-log-filter": "true" }</pre>
<code>--continuous-log-stream-prefix</code>	文字列	<p>連続 Amazon CloudWatch ログ記録が有効になっているジョブのカスタムログストリームプレフィックスを指定します。デフォルトは null です。</p> <p>例 :</p> <pre>%%configure { "--continuous-log-stream-prefix": "prefix" }</pre>

パラメータ	Type	説明
<code>--continuous-log-conversionPattern</code>	文字列	<p>連続ログ記録を有効にしたジョブのカスタム変換ログパターンを指定します。デフォルトは null です。</p> <p>例 :</p> <pre>%%configure { "--continuous-log-conversionPattern": "pattern" }</pre>
<code>--conf</code>	文字列	<p>Spark の設定パラメータを制御します。高度なユースケース向けです。各パラメータの前に <code>--conf</code> を使用してください。例 :</p> <pre>%%configure { "--conf": "spark.hadoop.hive.metastore.glue.catalogid=123456789012 --conf hive.metastore.client.factory.class=com.amazonaws.glue.catalog.metastore.AWSGlueDataCatalogHiveClientFactory --conf hive.metastore.schema.verification=false" }</pre>

Spark ジョブ (ETL とストリーミング) のマジック

名前	型	説明
<code>%worker_type</code>	文字列	Standard、G.1X、または G.2X を指定します。 <code>number_of_workers</code> も設定する必

名前	型	説明
		必要があります。デフォルトの worker_type は G.1X です。
%connections	リスト	セッションで使用する接続のカンマ区切りリストを指定します。 例： <pre>%connections my_rds_connection dy_f = glue_context.create_dynamic _frame.from_catalog(databas e='rds_tables', table_nam e='sales_table')</pre>
%extra_py_files	リスト	Amazon S3 からの追加の Python ファイルのカンマ区切りリスト。
%extra_jars	リスト	クラスターに含める追加の jar のカンマ区切りリスト
%spark_conf	文字列	セッション用のカスタム Spark 設定を指定します。例えば %spark_conf spark.serializer=org.apache.spark.serializer.KryoSerializer です。

Ray ジョブのマジック

名前	型	説明
%min_workers	Int	Ray ジョブに割り当てられるワーカーの最小数。デフォルト: 1。 例: %min_workers 2

名前	型	説明
<code>%object_memory_head</code>	Int	ウォームスタート後のインスタンスヘッドノードの空きメモリの割合。最小値: 0。最大値: 100。 例: <code>%object_memory_head 100</code>
<code>%object_memory_worker</code>	Int	ウォームスタート後のインスタンスワーカーノードの空きメモリの割合。最小値: 0。最大値: 100。 例: <code>%object_memory_worker 100</code>

Action Magics

名前	型	説明
<code>%%sql</code>	文字列	SQL コードを実行します。最初の <code>%%sql</code> マジックの後のすべての行は、SQL コードの一部として渡されます。 例: <code>%%sql select * from rds_tables.sales_table</code>
<code>%matplotlib</code>	Matplotlib Figure	Matplotlib ライブラリを使用してデータを可視化します。 例 : <pre>import matplotlib.pyplot as plt # Set X-axis and Y-axis values x = [5, 2, 8, 4, 9] y = [10, 4, 8, 5, 2] # Create a bar chart plt.bar(x, y)</pre>

名前	型	説明
		<pre># Show the plot %matplotlib plt</pre>
%plotly	Plotly Figure	<p>Plotly ライブラリを使用してデータを可視化します。</p> <p>例 :</p> <pre>import plotly.express as px #Create a graphical figure fig = px.line(x=["a","b","c"], y=[1,3,2], title="sample figure") #Show the figure %plotly fig</pre>

ネーミングセッション

AWS Glue インタラクティブセッションは AWS リソースであり、名前が必要です。名前はセッションごとに固有である必要があり、IAM 管理者によって制限される場合があります。詳細については、「[IAM を使用したインタラクティブセッション](#)」を参照してください。Jupyter カーネルは固有のセッション名を自動的に生成します。ただし、セッション名は、次の 2 つの方法で手動で名前を付けることができます。

1. にある AWS Command Line Interface 設定ファイルを使用します~.aws/config。 [「で AWS Config を設定する AWS Command Line Interface」](#) を参照してください。
2. %session_id_prefix マジックを使用する。 [Jupyter 向け AWS Glue インタラクティブセッションでサポートされているマジック](#) を参照してください。

セッション名は、次のように生成されます。

- プレフィックスと session_id が指定されている場合、セッション名は {prefix}-{UUID} になります。

- 何も指定されていない場合、セッション名は {UUID} になります。

セッション名のプレフィックスを付けると、セッションを AWS CLI または コンソールに一覧表示するときセッションを認識できます。

インタラクティブセッションの IAM ロールの指定

インタラクティブセッションで実行する ETL コードで使用する AWS Identity and Access Management (IAM) AWS Glue ロールを指定する必要があります。

ロールには、AWS Glue ジョブ実行に必要なアクセス許可と同じ IAM アクセス許可が付与されている必要があります。AWS Glue ジョブおよびインタラクティブセッションのロール作成の詳細については、「[AWS Glue 用に IAM ロールを作成する](#)」を参照してください。

IAM ロールは、次の 2 つの方法で指定できます。

- (`~/.aws/config` 推奨) にある AWS Command Line Interface 設定ファイルの使用。詳細については、「[~/.aws/config でのセッションの設定](#)」を参照してください。

Note

%profile マジックを使用すると、そのプロファイルの `glue_iam_role` の設定が優先されます。

- %iam_role マジックを使用する。(詳しくは、「[Jupyter 向け AWS Glue インタラクティブセッションでサポートされているマジック](#)」を参照してください。)

名前付きプロファイルを使用したセッションの設定

AWS Glue インタラクティブセッションは AWS Command Line Interface または boto3 と同じ認証情報を使用し、インタラクティブセッションは (Linux および MacOS) または `~/.aws/config` (`%USERPROFILE%\configWindows`) にある のような名前付きプロファイルを尊重して AWS CLI 使用します。MacOS 詳細については、「[名前を指定されたプロファイルを使用する](#)」を参照してください。

AWS Glue サービスロールとセッション ID プレフィックスをプロファイルで指定できるようにすることで、インタラクティブセッションが名前付きプロファイルを利用します。プロファイルロールを設定するには、以下に示すように、名前付きプロファイル

に `iam_role` キーまたは `session_id_prefix`、あるいはその両方の行を追加します。`session_id_prefix` は引用符を必要としません。例えば、`session_id_prefix` を追加する場合は、`session_id_prefix=myprefix` の値を入力します。

```
[default]
region=us-east-1
aws_access_key_id=AKIAIOSFODNN7EXAMPLE
aws_secret_access_key=wJalrXUtnFEMI/K7MDENG/bPxrFiCYEXAMPLEKEY
glue_iam_role=arn:aws:iam::<AccountID>:role/<GlueServiceRole>
session_id_prefix=<prefix_for_session_names>

[user1]
region=eu-west-1
aws_access_key_id=AKIAI44QH8DHBEXAMPLE
aws_secret_access_key=je7MtGbClwBF/2Zp9Utk/h3yCo8nvbEXAMPLEKEY
glue_iam_role=arn:aws:iam::<AccountID>:role/<GlueServiceRoleUser1>
session_id_prefix=<prefix_for_session_names_for_user1>
```

認証情報を生成するカスタムメソッドがある場合は、`~/.aws/config` ファイルの `credential_process` パラメーターを使用するようにプロファイルを設定することもできます。例:

```
[profile developer]
region=us-east-1
credential_process = "/Users/Dave/generate_my_credentials.sh" --username helen
```

`credential_process` パラメータを使用した認証情報の調達の詳細については、「[外部プロセスを使用した認証情報の調達](#)」を参照してください。

使用しているプロファイルにリージョンまたは `iam_role` が設定されていない場合は、最初に実行するセルの `%region` と `%iam_role` マジックを使用して設定する必要があります。

AWS Glue for Ray インタラクティブセッションの開始方法 (プレビュー)

Warning

AWS Glue for Ray インタラクティブセッションのプレビューは 2024 年 4 月 30 日に終了しました。AWS Glue for Ray で新しいインタラクティブセッションを作成できなくなります。

Note

AWS Glue for Ray は、米国東部 (バージニア北部)、米国東部 (オハイオ)、米国西部 (オレゴン)、アジアパシフィック (東京)、欧州 (アイルランド) で利用できます。

AWS Glue Studio コンソールでの Ray インタラクティブセッション

AWS Glue Studio コンソールのジョブページで、既存の Jupyter Notebook オプションを選択します。これにより、Notebook setup (ノートブックの設定) ページが開き、Kernel を選択できます。Ray カーネルを選択して、Ray のインタラクティブセッションを開始できます。インタラクティブセッションの詳細については、「[the section called “AWS Glue のインタラクティブセッションの開始方法”](#)」を参照してください。

AWS Glue Studio > Notebook setup

Notebook setup [Info](#)

Initial configuration

Job name

Enter a name for the job. This name will be used for the script and the notebook file.

IAM Role

Role assumed by the job with permission to access your data stores. Ensure that this role has permission to your Amazon S3 sources, targets, temporary directory, scripts, and any libraries used by the job.

Kernel

The kernel with which the notebook will be created.

[Cancel](#)[Start notebook](#)

Jupyter Kernel を使用した Ray のインタラクティブセッション

AWS Glue Studio コンソールの外部で Ray カーネルを使用するには、PyPI で公開されている `aws-glue-sessions` パッケージをインストールする必要があります。カーネルパッケージの使用についての詳細は、[the section called “AWS Glue のインタラクティブセッションの開始方法”](#) のドキュメントを参照してください。

カーネルを更新またはインストールするには、`pip install --upgrade aws-glue-sessions` を実行します。Ray カーネルを使用するには、バージョン `.37` 以上が必要です。

Ray のインタラクティブセッションでは、Ray ジョブと同じライブラリおよびバージョンの Ray にアクセスできます。プレビューでは、すべての Ray のインタラクティブセッションで Ray 2.4.0 を使用します。

Ray インタラクティブセッションのタイムアウトのデフォルト

- タイムアウト (セッション) のデフォルト: 8 時間

- アイドルタイムアウトのデフォルト: 1 時間

AWS Glue の Ray インタラクティブセッションでサポートされているマジック

AWS Glue Jupyter カーネルが Ray のインタラクティブセッションを動かすときのマジックは、Spark セッションのマジックと似ています。参考までに、「[the section called “AWS Glue Jupyter および AWS Glue Studio ノートブックのインタラクティブセッションの設定”](#)」を参照してください。

Sessions Magics

セッションのマジックは、AWS Glue for Ray のプレビュー前とほとんど同じです。このプレビュー以外のセッションのマジックに関する詳細は、「[the section called “Jupyter 向け AWS Glue インタラクティブセッションでサポートされているマジック”](#)」を参照してください。セッションタイプを AWS Glue for Ray に設定する新しいマジックを導入しました。

名前	型	説明
<code>%glue_ray</code>	文字列	セッションタイプを AWS Glue for Ray に変更します。

[AWS Glue Config Magics]

インタラクティブセッションで AWS Glue に設定するマジックは、セッションタイプによって異なる場合があります。現在、AWS Glue for Ray を使用する場合、既存のマジックの一部 (以下に示すもの) のみがサポートされています。

Warning

既知の問題: `additional_python_modules`

Ray のインタラクティブセッションのプレビューで `additional_python_modules` マジックを使用すると、ノートブックの保存時に問題が発生します。Ray セッション用の Python モジュールを設定するには、`%%configure` マジックを使用して、「[the section called “Ray ジョブのパラメータ”](#)」で定義されている `pip-install` パラメータを設定します。

名前	型	説明
<code>%%configure</code>	辞書	セッションのすべての構成パラメータで構成される JSON 形式のディクショナリを指定します。各パラメータは、ここで指定することも、個々のマジックを使って指定することもできます。
<code>%iam_role</code>	文字列	セッションを実行する IAM ロール ARN を指定します。~/aws/configure からのデフォルト
<code>%number_of_workers</code>	整数	ジョブの実行時に割り当てられる、定義された worker_type のワーカーの数。worker_type も設定する必要があります。
<code>%worker_type</code>	文字列	AWS Glue for Ray のプレビューでは、サポートされているワーカータイプは Z.2X だけです。
<code>%additional_python_modules</code>	リスト	クラスターに含める追加の Python モジュールのカンマ区切りリスト (Pypi または S3 からでも可)。

Action Magics

AWS Glue の Ray セッションは、アクションのマジックをサポートしていません。

IAM を使用したインタラクティブセッション

ここでは、AWS Glue インタラクティブセッションのセキュリティ上の考慮事項について説明します。

トピック

- [インタラクティブセッションで使用する IAM プリンシパル](#)
- [クライアントプリンシパルを設定する](#)

- [ランタイムロールを設定する](#)
- [以下の方法でセッションを非公開にします。 TagOnCreate](#)
- [IAM ポリシーに関する考慮事項](#)

インタラクティブセッションで使用する IAM プリンシパル

AWS Glue インタラクティブセッションでは、2 つの IAM プリンシパルを使用します。

- [Client principal] (クライアントプリンシパル): クライアントプリンシパル (ユーザーまたはロール) は、プリンシパルのアイデンティティベースの認証情報で構成された AWS Glue クライアントからのインタラクティブセッションの API オペレーションを認証します。例えば、AWS Glue コンソールにアクセスするために通常使用する IAM ロールのプリンシパルがこれに該当する場合があります。これは、認証情報がに使用される IAM のユーザーや AWS Command Line Interface、インタラクティブセッション Jupyter AWS Glue カーネルが使用するクライアントに付与されるロールでもかまいません。
- ランタイムロール: ランタイムロールは、クライアントプリンシパルがインタラクティブセッション API オペレーションに渡す IAM ロールです。AWS Glue がこのロールを使用して、セッションでステートメントを実行します。例えば、このロールは AWS Glue ETL ジョブの実行に使用される場合があります。

詳しくは、「[ランタイムロールを設定する](#)」を参照してください。

クライアントプリンシパルを設定する

インタラクティブセッション API を呼び出すには、アイデンティティポリシーをクライアントプリンシパルにアタッチする必要があります。このロールには、CreateSession などのインタラクティブセッション API に渡す実行ロールへの iam:PassRole アクセス権が必要です。たとえば、AWSGlueConsoleFullAccess 管理ポリシーを IAM ロールにアタッチすると、ポリシーがアタッチされたアカウントのユーザーが、アカウントで作成されたすべてのセッション (ランタイムステートメント、キャンセルステートメントなど) にアクセスできるようになります。

セッションを保護して、セッションを作成したユーザーに関連するものなど、特定の IAM ロールのみ限定したい場合は、AWS Glue Interactive Session のタグベース承認制御と呼ばれるタグベース承認コントロールを使用できます。TagOnCreate 詳細については、[以下の方法でセッションを非公開にします。 TagOnCreate](#) オーナータグベースのスコープダウン管理ポリシーでセッションをプライベートにする方法をご覧ください。TagOnCreate アイデンティティベースポリシーの詳細については、「[AWS Glue のアイデンティティベースポリシー](#)」を参照してください。

ランタイムロールを設定する

AWS Glue インタラクティブセッションでステートメントを引き受けて実行できるようにするには、IAM ロールを `CreateSession` API オペレーションに渡す必要があります。ロールには、通常の AWS Glue ジョブの実行に必要なアクセス許可と同じ IAM アクセス許可が必要です。たとえば、`AWSGlueServiceRole` AWS Glue AWS ユーザーに代わってサービスを呼び出すことを許可するポリシーを使用してサービスロールを作成できます。AWS Glue コンソールを使用すると、ユーザーに代わってサービスロールが自動的に作成されるか、既存のものを使用します。独自の IAM ロールを作成し、独自の IAM ポリシーをアタッチして、同様の許可を許可することもできます。

セッションを保護して、そのセッションを作成したユーザーのみに公開したい場合は、AWS Glue Interactive Session のタグベース認証制御と呼ばれるものを使用できます `TagOnCreate`。詳しくは、[以下の方法でセッションを非公開にします。TagOnCreate](#) 所有者タグベースのスコープダウン管理ポリシーでセッションを非公開にする方法をご覧ください。 `TagOnCreate` アイデンティティベースポリシーの詳細については、「[AWS Glue のアイデンティティベースのポリシー](#)」を参照してください。IAM コンソールから実行ロールを自分で作成していて、`TagOnCreate` サービスを機能付きでプライベートにしたい場合は、以下の手順に従ってください。

1. ロールタイプを `Glue` に設定して IAM ロールを作成します。
2. AWS Glue この管理ポリシーを添付してください。 `AwsGlueSessionUserRestrictedServiceRole`
3. `AwsGlueSessionUserRestrictedServiceRole` ロール名の前にポリシー名を付けます。たとえば、 `AwsGlueSessionUserRestrictedServiceRole-myrole` という名前のロールを作成し、AWS Glue 管理ポリシーをアタッチできます。 `AwsGlueSessionUserRestrictedServiceRole`
4. 次のような信頼ポリシーをアタッチして、AWS Glue がロールを引き受けることを許可します。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": [
          "glue.amazonaws.com"
        ]
      },
      "Action": [
        "sts:AssumeRole"
      ]
    }
  ]
}
```

```
]
}
```

インタラクティブセッション Jupyter カーネルの場合は、プロファイルでキーを指定できます。iam_role AWS Command Line Interface 詳細については、「[~/.aws/config でのセッションの設定](#)」を参照してください。AWS Glue ノートブックを使用してインタラクティブセッションと対話している場合は、実行する最初のセルで、%iam_role マジックの実行ロールを渡すことができます。

以下の方法でセッションを非公開にします。 TagOnCreate

AWS Glue インタラクティブセッション は、名前付きリソースとして、インタラクティブセッション向けにタグ付けとタグベースの認証 (TBAC) をサポートしています。TBAC TagResource を使用した UntagResource API に加えて、AWS Glueインタラクティブセッションでは、操作を伴うセッションの作成中にのみ、セッションに特定のタグを「タグ付け」 TagOnCreate する機能がサポートされています。CreateSession つまり、aka DeleteSession ではそれらのタグが削除されるということです。 UntagOnDelete

TagOnCreate セッションをセッションの作成者に非公開にする強力なセキュリティメカニズムを提供します。たとえば、「owner」 RequestTag と $\{aws: UserID\}$ の値を持つ IAM ポリシーをクライアントプリンシパル (ユーザーなど) にアタッチして、呼び出し元の UserID の値と一致する「owner」タグがリクエストで userId タグとして提供された場合にのみセッションを作成できるようにすることができます。CreateSession このポリシーにより、AWS Glue インタラクティブセッションはセッションリソースを作成し、セッションの作成時にのみ、userId タグでセッションにタグ付けすることができます。それに加えて、実行中に渡した実行ロールに「owner」 ResourceTag の付いたIAMポリシーをアタッチすることで、セッションへのアクセス (ステートメントの実行など) をセッションの作成者 (つまり $\{aws: UserId\}$ という値の所有者タグ) のみに限定できます CreateSession。

TagOnCreate セッションをセッション作成者のみに公開する機能をより使いやすくするために、AWS Glue専用の管理ポリシーとサービスロールが用意されています。

IAM AssumeRole プリンシパルを使用して (つまり、IAM ロールを引き受けた認証情報を使用して) AWS Glue インタラクティブセッションを作成し、そのセッションを作成者専用にした場合は、それぞれおよび同様のポリシーを使用してください。AWSGlueSessionUserRestrictedNotebookPolicyAWSGlueSessionUserRestrictedNotebookServiceRoleこれらのポリシーによりAWS Glue、 $\{aws: PrincipalTag\}$ を使用して所有者のタグ値を抽出できます。これには、引き受けロール認証情報と同様に SessionTag 、値が $\{aws: UserId\}$ の UseruserId タグを渡す必要があります。「[ID session tags](#)」 (ID セッションタグの受け渡し) を参照してください

さい。認証情報を販売するインスタンスプロファイルで Amazon EC2 インスタンスを使用している、Amazon EC2 インスタンス内からセッションを作成したり、セッションを操作したりする場合は、割り当てロール認証情報と同様に SessionTag \$ {aws: userId} という値の UserID タグを渡す必要があります。

たとえば、IAM AssumeRole プリンシパル認証情報を使用してセッションを作成していて、TagOnCreate サービスを機能付きでプライベートにしたい場合は、以下の手順に従ってください。

1. IAM コンソールからランタイムロールを自分で作成します。AWS GlueAwsGlueSessionUserRestrictedNotebookServiceRoleこの管理ポリシーを添付し、ロール名の前にポリシー名を付けてください。AwsGlueSessionUserRestrictedNotebookServiceRoleたとえば、AwsGlueSessionUserRestrictedNotebookServiceRole-myrole という名前のロールを作成し、AWS Glue管理ポリシーをアタッチできます。AwsGlueSessionUserRestrictedNotebookServiceRole
2. AWS Glue が上記のロールを引き受けることができるように、以下のような信頼ポリシーを適用します。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": [
          "glue.amazonaws.com"
        ]
      },
      "Action": [
        "sts:AssumeRole"
      ]
    }
  ]
}
```

3. プレフィックスを付けた別のロールを作成しAwsGlueSessionUserRestrictedNotebookPolicy、AWS GlueAwsGlueSessionUserRestrictedNotebookPolicy管理ポリシーをアタッチしてセッションをプライベートにします。管理ポリシーに加えて、ステップ 1 で作成したロールに iam: PassRole を許可する以下のインラインポリシーをアタッチしてください。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "iam:PassRole"
      ],
      "Resource": [
        "arn:aws:iam::*:role/
        AwsGlueSessionUserRestrictedNotebookServiceRole*"
      ],
      "Condition": {
        "StringLike": {
          "iam:PassedToService": [
            "glue.amazonaws.com"
          ]
        }
      }
    }
  ]
}
```

4. 上記の IAM AWS Glue に次のような信頼ポリシーをアタッチして、ロールを引き受けます。

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Effect": "Allow",
    "Principal": {
      "Service": [
        "glue.amazonaws.com"
      ]
    },
    "Action": [
      "sts:AssumeRole",
      "sts:TagSession"
    ]
  }]
}
```

Note

オプションで、1つのロール(ノートブックロールなど)を使用して、`AwsGlueSessionUserRestrictedNotebookServiceRole` `AwsGlueSessionUserRestrictedNotebookServiceRole` の管理ポリシーとの両方をアタッチできます。また、追加のインラインポリシーをアタッチして、自分のロールの `iam:passrole` を AWS Glue に許可します。最後に、上記の信頼ポリシーをアタッチして、`sts:AssumeRole` と `sts:TagSession` を許可します。

AWSGlueSessionUserRestrictedNotebookPolicy

`AWSGlueSessionUserRestrictedNotebookPolicy` は、タグキーが「owner」で、値がプリンシパル(ユーザーまたはロール) AWS のユーザー ID と一致する場合のみ、AWS Glue ノートブックからインタラクティブセッションを作成できます。詳細については、「[ポリシー変数を使用する場所](#)」を参照してください。このポリシーは、AWS Glue Studio から AWS Glue インタラクティブセッション ノートブックを作成するプリンシパル(ユーザーまたはロール) にアタッチされます。また、このポリシーでは、AWS プリンシパルのユーザー ID と一致する「owner」AWS Glue Studio AWS Glue Studio タグ値で作成されたインタラクティブセッションリソースを操作するためのノートブックへの十分なアクセスも許可されます。このポリシーは、セッションが作成された後に、AWS Glue セッションリソースから「owner」タグを変更または削除する許可を拒否します。

AWSGlueSessionUserRestrictedNotebookServiceRole

は、AWS Glue Studio ノートブック作成者のプリンシパル(ユーザーまたはロール) AWS のユーザー ID と一致する「owner」`AWSGlueSessionUserRestrictedNotebookServiceRole` AWS Glue タグ値で作成されたインタラクティブセッションリソースとやり取りするためのノートブックへの十分なアクセスを提供します。詳細については、「[ポリシー変数を使用する場所](#)」を参照してください。このサービスロールポリシーは、ノートブックにマジックロールとして渡されるロール、または実行ロールとして API に渡されるロールにアタッチされます。CreateSession このポリシーでは、タグキーが「owner」で、AWS 値がプリンシパルのユーザー ID と一致する場合にのみ、AWS Glue ノートブックからインタラクティブセッションを作成することも許可されます。このポリシーは、セッションが作成された後に、AWS Glue セッションリソースから「owner」タグを変更または削除する許可を拒否します。このポリシーには、Amazon S3 バケットへの書き込みと読み取り、CloudWatch ログの書き込み、AWS Glue が使用する Amazon EC2 リソースのタグの作成と削除を行う権限も含まれています。

ユーザーポリシーを使用してセッションをプライベートにする

をアカウント内の各ユーザーにアタッチされている IAM ロールにアタッチして、ユーザー自身の `{aws: UserId}` と一致する値の所有者タグのみを使用してセッションを作成できないように制限できます。AWSGlueSessionUserRestrictedPolicyとを使用する代わりに、AWSGlueSessionUserRestrictedNotebookPolicyAWSGlueSessionUserRestrictedNotebookServiceRoleそれぞれと同様のポリシーを使用する必要があります。AWSGlueSessionUserRestrictedPolicyAWSGlueSessionUserRestrictedServiceRole詳細については、「[アイデンティティベースポリシーを使用する](#)」を参照してください。このポリシーは、作成者、つまり自分の `{aws:userId}` を持つ owner タグでセッションを作成したユーザーの `{aws:userId}` のみに、セッションへのアクセス範囲を絞り込みます。の手順に従って IAM コンソールを使用して自分で実行ロールを作成した場合は[ランタイムロールを設定する](#)、AwsGlueSessionUserRestrictedPolicy管理ポリシーをアタッチするだけでなく、次のインラインポリシーをアカウント内の各ユーザーにアタッチして、iam:PassRole前に作成した実行ロールを許可します。

```
{
  "Version": "2012-10-17",
  "Statement": [{
    "Effect": "Allow",
    "Action": [
      "iam:PassRole"
    ],
    "Resource": [
      "arn:aws:iam::*:role/AwsGlueSessionUserRestrictedServiceRole*"
    ],
    "Condition": {
      "StringLike": {
        "iam:PassedToService": [
          "glue.amazonaws.com"
        ]
      }
    }
  ]
}
```

AWSGlueSessionUserRestrictedPolicy

は、タグキー「owner」AWS とユーザー ID と一致する値が提供されている場合のみ、CreateSession API AWSGlueSessionUserRestrictedPolicyAWS Glueを使用してインタラクティブ

セッションを作成するためのアクセスを提供します。この ID ポリシーは API を呼び出すユーザーにアタッチされます `CreateSession`。このポリシーでは、ユーザー ID と一致する「owner」AWS Glue タグと値を使用して作成されたインタラクティブセッションリソースとのやり取りも許可されます。AWS このポリシーは、セッションが作成された後に、AWS Glue セッションリソースから「owner」タグを変更または削除する許可を拒否します。

`AWSGlueSessionUserRestrictedServiceRole`

は、`AWSGlueSessionUserRestrictedServiceRole` AWS Glue セッションを除くすべてのリソースへのフルアクセスを提供し、ユーザーに関連付けられているインタラクティブセッションのみを作成して使用することを許可します。このポリシーには、AWS 他のサービスの Glue AWS Glue リソースを管理するために必要なその他の権限も含まれています。このポリシーでは、AWS Glue AWS 他のサービスのリソースにタグを追加することもできます。

IAM ポリシーに関する考慮事項

インタラクティブセッションは AWS Glue の IAM リソースです。これらは IAM リソースであるため、セッションへのアクセスとインタラクションは IAM ポリシーによって管理されます。管理者によって設定されたクライアントプリンシパルまたは実行ロールに添付された IAM ポリシーに基づいて、クライアントプリンシパル (ユーザーまたはロール) は新しいセッションを作成し、独自のセッションや他のセッションと対話することができます。

管理者が、`AWSGlueConsoleFullAccess` `AWSGlueServiceRole` AWS Glue そのアカウントのすべてのリソースへのアクセスを許可するやなどの IAM ポリシーをアタッチした場合、クライアントプリンシパルは相互に連携できます。例えば、ポリシーで許可されていれば、ユーザーは他のユーザーによって作成されたセッションと対話することができます。

特定のニーズに合わせたポリシーを設定したい場合は、[ポリシーのリソースの設定に関する IAM ドキュメント](#)を参照してください。たとえば、ユーザーに属するセッションを分離するには、`TagOnCreate` AWS Glue インタラクティブセッションでサポートされている機能を使用できます。[以下の方法でセッションを非公開にします。](#) `TagOnCreate` を参照してください。

インタラクティブセッションは、特定の VPC 条件に基づくセッション作成の制限をサポートします。「[条件キーを使って設定を制御するポリシーを制御する](#)」を参照してください。

スクリプトまたはノートブックの AWS Glue ジョブジョブへの変換

スクリプトまたはノートブックを AWS Glue ジョブに変換する方法は2つあります。

- nbconvertを使用して、Jupyter .ipynb Notebook ドキュメントファイルを .py ファイルに変換します。詳細については、「[nbconvert: Convert Notebooks to other formats](#)」(nbconvert : ノートブックを他の形式に変換する) を参照してください。
- ファイルを AWS Glue Studio ノートブックにアップロードします。
 - AWS Glue Studio コンソールで、ナビゲーションメニューから ジョブを選択します。
 - ジョブの作成セクションで、Jupyter Notebookを選択します。
 - オプションセクションで、[Upload] (アップロード) を選択し、既存のノートブックを編集します。
 - Choose file (ファイルの選択)を選択して、.ipynb ファイルをアップロードします。

AWS Glue ストリーミングのためのインタラクティブセッション

ストリーミングセッションタイプの切り替え

AWS Glue インタラクティブセッション設定マジックである `%streaming` を使用して、実行しているジョブを定義し、ストリーミングインタラクティブセッションを初期化します。

インタラクティブな開発のためのサンプリング入カストリーム

AWS Glueインタラクティブセッションでのインタラクティブ体験を向上させるために開発したツールの1つは、`GlueContext DynamicFrame`静的なストリームのスナップショットを取得する新しいメソッドを新たに追加したことです。GlueContextワークフローを検証、操作、実装することができます。

GlueContext クラスインスタンスを使用して、メソッド `getSampleStreamingDynamicFrame` を見つけることができます。このメソッドに必要な引数は以下のとおりです。

- `dataFrame`: Spark ストリーミング DataFrame
- `options`: 以下の利用可能なオプションを参照してください。

以下のオプションを使用できます。

- `windowSize`: これは Microbatch Duration と呼ばれます。このパラメータは、前回のバッチがトリガーされてからストリーミングクエリが待機する時間を特定します。このパラメータの値は、`pollingTimeInMs` より小さくしてください。

- `pollingTimeInMs`: メソッドが実行される合計時間。入カストリームからサンプルレコードを取得するために、少なくとも1つのマイクロバッチが起動されます。
- `recordPollingLimit`: このパラメータは、ストリームからポーリングするレコードの総数を制限するのに役立ちます。
- (オプション) `writeStreamFunction` を使用して、このカスタム機能をすべてのレコードサンプリング機能に適用することもできます。以下は、Scala と Python のサンプルコードです。

Scala

```
val sampleBatchFunction = (batchDF: DataFrame, batchId: Long) => {//Optional but you can replace your own forEachBatch function here
val jsonString: String = s""""{"pollingTimeInMs": "10000", "windowSize": "5 seconds"}""""
val dynFrame = glueContext.getSampleStreamingDynamicFrame(YOUR_STREAMING_DF,
  JsonOptions(jsonString), sampleBatchFunction)
dynFrame.show()
```

Python

```
def sample_batch_function(batch_df, batch_id):
    //Optional but you can replace your own forEachBatch function here
    options = {
        "pollingTimeInMs": "10000",
        "windowSize": "5 seconds",
    }
    glue_context.getSampleStreamingDynamicFrame(YOUR_STREAMING_DF, options,
        sample_batch_function)
```

Note

サンプリングされた `DynFrame` が空の場合、いくつかの原因が考えられます。

- ストリーミングソースが「Latest」(最新) に設定されており、サンプリング期間中に新しいデータは取り込まれていません。

- ポーリング時間が十分ではなく、取り込まれたレコードを処理できません。バッチ全体が処理されない限り、データは表示されません。

インタラクティブセッションでのストリーミングアプリケーションの実行

AWS Glue インタラクティブセッションでは、AWS Glue コンソールでストリーミングアプリケーションを作成するのと同じように、AWS Glue ストリーミングアプリケーションを実行できます。インタラクティブセッションはセッションベースのため、ランタイムで例外が発生してもセッションが停止することはありません。バッチ関数を反復的に開発するというメリットが追加されました。例:

```
def batch_function(data_frame, batch_id):
    log.info(data_frame.count())
    invalid_method_call()
glueContext.forEachBatch(frame=streaming_df, batch_function = batch_function, options =
{**})
```

上記の例にはメソッドの無効な使用が含まれており、アプリケーション全体が終了する通常の AWS Glue ジョブとは異なります。ユーザーのコーディングコンテキストと定義は完全に保持され、セッションは引き続き動作可能な状態になっています。新しいクラスターをブートストラップして、先行するすべての変換を再実行する必要はありません。これにより、バッチ関数の実装をすばやくに反復でき、望ましい結果を得ることができます。

インタラクティブセッションは、セッションが一度に 1 つのステートメントだけを実行するように、各ステートメントをブロックする方法で評価することに注意してください。ストリーミングクエリは継続的で終了することがないため、アクティブなストリーミングクエリを含むセッションは、中断されない限り後続のステートメントを処理できません。Jupyter Notebook から直接中断のコマンドを発行すると、カーネルがキャンセルを処理します。

例として、実行を待機するシーケンシャルなステートメントを考えてみましょう。

```
Statement 1:
    val number = df.count()
    #Spark Action with deterministic result
    Result: 5

Statement 2:
    streamingQuery.start().awaitTermination()
```

```
#Spark Streaming Query that will be executing continuously  
Result: Constantly updated with each microbatch
```

Statement 3:

```
val number2 = df.count()  
#This will not be executed as previous statement will be running indefinitely
```

AWS Glueスクリプトのローカルでの開発およびテスト

Spark ジョブスクリプト用の AWS Glue の開発およびテストする場合、以下のような複数のオプションが用意されています。

- AWS Glue Studio コンソール
 - Visual editor (ビジュアルエディタ)
 - スクリプトエディタ
 - AWS Glue Studio ノートブック
- インタラクティブなセッション
 - Jupyter Notebook
- Docker イメージ
 - ローカル開発
 - リモート開発
- AWS Glue Studio ETL ライブラリ
 - ローカル開発

ユーザーは、要件に基づいて、上記のオプションのいずれかを選択できます。

コードを記述したくない、あるいは記述する量を削減したい場合は、AWS Glue Studio のビジュアルエディタが便利です。

ノートブックのインタラクティブな操作をご希望の場合には、AWS Glue Studio のノートブックが適しています。詳細については、「[AWS Glue Studio と AWS Glue でのノートブックの使用](#)」を参照してください。独自のローカル環境を使用するユーザーには、インタラクティブセッションが適しています。詳細については、「[AWS Glue でインタラクティブセッションを使用する](#)」を参照してください。

開発をローカル/リモートの両方で行いたい場合は、Docker イメージが便利です。これにより、Spark ジョブスクリプト用の AWS Glue の開発とテストが好きな場所で行うことができ、AWS Glue のコストもかかりません。

Docker を使用せずローカルで開発を行いたい場合は、AWS Glue ETL ライブラリディレクトリを、ローカルにインストールします。

を使用した開発AWS Glue Studio

AWS Glue Studio のビジュアルエディターは、AWS Glue での抽出、変換、ロード (ETL) ジョブの作成、実行、およびモニタリングが簡単に行えるグラフィカルなインターフェイスです。データ変換ワークフローを視覚的に作成し、AWS Glue の Apache Spark ベースのサーバーレス ETL エンジン上で、それらをシームレスに実行することができます。ジョブの各ステップでスキーマとデータの結果を調べることができます。詳細については、[AWS Glue Studio ユーザーガイド](#)を参照してください。

インタラクティブセッションによる開発

インタラクティブセッションでは、お好みの環境からアプリケーションを構築し、そのテストが行えます。詳細については、「[AWS Glue でインタラクティブセッションを使用する](#)」を参照してください。

Docker イメージによる開発

Note

このセクションでの手順は、Microsoft Windows オペレーティングシステムではテストされていません。

Windows プラットフォームでのローカル開発およびテストについては、ブログ記事「[Building an AWS Glue ETL pipeline locally without an AWS account](#)」を参照してください。

本番環境対応のデータプラットフォームの場合、AWS Glue ジョブ向けの開発プロセスや CI/CD パイプラインで検討すべきです。Docker コンテナでは、AWS Glue のジョブの開発やテストが柔軟に行えます。AWS Glue は Docker Hub に Docker イメージをホストすることで、追加のユーティリティを使用する開発環境をセットアップします。IDE、ノートブック、または、AWS Glue ETL ライブラリを使用する REPL などを、ご自分で選択して使用できます。このトピックでは、Docker コン

テナ上で Docker イメージを使用して、AWS Glue ジョブのバージョン 4.0 を開発およびテストする方法について説明します。

Docker Hub には、以下のように AWS Glue で利用可能な Docker イメージが用意されています。

- AWS Glue バージョン 4.0 の場合: `amazon/aws-glue-libs:glue_libs_4.0.0_image_01`
- AWS Glue バージョン 3.0 の場合: `amazon/aws-glue-libs:glue_libs_3.0.0_image_01`
- AWS Glue バージョン 2.0 の場合: `amazon/aws-glue-libs:glue_libs_2.0.0_image_01`

これらのイメージは x86_64 用です。このアーキテクチャでテストすることをお勧めします。ただし、サポートされていないベースイメージでローカル開発ソリューションを作り直すことは可能な場合があります。

この例では、ローカルマシンでの `amazon/aws-glue-libs:glue_libs_4.0.0_image_01` の使用方法とコンテナの実行方法について説明します。このコンテナイメージは、AWS Glue Spark ジョブのバージョン 3.3 でテスト済みです。このグループには、以下が含まれます。

- Amazon Linux
- AWS Glue ETL ライブラリ ([aws-glue-libs](#))
- Apache Spark 3.3.0
- Spark 履歴サーバー
- Jupyter ラボ
- Livy
- その他のライブラリ依存関係 (AWS Glue ジョブシステムのライブラリセットと同じです)

要件に応じて、次のいずれかのセクションを完了します。

- `spark-submit` を使用するようにコンテナをセットアップする
- REPL シェル (PySpark) を使用するようにコンテナをセットアップする
- Pytest を使用するようにコンテナをセットアップする
- Jupyter Lab を使用するようにコンテナをセットアップする
- Visual Studio Code を使用するようにコンテナをセットアップする

前提条件

作業を開始する前に、Docker がインストール済みであり、Docker デーモンが実行中であることを確認します。インストールの手順については、[Mac](#) または [Linux](#) 向けの Docker ドキュメントを参照してください。AWS Glue コンテナは、Docker を実行しているマシンでホストされます。また、Docker を実行しているホスト上のイメージ用として、少なくとも 7 GB のディスク領域が必要です。

AWS Glue のコードをローカルで開発する際の制限については、「[ローカル開発の制限](#)」を参照してください。

AWS の設定

コンテナから AWS API コールを有効にするには、以下の手順に従って AWS 認証情報をセットアップします。以下のセクションでは、この AWS により名前が指定されたプロファイルを使用します。

1. AWS CLI を設定し、名前付きプロファイルを設定します。AWS CLI 構成の詳細については、AWS CLI ドキュメントの「[構成と認証情報ファイルの設定](#)」を参照してください。
2. ターミナルで次のコマンドを実行します。

```
PROFILE_NAME="<your_profile_name>"
```

リクエストの送信先となる AWS リージョンを指定するとき、AWS_REGION 環境変数の設定も必要になる場合があります。

コンテナの設定と実行

spark-submit コマンドから PySpark コードを実行できるようにコンテナを設定するには、以下のような高度な手順が必要です。

1. Docker Hub からイメージを取得します。
2. コンテナを実行します。

Docker Hub からのイメージの取得

Docker Hub からイメージを取得するには、次のコマンドを実行します。

```
docker pull amazon/aws-glue-libs:glue_libs_4.0.0_image_01
```

コンテナの実行

ここまでの操作により、取得したイメージを使用してコンテナを実行できるようになります。ご自身の要件に応じて、以下のいずれかを選択できます。

spark-submit

コンテナで `spark-submit` コマンドを使用して、AWS Glue のジョブスクリプトを実行できます。

1. スクリプトを記述し、`/local_path_to_workspace` ディレクトリの下に `sample1.py` として保存します。このトピックの付録に、このためのサンプルコードが含まれています。

```
$ WORKSPACE_LOCATION=/local_path_to_workspace
$ SCRIPT_FILE_NAME=sample.py
$ mkdir -p ${WORKSPACE_LOCATION}/src
$ vim ${WORKSPACE_LOCATION}/src/${SCRIPT_FILE_NAME}
```

2. 以下のコマンドにより、コンテナで `spark-submit` コマンドを使用して、新しい Spark アプリケーションをサブミットします。

```
$ docker run -it -v ~/.aws:/home/glue_user/.aws -v $WORKSPACE_LOCATION:/
home/glue_user/workspace/ -e AWS_PROFILE=$PROFILE_NAME -e DISABLE_SSL=true
--rm -p 4040:4040 -p 18080:18080 --name glue_spark_submit amazon/aws-glue-
libs:glue_libs_4.0.0_image_01 spark-submit /home/glue_user/workspace/src/
$SCRIPT_FILE_NAME
...22/01/26 09:08:55 INFO DAGScheduler: Job 0 finished: fromRDD at
DynamicFrame.scala:305, took 3.639886 s
root
|-- family_name: string
|-- name: string
|-- links: array
| |-- element: struct
| | |-- note: string
| | |-- url: string
|-- gender: string
|-- image: string
|-- identifiers: array
| |-- element: struct
| | |-- scheme: string
| | |-- identifier: string
|-- other_names: array
| |-- element: struct
| | |-- lang: string
```

```

| | |-- note: string
| | |-- name: string
|-- sort_name: string
|-- images: array
| |-- element: struct
| | |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
|-- contact_details: array
| |-- element: struct
| | |-- type: string
| | |-- value: string
|-- death_date: string

...

```

3. (オプション) 環境に合わせて `spark-submit` を構成します。たとえば、`--jars` 設定を使用して依存関係を渡すことができます。詳細については、Spark ドキュメントの [spark-submit を使用したアプリケーションの起動](#) を参照してください。

REPL シェル (Pyspark)

REPL (read-eval-print loops) シェルを実行すると、インタラクティブな開発を行うことができます。

次のコマンドにより、コンテナ上で PySpark コマンドを実行し、REPL シェルを起動します。

```

$ docker run -it -v ~/.aws:/home/glue_user/.aws -e AWS_PROFILE=$PROFILE_NAME -e
  DISABLE_SSL=true --rm -p 4040:4040 -p 18080:18080 --name glue_pyspark amazon/aws-glue-
  libs:glue_libs_4.0.0_image_01 pyspark

```

```
...
```

```

  ____
 /  _/  _  _  _  _  _  /  _
_ \  \  _  \  _  \  _  \  _
/  _/  .  _  \  _  \  _  \  _  version 3.1.1-amzn-0
/_/

```

```

Using Python version 3.7.10 (default, Jun 3 2021 00:02:01)
Spark context Web UI available at http://56e99d000c99:4040
Spark context available as 'sc' (master = local[*], app id = local-1643011860812).
SparkSession available as 'spark'.
>>>

```

Pytest

ユニットテストの場合、AWS Glue のSpark ジョブスクリプトで pytest を使用できます。

以下のコマンドを実行し、準備を行います。

```
$ WORKSPACE_LOCATION=/local_path_to_workspace
$ SCRIPT_FILE_NAME=sample.py
$ UNIT_TEST_FILE_NAME=test_sample.py
$ mkdir -p ${WORKSPACE_LOCATION}/tests
$ vim ${WORKSPACE_LOCATION}/tests/${UNIT_TEST_FILE_NAME}
```

テストスイートで pytest を実行するため、以下のコマンドを実行します。

```
$ docker run -it -v ~/.aws:/home/glue_user/.aws -v $WORKSPACE_LOCATION:/home/glue_user/
workspace/ -e AWS_PROFILE=$PROFILE_NAME -e DISABLE_SSL=true --rm -p 4040:4040 -p
18080:18080 --name glue_pytest amazon/aws-glue-libs:glue_libs_4.0.0_image_01 -c
"python3 -m pytest"
starting org.apache.spark.deploy.history.HistoryServer,
 logging to /home/glue_user/spark/logs/spark-glue_user-
org.apache.spark.deploy.history.HistoryServer-1-5168f209bd78.out
*=====  
=====  
*platform linux -- Python 3.7.10, pytest-6.2.3, py-1.11.0, pluggy-0.13.1
rootdir: /home/glue_user/workspace
plugins: anyio-3.4.0
*collected 1 item *

tests/test_sample.py . [100%]

=====  
=====  
tests/test_sample.py::test_counts
/home/glue_user/spark/python/pyspark/sql/context.py:79: DeprecationWarning: Deprecated
in 3.0.0. Use SparkSession.builder.getOrCreate() instead.
DeprecationWarning)

-- Docs: https://docs.pytest.org/en/stable/warnings.html
=====  
=====  
1 passed, *1 warning* in
21.07s
```

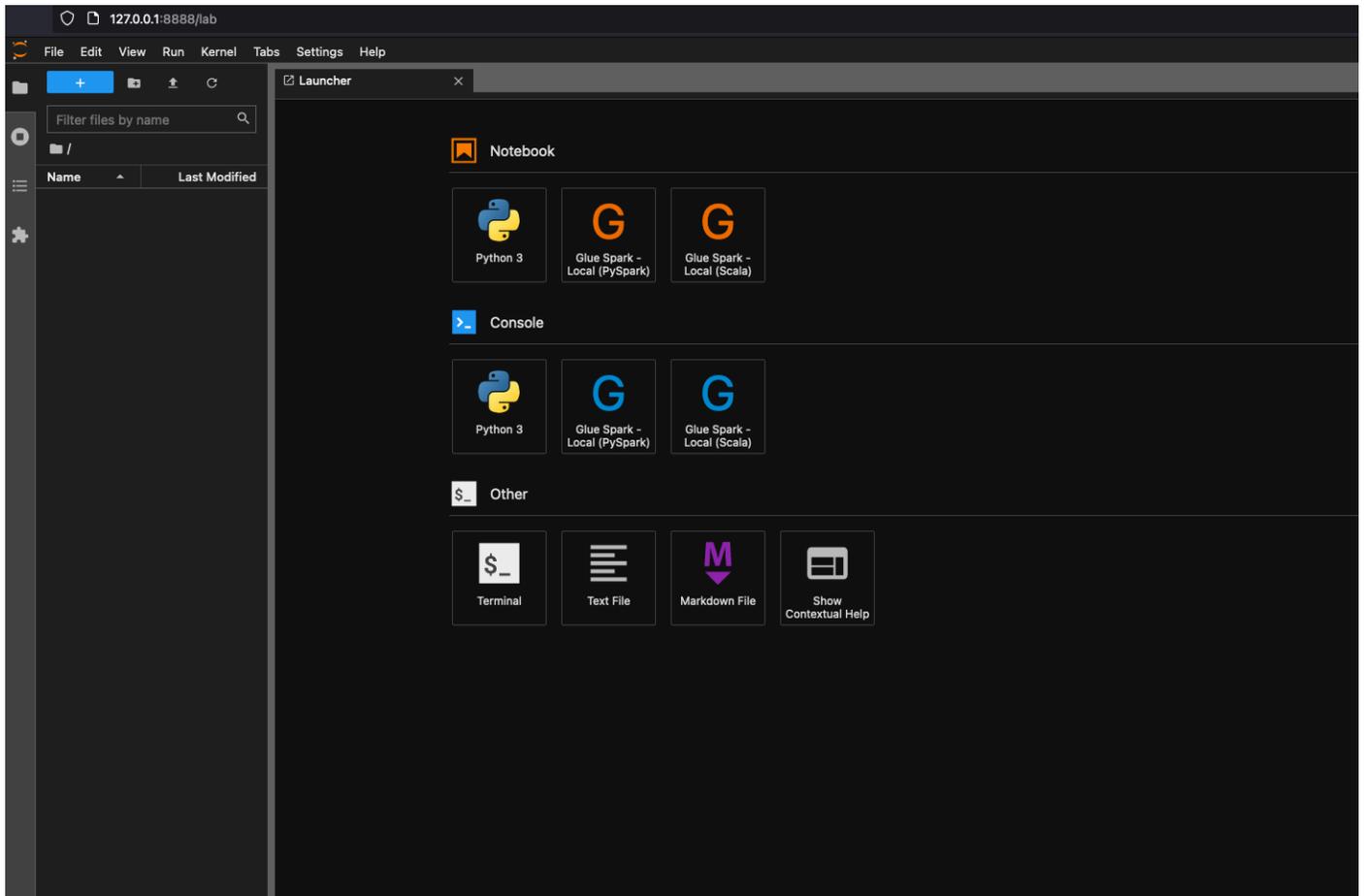
Jupyter ラボ

Jupyter を起動すると、ノートブックでインタラクティブな開発を行ったりアドホックなクエリを実行できます。

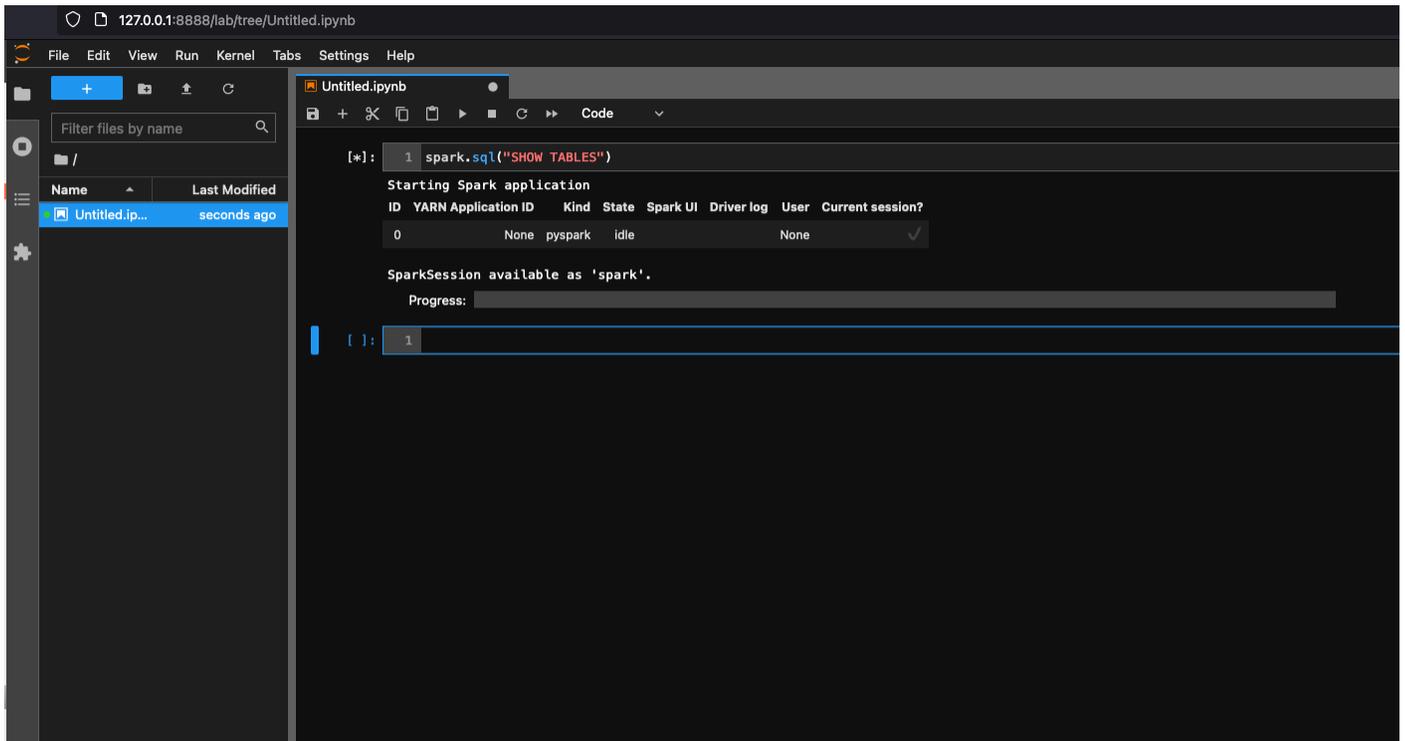
1. 次のコマンドを実行して、Jupyter ラボを起動します。

```
$ JUPYTER_WORKSPACE_LOCATION=/local_path_to_workspace/jupyter_workspace/  
$ docker run -it -v ~/.aws:/home/glue_user/.aws -v $JUPYTER_WORKSPACE_LOCATION:/home/glue_user/workspace/jupyter_workspace/ -e AWS_PROFILE=$PROFILE_NAME -e DISABLE_SSL=true --rm -p 4040:4040 -p 18080:18080 -p 8998:8998 -p 8888:8888 --name glue_jupyter_lab amazon/aws-glue-libs:glue_libs_4.0.0_image_01 /home/glue_user/jupyter/jupyter_start.sh  
...  
[I 2022-01-24 08:19:21.368 ServerApp] Serving notebooks from local directory: /home/glue_user/workspace/jupyter_workspace  
[I 2022-01-24 08:19:21.368 ServerApp] Jupyter Server 1.13.1 is running at:  
[I 2022-01-24 08:19:21.368 ServerApp] http://faa541f8f99f:8888/lab  
[I 2022-01-24 08:19:21.368 ServerApp] or http://127.0.0.1:8888/lab  
[I 2022-01-24 08:19:21.368 ServerApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).
```

2. ローカルマシンのウェブブラウザで <http://127.0.0.1:8888/lab> を開き、Jupyter ラボの UI を表示します。



3. [Notebook] (ノートブック) の下にある、[Glue Spark Local (PySpark)] (Glue スパークローカル (PySpark)) をクリックします。これで、インタラクティブな Jupyter Notebook UI によるコードの開発を開始できます。



Visual Studio Code を使用するためのコンテナのセットアップ

前提条件:

1. Visual Studio Code をインストールします。
2. [Python](#) をインストールします。
3. [Visual Studio Code Remote – コンテナ](#) をインストールします。
4. Visual Studio Code でワークスペースフォルダを開きます。
5. [設定] を選択します。
6. [Workspace] (ワークスペース) を選択します。
7. [Open Settings (JSON)] (設定を開く (JSON)) をクリックします。
8. 次の JSON を貼り付け、保存します。

```
{
  "python.defaultInterpreterPath": "/usr/bin/python3",
  "python.analysis.extraPaths": [
    "/home/glue_user/aws-glue-libs/PyGlue.zip:/home/glue_user/spark/python/lib/
py4j-0.10.9-src.zip:/home/glue_user/spark/python/"
  ]
}
```

```
}
```

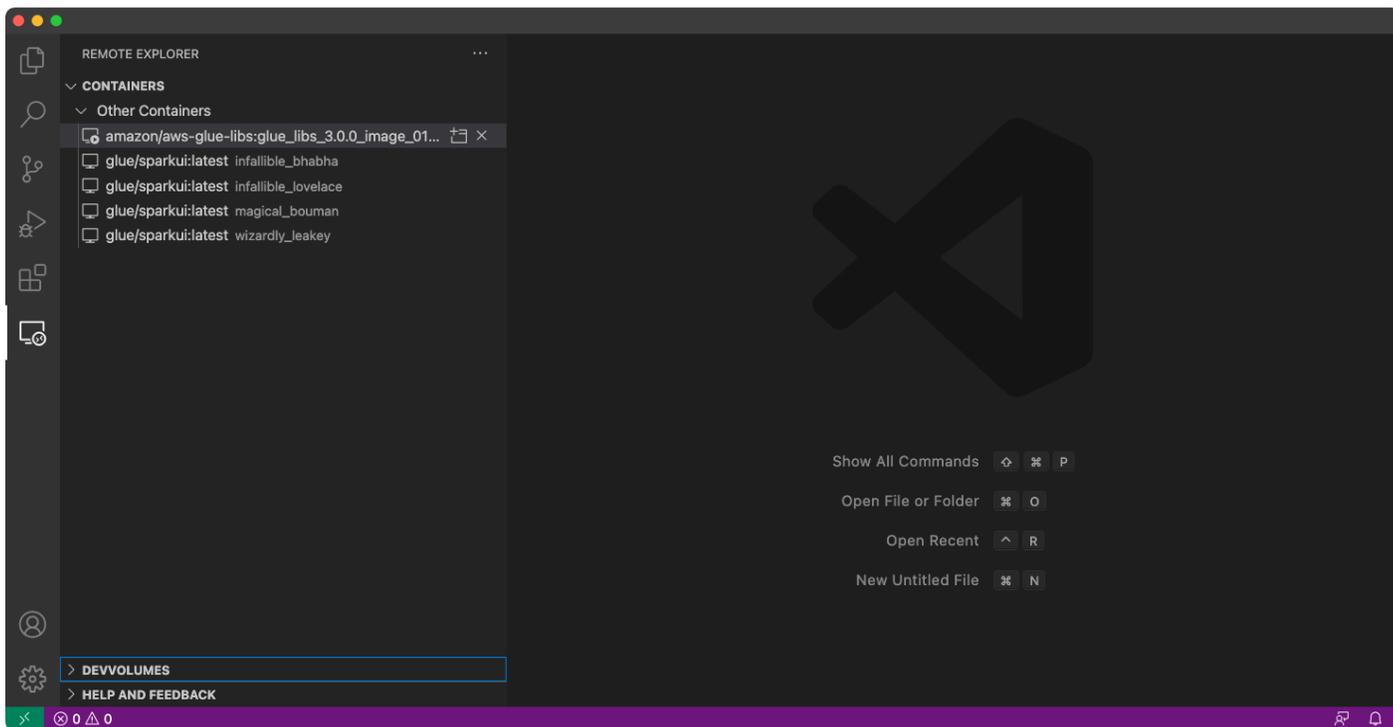
ステップ:

1. Docker コンテナを実行します。

```
$ docker run -it -v ~/.aws:/home/glue_user/.aws -v $WORKSPACE_LOCATION:/home/glue_user/workspace/ -e AWS_PROFILE=$PROFILE_NAME -e DISABLE_SSL=true --rm -p 4040:4040 -p 18080:18080 --name glue_pyspark amazon/aws-glue-libs:glue_libs_4.0.0_image_01 pyspark
```

2. Visual Studio Code を起動します。

3. 左側のメニューで [Remote Explorer] (リモートエクスプローラー) を選択した上で、amazon/aws-glue-libs:glue_libs_4.0.0_image_01 を選択します。



4. 右クリックし、[Attach to Container] (コンテナにアタッチ) を選択します。ダイアログが表示されたら、[Got it] (確認) をクリックします。

5. /home/glue_user/workspace/ を開きます。

6. Glue PySpark スクリプトを作成し、[Run] (実行) をクリックします。

このスクリプトは正常に実行されたことが表示されます。

```

6
7
8 class GluePythonSampleTest:
9     def __init__(self):
10         params = []
11         if '--JOB_NAME' in sys.argv:
12             params.append('JOB_NAME')
13         args = getResolvedOptions(sys.argv, params)
14
15         self.context = GlueContext(SparkContext.getOrCreate())
16         self.job = Job(self.context)
17
18         if 'JOB_NAME' in args:
19             jobname = args['JOB_NAME']
20         else:
21             jobname = "test"
22         self.job.init(jobname, args)
23
24     def run(self):
25         dyf = read_json(self.context, "s3://awsglue-datasets/examples/us-legislators/all/persons.json")
26         dyf.printSchema()

```

```

|-- name: string
|-- sort_name: string
|-- images: array
|   |-- element: struct
|       |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
|-- contact_details: array
|   |-- element: struct
|       |-- type: string
|       |-- value: string
|-- death_date: string

```

付録: テスト用の AWS Glue ジョブのサンプルコード

この付録では、テスト用の AWS Glue ジョブのサンプルコードとしてスクリプトを示します。

sample.py: Amazon S3 API コールにより AWS Glue ETL ライブラリを利用するためのサンプルコード

```

import sys
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from awsglue.utils import getResolvedOptions

```

```

class GluePythonSampleTest:
    def __init__(self):
        params = []
        if '--JOB_NAME' in sys.argv:
            params.append('JOB_NAME')
        args = getResolvedOptions(sys.argv, params)

```

```
self.context = GlueContext(SparkContext.getOrCreate())
self.job = Job(self.context)

if 'JOB_NAME' in args:
    jobname = args['JOB_NAME']
else:
    jobname = "test"
self.job.init(jobname, args)

def run(self):
    dyf = read_json(self.context, "s3://awsglue-datasets/examples/us-legislators/
all/persons.json")
    dyf.printSchema()

    self.job.commit()

def read_json(glue_context, path):
    dynamicframe = glue_context.create_dynamic_frame.from_options(
        connection_type='s3',
        connection_options={
            'paths': [path],
            'recurse': True
        },
        format='json'
    )
    return dynamicframe

if __name__ == '__main__':
    GluePythonSampleTest().run()
```

上記のコードには、AWS IAM 内に、Amazon S3 へのアクセス許可が必要です。IAM マネージドのポリシー `arn:aws:iam::aws:policy/AmazonS3ReadOnlyAccess` を付与するか、Amazon S3 パスに対し `ListBucket` および `GetObject` を呼び出せるようにする IAM のカスタムポリシーを付与します。

`test_sample.py`: `sample.py` でのユニットテスト向けのサンプルコード。

```
import pytest
from pyspark.context import SparkContext
from awsglue.context import GlueContext
```

```
from awsglue.job import Job
from awsglue.utils import getResolvedOptions
import sys
from src import sample

@pytest.fixture(scope="module", autouse=True)
def glue_context():
    sys.argv.append('--JOB_NAME')
    sys.argv.append('test_count')

    args = getResolvedOptions(sys.argv, ['JOB_NAME'])
    context = GlueContext(SparkContext.getOrCreate())
    job = Job(context)
    job.init(args['JOB_NAME'], args)

    yield(context)

    job.commit()

def test_counts(glue_context):
    dyf = sample.read_json(glue_context, "s3://awsglue-datasets/examples/us-
legislators/all/persons.json")
    assert dyf.toDF().count() == 1961
```

AWS Glue ETL ライブラリを使用した開発

AWS Glue ETL ライブラリはパブリックな Amazon S3 バケットに加えて、Apache Maven のビルドシステムでも使用が可能です。これにより、ネットワーク接続を必要とせずに Python および Scala の抽出、変換、ロード (ETL) スクリプトをローカルで開発およびテストできます。Docker イメージを使用したローカル開発が推奨されています。これを使用すると、このライブラリを使用するための環境が適切に構成されるからです。

ローカル開発は、AWS Glue バージョン0.9、1.0、2.0 以降を含む、すべての AWS Glue バージョンで使用可能です。AWS Glue で使用できる Python および Apache Spark のバージョンについては、「[Glue version job property](#)」を参照してください。

ライブラリは、Amazon ソフトウェアライセンス (<https://aws.amazon.com/asl>) とともにリリースされています。

ローカル開発の制限

AWS Glue Scala ライブラリを使用してローカルで開発する場合は、次の制限事項に注意してください。

- AWS Glue ライブラリでアセンブリ jar (「fat jar」または「uber jar」) を作成しないでください。作成すると、次の機能が無効になるためです。
 - [ジョブのブックマーク](#)
 - AWS Glue Parquet ライター ([AWS Glue で Parquet 形式を使用する](#))
 - FillMissingValues 変換 ([Scala](#) または [Python](#))

これらの機能は、AWS Glue ジョブシステム内でのみ使用できます。

- [FindMatches 変換](#) はローカル開発ではサポートされていません。
- [ベクトル化された SIMD CSV リーダー](#) はローカル開発ではサポートされていません。
- S3 パスから JDBC ドライバーをロードするための [customJdbcDriverS3Path](#) プロパティは、ローカル開発ではサポートされていません。代わりに、ローカルに JDBC ドライバーをダウンロードして、そこからロードすることはできます。
- [Glue Data Quality](#) はローカル開発ではサポートされていません。

Python を使用したローカルでの開発

いくつかの前提条件ステップを完了してから、AWS Glue ユーティリティを使用して Python ETL スクリプトをテストして送信します。

ローカル Python 開発の前提条件

ローカル Python 開発に備えるには、以下の手順を実行します。

1. GitHub (AWS Glue) から <https://github.com/aws-labs/aws-glue-libs> の Python リポジトリをクローンします。
2. 次のいずれかを行います。
 - AWS Glue バージョン 0.9 の場合は、ブランチ glue-0.9 を確認します。
 - AWS Glue バージョン 1.0 の場合は、ブランチ glue-1.0 を確認します。AWS Glue 0.9 以降のすべてのバージョンは、Python 3 をサポートしています。
 - AWS Glue バージョン 2.0 の場合は、ブランチ glue-2.0 を確認します。
 - AWS Glue バージョン 3.0 の場合は、ブランチ glue-3.0 を確認します。

- AWS Glue バージョン 4.0 の場合は、master ブランチを確認します。
3. Apache Maven を <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-common/apache-maven-3.6.0-bin.tar.gz> からインストールします。
 4. Apache Spark ディストリビューションを次のいずれかの場所からインストールします。
 - AWS Glue バージョン 0.9 の場合: <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-0.9/spark-2.2.1-bin-hadoop2.7.tgz>
 - AWS Glue バージョン 1.0 の場合: <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-1.0/spark-2.4.3-bin-hadoop2.8.tgz>
 - AWS Glue バージョン 2.0 の場合: <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-2.0/spark-2.4.3-bin-hadoop2.8.tgz>
 - AWS Glue バージョン 3.0 の場合: <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-3.0/spark-3.1.1-amzn-0-bin-3.2.1-amzn-3.tgz>
 - AWS Glue バージョン 4.0 の場合: <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-4.0/spark-3.3.0-amzn-1-bin-3.3.3-amzn-0.tgz>
 5. SPARK_HOME 環境変数をエクスポートし、Spark アーカイブから抽出したルートの場合に設定します。例:
 - AWS Glue バージョン 0.9 の場合: `export SPARK_HOME=/home/$USER/spark-2.2.1-bin-hadoop2.7`
 - AWS Glue バージョン 1.0 および 2.0 の場合: `export SPARK_HOME=/home/$USER/spark-2.4.3-bin-spark-2.4.3-bin-hadoop2.8`
 - AWS Glue バージョン 3.0 の場合: `export SPARK_HOME=/home/$USER/spark-3.1.1-amzn-0-bin-3.2.1-amzn-3`
 - AWS Glue バージョン 4.0 の場合: `export SPARK_HOME=/home/$USER/spark-3.3.0-amzn-1-bin-3.3.3-amzn-0`

Python ETL スクリプトの実行

ローカル開発に使用できる AWS Glue jar ファイルを使用すると、AWS Glue Python パッケージをローカルで実行できます。

Python スクリプトをテストして実行するには、次のユーティリティとフレームワークを使用します。次の表に示すコマンドは、[AWS Glue Python パッケージ](#)のルートディレクトリから実行されます。

ユーティリティ	コマンド	説明
AWS Glue シェル	<code>./bin/gluepyspark</code>	AWS Glue ETL ライブラリと統合されるシェルで Python スクリプトを入力して実行します。
AWS Glue [Submit (送信)]	<code>./bin/gluesparksubmit</code>	実行のために完全な Python スクリプトを送信します。
Pytest	<code>./bin/gluepytest</code>	Python コードのユニットテストを記述して実行します。pytest モジュールが PATH にインストールされ、使用可能になっている必要があります。詳細については、 pytest のドキュメント を参照してください。

Scala を使用したローカルでの開発

いくつかの前提条件ステップを完了してから、Maven コマンドを発行して Scala ETL スクリプトをローカルで実行します。

ローカル Scala 開発の前提条件

ローカルの Scala 開発に備えるには、以下のステップを実行します。

ステップ 1: ソフトウェアをインストールする

このステップでは、ソフトウェアをインストールし、必要な環境変数を設定します。

1. Apache Maven を <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-common/apache-maven-3.6.0-bin.tar.gz> からインストールします。
2. Apache Spark ディストリビューションを次のいずれかの場所からインストールします。
 - AWS Glue バージョン 0.9 の場合: <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-0.9/spark-2.2.1-bin-hadoop2.7.tgz>
 - AWS Glue バージョン 1.0 の場合: <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-1.0/spark-2.4.3-bin-hadoop2.8.tgz>
 - AWS Glue バージョン 2.0 の場合: <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-2.0/spark-2.4.3-bin-hadoop2.8.tgz>

- AWS Glue バージョン 3.0 の場合: <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-3.0/spark-3.1.1-amzn-0-bin-3.2.1-amzn-3.tgz>
 - AWS Glue バージョン 4.0 の場合: <https://aws-glue-etl-artifacts.s3.amazonaws.com/glue-4.0/spark-3.3.0-amzn-1-bin-3.3.3-amzn-0.tgz>
3. SPARK_HOME 環境変数をエクスポートし、Spark アーカイブから抽出したルートの場合に設定します。例:
- AWS Glue バージョン 0.9 の場合: `export SPARK_HOME=/home/$USER/spark-2.2.1-bin-hadoop2.7`
 - AWS Glue バージョン 1.0 および 2.0 の場合: `export SPARK_HOME=/home/$USER/spark-2.4.3-bin-spark-2.4.3-bin-hadoop2.8`
 - AWS Glue バージョン 3.0 の場合: `export SPARK_HOME=/home/$USER/spark-3.1.1-amzn-0-bin-3.2.1-amzn-3`
 - AWS Glue バージョン 4.0 の場合: `export SPARK_HOME=/home/$USER/spark-3.3.0-amzn-1-bin-3.3.3-amzn-0`

ステップ 2: Maven プロジェクトを設定する

AWS Glue Scala アプリケーションのテンプレートとして次の `pom.xml` ファイルを使用します。これには、必須の `dependencies`、`repositories`、`plugins` 要素が含まれています。Glue version を以下のいずれかに置き換えます。

- AWS Glue バージョン 4.0 の場合は 4.0.0
- 3.0.0 バージョン 3.0 の場合の AWS Glue
- 1.0.0 バージョン 1.0 または 2.0 の場合は AWS Glue
- 0.9.0 バージョン 0.9 の場合の AWS Glue

```
<project xmlns="http://maven.apache.org/POM/4.0.0" xmlns:xsi="http://www.w3.org/2001/
XMLSchema-instance" xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http://
maven.apache.org/maven-v4_0_0.xsd">
  <modelVersion>4.0.0</modelVersion>
  <groupId>com.amazonaws</groupId>
  <artifactId>AWSGlueApp</artifactId>
  <version>1.0-SNAPSHOT</version>
  <name>${project.artifactId}</name>
  <description>AWS ETL application</description>
```

```
    <properties>
      <scala.version>2.11.1 for AWS Glue 2.0 or below, 2.12.7 for AWS Glue 3.0
and 4.0</scala.version>
      <glue.version>Glue version with three numbers (as mentioned earlier)</
glue.version>
    </properties>
  <dependencies>
    <dependency>
      <groupId>org.scala-lang</groupId>
      <artifactId>scala-library</artifactId>
      <version>${scala.version}</version>
<!-- A "provided" dependency, this will be ignored when you package your application
-->
      <scope>provided</scope>
    </dependency>
    <dependency>
      <groupId>com.amazonaws</groupId>
      <artifactId>AWSGlueETL</artifactId>
<version>${glue.version}</version>
      <!-- A "provided" dependency, this will be ignored when you package your
application -->
      <scope>provided</scope>
    </dependency>
  </dependencies>

  <repositories>
    <repository>
      <id>aws-glue-etl-artifacts</id>
      <url>https://aws-glue-etl-artifacts.s3.amazonaws.com/release/</url>
    </repository>
  </repositories>
  <build>
    <sourceDirectory>src/main/scala</sourceDirectory>
    <plugins>
      <plugin>
        <!-- see http://davidb.github.com/scala-maven-plugin -->
        <groupId>net.alchim31.maven</groupId>
        <artifactId>scala-maven-plugin</artifactId>
        <version>3.4.0</version>
        <executions>
          <execution>
            <goals>
              <goal>compile</goal>
              <goal>testCompile</goal>
            </goals>
          </execution>
        </executions>
      </plugin>
    </plugins>
  </build>
</project>
```

```
        </goals>
      </execution>
    </executions>
  </plugin>
<plugin>
  <groupId>org.codehaus.mojo</groupId>
  <artifactId>exec-maven-plugin</artifactId>
  <version>1.6.0</version>
  <executions>
    <execution>
      <goals>
        <goal>java</goal>
      </goals>
    </execution>
  </executions>
  <configuration>
    <systemProperties>
      <systemProperty>
        <key>spark.master</key>
        <value>local[*]</value>
      </systemProperty>
      <systemProperty>
        <key>spark.app.name</key>
        <value>localrun</value>
      </systemProperty>
      <systemProperty>
        <key>org.xerial.snappy.lib.name</key>
        <value>libsnappyjava.jnilib</value>
      </systemProperty>
    </systemProperties>
  </configuration>
</plugin>
<plugin>
  <groupId>org.apache.maven.plugins</groupId>
  <artifactId>maven-enforcer-plugin</artifactId>
  <version>3.0.0-M2</version>
  <executions>
    <execution>
      <id>enforce-maven</id>
      <goals>
        <goal>enforce</goal>
      </goals>
      <configuration>
        <rules>
```

```
        <requireMavenVersion>
            <version>3.5.3</version>
        </requireMavenVersion>
    </rules>
</configuration>
</execution>
</executions>
</plugin>
<!-- The shade plugin will be helpful in building a uberjar or fatjar.
You can use this jar in the AWS Glue runtime environment. For more information, see
https://maven.apache.org/plugins/maven-shade-plugin/ -->
<plugin>
    <groupId>org.apache.maven.plugins</groupId>
    <artifactId>maven-shade-plugin</artifactId>
    <version>3.2.4</version>
    <configuration>
        <!-- any other shade configurations -->
    </configuration>
    <executions>
        <execution>
            <phase>package</phase>
            <goals>
                <goal>shade</goal>
            </goals>
        </execution>
    </executions>
</plugin>
</plugins>
</build>
</project>
```

Scala ETL スクリプトの実行

Maven プロジェクトのルートディレクトリから次のコマンドを実行して、Scala ETL スクリプトを実行します。

```
mvn exec:java -Dexec.mainClass="mainClass" -Dexec.args="--JOB-NAME jobName"
```

mainClass を、スクリプトのメインクラスの完全修飾クラス名に置き換えます。*jobName* を目的のジョブ名に置き換えます。

テスト環境の設定

ローカルテスト環境の設定例については、以下のブログ記事を参照してください。

- [AWS Glueアカウントなしでローカルに AWS ETL パイプラインを構築する](#)
- [コンテナを利用したローカルでの AWS Glue ETLジョブの開発](#)

ETL スクリプトをテストするために、開発エンドポイントまたはノートブックを使用する必要がある場合は、「[開発エンドポイントを使用してスクリプトを開発する](#)」を参照してください。

Note

開発エンドポイントは、AWS Glue バージョン 2.0 ジョブでの使用がサポートされていません。詳細については、「[Running Spark ETL Jobs with Reduced Startup Times](#)」を参照してください。

開発エンドポイント

Note

開発エンドポイントのコンソール体験は、2023 年 3 月 31 日に削除されました。開発エンドポイントの作成、更新、モニタリングは、[開発エンドポイント API](#) および [AWS Glue CLI](#) を介して引き続きご使用いただけます。

以下に示す理由により、開発エンドポイントからインタラクティブセッションに移行することを強くお勧めします。開発エンドポイントからインタラクティブセッションへの移行に必要なアクションについては、「[開発エンドポイントからインタラクティブセッションへ移行する](#)」を参照してください。

説明	開発者エンドポイント	インタラクティブセッション
Glue バージョンサポート	AWS Glue バージョン 0.9 および 1.0 をサポートしています	AWS Glue バージョン 2.0 以降をサポートしています

説明	開発者エンドポイント	インタラクティブセッション
<p>開発エンドポイントは、アジアパシフィック (ジャカルタ) (ap-southeast-3)、中東 (UAE) (me-central-1)、欧州 (スペイン) (eu-south-2)、欧州 (チューリッヒ) (eu-central-2) や、今後開設される新しいリージョンでは利用できません。</p>	<p>現在、インタラクティブセッションは中東 (UAE) (me-central-1) リージョンでは利用できませんが、今後は利用できるようになる可能性があります</p>	
<p>Spark クラスターへのアクセス方法</p>	<p>SSH、REPL シェル、Jupyter Notebook、および IDE (PyCharm など) をサポートしています</p>	<p>AWS Glue Studio ノートブック、Jupyter Notebook、さまざまな IDE (Visual Studio Code、PyCharm など)、および SageMaker ノートブックをサポートしています</p>
<p>先頭クエリまでの時間</p>	<p>Spark クラスターのセットアップには 10~15 分かかります</p>	<p>エフェメラル Spark クラスターの設定には、最大 1 分程度かかる場合があります</p>

説明	開発者エンドポイント	インタラクティブセッション
料金モデル	AWS は、エンドポイントがプロビジョニングされた時間と DPU の数に基づいて、開発エンドポイントの料金を請求します。開発エンドポイントはタイムアウトしません。プロビジョニングされた開発エンドポイントごとに 10 分の最小課金時間が適用されます。さらに、AWS は、Amazon EC2 インスタンスでの Jupyter Notebook と、開発エンドポイントで設定した場合の SageMaker ノートブックに対して料金を請求します。	AWS は、セッションがアクティブになっている時間と DPU の数に基づいて、インタラクティブセッションの料金を請求します。インタラクティブセッションには、設定可能なアイドルタイムアウトがあります。AWS Glue Studio ノートブックには、インタラクティブなセッション用の組み込みインターフェイスが用意されており、追加費用なしで提供されます。インタラクティブセッションごとに 1 分の最小課金時間が適用されます。AWS Glue Studio ノートブックには、インタラクティブなセッション用の組み込みインターフェイスが用意されているため、追加コストなしで提供されます
コンソールエクスペリエンス	CLI および API 経由でのみ使用可能	AWS Glue コンソール、CLI、および API を介して使用可能

開発エンドポイントからインタラクティブセッションへ移行する

次のチェックリストを使用して、開発エンドポイントからインタラクティブセッションに移行する適切な方法を決定してください。

スクリプトは AWS Glue 0.9 または 1.0 の特定の機能 (HDFS、YARN など) に依存していますか？

答えが「はい」の場合は、「[AWS Glue ジョブの AWS Glue バージョン 3.0 への移行](#)」を参照して、Glue 0.9 または 1.0 から Glue 3.0 以降に移行する方法を確認してください。

開発エンドポイントへのアクセスにはどの方法を使用していますか？

この方法を使用する場合	こちらを実行してください
SageMaker ノートブック、Jupyter Notebook、または JupyterLab	Jupyter で .ipynb ファイルをダウンロードして AWS Glue Studio ノートブック に移行し、.ipynb ファイルをアップロードして新しい AWS Glue Studio ノートブックジョブを作成します。または、 SageMaker Studio を使用して AWS Glue カーネルを選択することもできます。
Zeppelin ノートブック	ノートブックの Jupyter Notebook への変換は、コードをコピーして貼り付ける操作を手動で行うか、ze2nb などのサードパーティコンバータを使用して自動で行います。次に、ノートブックを AWS Glue Studio ノートブックまたは SageMaker Studio で使用します。
IDE	「 AWS Glue インタラクティブセッションを使用した PyCharm による AWS Glue ジョブのオーサリング 」、または「 Microsoft Visual Studio Code によるインタラクティブセッションの使用 」を参照してください。
REPL	<p>aws-glue-session package をローカルでインストールしたら、次のコマンドを実行します。</p> <ul style="list-style-type: none">• Python の場合: <code>jupyter console --kernal glue_pyspark</code>• Scala の場合: <code>jupyter console --kernal glue_spark</code>
SSH	インタラクティブセッションでは対応するオプションがありません。または、Docker イメージを使用することもできます。詳細は「 Docker イメージによる開発 」をご覧ください。

以下のセクションでは、デベロッパーエンドポイントを使用してジョブを開発する方法の情報を提供します。AWS Glueバージョン 1.0

トピック

- [開発エンドポイントを使用してスクリプトを開発する](#)
- [ノートブックの管理](#)

開発エンドポイントを使用してスクリプトを開発する

Note

開発エンドポイントは、AWS Glue 2.0 より前のバージョンでのみサポートされます。ETL スクリプトを作成およびテストできるインタラクティブな環境の場合は、[AWS Glue Studio](#) でノートブックを使用します。

AWS Glue では、開発エンドポイントと呼ばれる環境を作成できます。この環境を使用して、抽出、変換、ロード (ETL) スクリプトを反復的に開発およびテストできます。開発エンドポイントを作成、編集、削除するには、AWS Glue コンソールまたは API を使用します。

開発環境の管理

開発エンドポイントを作成するときは、開発環境をプロビジョニングするための設定値を指定します。これらの値は、ユーザーが開発エンドポイントに安全にアクセスできるように、またエンドポイントがデータストアにアクセスできるように、ネットワークを設定する方法を AWS Glue に指示します。

次に、開発エンドポイントに接続するノートブックを作成し、このノートブックを使用して ETL スクリプトを作成およびテストできます。開発プロセスの結果に満足したら、スクリプトを実行する ETL ジョブを作成します。このプロセスにより、インタラクティブな方法で機能を追加してスクリプトをデバッグできます。

このセクションのチュートリアルに従って、ノートブックで開発エンドポイントを使用する方法を学習します。

トピック

- [開発エンドポイントのワークフロー](#)

- [AWS Glue 開発エンドポイントで SageMaker ノートブックを使用する方法](#)
- [開発エンドポイントの追加](#)
- [開発エンドポイントへのアクセス](#)
- [チュートリアル: JupyterLab で Jupyter Notebook をセットアップして ETL スクリプトをテストおよびデバッグする](#)
- [チュートリアル: 開発エンドポイントで Amazon SageMaker ノートブックを使用する](#)
- [チュートリアル: 開発エンドポイントで REPL シェルを使用する](#)
- [チュートリアル: 開発エンドポイントで PyCharm Professional をセットアップする](#)
- [高度な設定: 複数のユーザー間で開発エンドポイントを共有する](#)

開発エンドポイントのワークフロー

AWS Glue の開発エンドポイントを使用するには、次のワークフローに従います。

1. API を使用して開発エンドポイントを作成します。このエンドポイントは、定義したセキュリティグループを使用して、Virtual Private Cloud (VPC) 内で起動されます。
2. API は、開発エンドポイントがプロビジョニングされて使用可能になるまで、開発エンドポイントをポーリングします。使用可能になると、次のいずれかの方法を使用して開発エンドポイントに接続し、AWS Glue スクリプトを作成してテストします。
 - アカウントに SageMaker ノートブックを作成します。ノートブックの作り方の詳細については、「[the section called “AWS Glue Studio ノートブックによるコードの作成”](#)」を参照してください。
 - ターミナルウィンドウを開いて、開発エンドポイントに直接接続します。
 - JetBrains [PyCharm Python IDE](#) の Professional エディションがある場合は、これを開発エンドポイントに接続し、これを使用してインタラクティブに開発を行います。スクリプトに `pydevd` ステートメントを挿入すると、PyCharm でリモートブレイクポイントをサポートできます。
3. 開発エンドポイントでのデバッグとテストが完了したら、削除することができます。

AWS Glue 開発エンドポイントで SageMaker ノートブックを使用する方法

開発エンドポイントにアクセスする一般的な方法の 1 つは、SageMaker ノートブックの [Jupyter](#) を使用することです。Jupyter Notebook は、可視化、分析、機械学習などで広く使用されているオープンソースのウェブアプリケーションです。AWS Glue SageMaker ノートブックにより、AWS Glue

開発エンドポイントを使用した Jupyter Notebook を使用できます。AWS Glue SageMaker ノートブックでは、Jupyter Notebook 環境は [SparkMagic](#) (リモートの Spark クラスターに Spark ジョブを送信するオープンソースの Jupyter プラグイン) で事前設定されています。[Apache Livy](#)は、REST API を介してリモート Spark クラスターとのやり取りを可能にするサービスです。AWS Glue SageMaker ノートブックでは、SparkMagic は、AWS Glue 開発エンドポイントで実行されている Livy サーバーに対して REST API を呼び出すように設定されています。

各コンポーネントの動作について、次のテキストフローで説明します。

AWS Glue SageMaker ノートブック: (Jupyter → SparkMagic) → (ネットワーク) → AWS Glue 開発エンドポイント: (Apache Livy → Apache Spark)

Jupyter Notebook 上の各段落に記述された Spark スクリプトを実行すると、Spark コードが SparkMagic 経由で Livy サーバーに送信され、「livy-session-N」という名前の Spark ジョブが Spark クラスター上で実行されます。このジョブは、Livy セッションと呼ばれます。ノートブックセッションが存続している間、Spark ジョブは実行されます。Spark ジョブは、ノートブックから Jupyter カーネルをシャットダウンしたとき、またはセッションがタイムアウトしたときに終了します。ノートブック (.ipynb) ファイルごとに 1 つの Spark ジョブが起動します。

単一の AWS Glue 開発エンドポイントで複数の SageMaker ノートブックインスタンスを使用できます。SageMaker ノートブックインスタンスごとに、複数のノートブックファイルを作成できます。各ノートブックファイルを開いて段落を実行すると、SparkMagic 経由で Spark クラスター上のノートブックファイルごとに Livy セッションが起動します。各 Livy セッションは、単一の Spark ジョブに対応します。

AWS Glue 開発エンドポイントと SageMaker ノートブックのデフォルトの動作

Spark ジョブは、[\[Spark configuration\]](#) (Spark の設定) に基づいて実行されます。Spark の設定には複数の方法があります (Spark クラスター設定、SparkMagic の設定など)。

デフォルトでは、Spark は Spark クラスター設定に基づいて Livy セッションにクラスターリソースを割り当てます。AWS Glue 開発エンドポイントでは、クラスター設定はワーカーのタイプによって異なります。ワーカータイプごとの共通設定について説明した表を次に示します。

	規格	G.1X	G.2X
spark.dri ver.memor y	5G	10G	20G

	規格	G.1X	G.2X
<code>spark.executor.memory</code>	5G	10G	20G
<code>spark.executor.cores</code>	4	8	16
<code>spark.dynamicAllocation.enabled</code>	TRUE	TRUE	TRUE

Spark エグゼキュターの最大数は、DPU (または `NumberOfWorkers`) とワーカータイプの組み合わせによって自動的に計算されます。

	規格	G.1X	G.2X
Spark エグゼキュターの最大数	$(\text{DPU} - 1) * 2 - 1$	$(\text{NumberOfWorkers} - 1)$	$(\text{NumberOfWorkers} - 1)$

例えば、開発エンドポイントに 10 ワーカーあり、ワーカータイプが G.1X の場合、Spark エグゼキュターは 9 つになり、各エグゼキュターは 10G のメモリを持つため、クラスター全体のエグゼキュターメモリは 90G になります。

指定されたワーカータイプに関係なく、Spark 動的リソース割り当てが有効になります。データセットが十分に大きい場合、`spark.dynamicAllocation.maxExecutors` はデフォルトで設定されていないため、Spark はすべてのエグゼキュターを単一の Livy セッションに割り当てることができます。つまり、同じ開発エンドポイント上の他の Livy セッションは、新しいエグゼキュターの起動を待つことになります。データセットが小さい場合、Spark は同時に複数の Livy セッションにエグゼキュターを割り当てることができます。

Note

さまざまなユースケースでリソースがどのように割り当てられるか、および動作を変更するための設定方法の詳細については、「[高度な設定: 複数のユーザー間で開発エンドポイントを共有する](#)」を参照してください。

開発エンドポイントの追加

開発エンドポイントを使用して、AWS Glue で抽出、変換、ロード (ETL) スクリプトを反復的に開発およびテストします。開発エンドポイントは、AWS Command Line Interface を介してのみ使用できます。

1. コマンドラインウィンドウで、次のようなコマンドを入力します。

```
aws glue create-dev-endpoint --endpoint-name "endpoint1" --role-arn
"arn:aws:iam::account-id:role/role-name" --number-of-nodes "3" --glue-version
"1.0" --arguments '{"GLUE_PYTHON_VERSION": "3"}' --region "region-name"
```

このコマンドでは、AWS Glue バージョン 1.0 を指定します。このバージョンは Python 2 と Python 3 の両方をサポートしているため、必要な Python バージョンを指定するには arguments パラメータを使用します。glue-version パラメータを省略すると、AWS Glue バージョン 0.9 が想定されます。AWS Glue のバージョンの詳細については、「[Glue version job property](#)」を参照してください。

追加のコマンドラインパラメータについては、AWS CLI コマンドリファレンスの「[create-dev-endpoint](#)」を参照してください。

2. (オプション) 次のコマンドを入力して、開発エンドポイントのステータスを確認します。ステータスが READY に変わったら、開発エンドポイントの使用を開始できます。

```
aws glue get-dev-endpoint --endpoint-name "endpoint1"
```

開発エンドポイントへのアクセス

Virtual Private Cloud (VPC) に開発エンドポイントを作成すると、AWS Glue はプライベート IP アドレスのみを返します。パブリック IP アドレスフィールドは返されません。VPC 以外の開発エンドポイントを作成する場合、AWS Glue はパブリック IP アドレスのみを返します。

開発エンドポイントにパブリックアドレスがある場合は、次の例のように、開発エンドポイントの SSH プライベートキーを使用して、このアドレスに到達できることを確認します。

```
ssh -i dev-endpoint-private-key.pem glue@public-address
```

開発エンドポイントにプライベートアドレスがあり、VPC サブネットがパブリックインターネットからルーティング可能であり、そのセキュリティグループがクライアントからのインバウンドアクセスを許可するとします。この場合は、以下のステップに従って Elastic IP アドレスを開発エンドポイントにアタッチし、インターネットからのアクセスを許可します。

Note

Elastic IP アドレスを使用する場合、使用されているサブネットには、ルートテーブルを通じて関連付けられたインターネットゲートウェイが必要です。

Elastic IP アドレスをアタッチして開発エンドポイントにアクセスするには

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。
2. ナビゲーションペインで、[開発エンドポイント] を選択し、開発エンドポイントの詳細ページに移動します。次のステップで使用するためにプライベートアドレスを書き留めます。
3. Amazon EC2 コンソール (<https://console.aws.amazon.com/ec2/>) を開きます。
4. ナビゲーションペインの [Network & Security] で、[ネットワークインターフェイス] を選択します。
5. AWS Glue コンソールの [development endpoint details] (開発エンドポイントの詳細) ページで、プライベートアドレスに対応するプライベート DNS (IPv4) を検索します。

必要に応じて Amazon EC2 コンソールで表示する列を変更します。このアドレスのネットワークインターフェイス ID (ENI) を書き留めます (例: eni-12345678)。

6. Amazon EC2 コンソールの [Network & Security] (ネットワーク & セキュリティ) で、[Elastic IPs] (Elastic IP) をクリックします。
7. [新しいアドレスの割り当て]、[割り当て] の順に選択して、新しい Elastic IP アドレスを割り当てます。
8. [Elastic IP] ページで、新しく割り当てた Elastic IP を選択します。[アクション]、[アドレスの関連付け] の順に選択します。

9. [アドレスの関連付け] ページで、以下の操作を行います。
 - [リソースタイプ] で、[ネットワークインターフェイス] を選択します。
 - [ネットワークインターフェイス] ボックスに、プライベートアドレスのネットワークインターフェイス ID (ENI) を入力します。
 - [関連付ける] を選択します。
10. 次の例に示すように、開発エンドポイントに関連付けられた SSH プライベートキーを使用して、新しく関連付けられた Elastic IP アドレスに到達できることを確認します。

```
ssh -i dev-endpoint-private-key.pem glue@elastic-ip
```

踏み台ホストを使用して、開発エンドポイントのプライベートアドレスへの SSH アクセスを取得する方法については、AWS セキュリティブログの投稿「[Securely Connect to Linux Instances Running in a Private Amazon VPC](#)」を参照してください。

チュートリアル: JupyterLab で Jupyter Notebook をセットアップして ETL スクリプトをテストおよびデバッグする

このチュートリアルでは、ローカルマシンで実行されている JupyterLab の Jupyter Notebook を開発エンドポイントに接続します。これは、デプロイする前に、AWS Glue 抽出、変換、ロード (ETL) スクリプトを対話的に実行、デバッグ、およびテストできるようにするためです。このチュートリアルでは、Secure Shell (SSH) ポートフォワーディングを使用して、ローカルマシンを AWS Glue 開発エンドポイントに接続します。詳細については、Wikipedia の「[Port forwarding](#)」を参照してください。

ステップ1: JupyterLab と Sparkmagic をインストールする

JupyterLab は conda または pip を使用してインストールできます。conda は、Windows、macOS、および Linux で稼働するオープンソースのパッケージ管理システムおよび環境管理システムです。pip は Python のパッケージインストーラです。

macOS にインストールする場合は、Sparkmagic をインストールする前に Xcode をインストールしておく必要があります。

1. JupyterLab、Sparkmagic、および関連する拡張機能をインストールします。

```
$ conda install -c conda-forge jupyterlab
```

```
$ pip install sparkmagic
$ jupyter nbextension enable --py --sys-prefix widgetsnbextension
$ jupyter labextension install @jupyter-widgets/jupyterlab-manager
```

2. Location から sparkmagic ディレクトリをチェックします。

```
$ pip show sparkmagic | grep Location
Location: /Users/username/.pyenv/versions/anaconda3-5.3.1/lib/python3.7/site-packages
```

3. ディレクトリを、Location に対して返されたディレクトリに変更し、Scala と PySpark 用のカーネルをインストールします。

```
$ cd /Users/username/.pyenv/versions/anaconda3-5.3.1/lib/python3.7/site-packages
$ jupyter-kernelspec install sparkmagic/kernels/sparkkernel
$ jupyter-kernelspec install sparkmagic/kernels/pysparkkernel
```

4. サンプルの config ファイルをダウンロードします。

```
$ curl -o ~/.sparkmagic/config.json https://raw.githubusercontent.com/jupyter-incubator/sparkmagic/master/sparkmagic/example_config.json
```

この設定ファイルでは、driverMemory や executorCores などの Spark 関連のパラメータを設定できます。

ステップ 2: JupyterLab を起動する

JupyterLab を起動すると、デフォルトのウェブブラウザが自動的に開き、`http://localhost:8888/lab/workspaces/{workspace_name}` という URL が表示されます。

```
$ jupyter lab
```

ステップ 3: 開発エンドポイントに接続するための SSH ポート転送を開始する

次に、SSH ローカルポート転送を使用して、ローカルポート (ここでは 8998) を AWS Glue で定義されたリモート送信先 (169.254.76.1:8998) に転送します。

1. SSH へのアクセスができる別のターミナルウィンドウを開きます。Microsoft Windows の場合、[Git for Windows](#) に用意されている BASH シェルを使用するか、[Cygwin](#) をインストールすることができます。

2. 次のように変更した、以下の SSH コマンドを実行します。

- `private-key-file-path` を、開発エンドポイントを作成するのに使用したパブリックキーに対応するプライベートキーを含む `.pem` ファイルへのパスに置き換えます。
- 8998 とは異なるポートを転送している場合は、8998 をローカルで実際に使用しているポート番号に置き換えます。この `169.254.76.1:8998` というアドレスはリモートポートです。変更することはできません。
- `dev-endpoint-public-dns` を開発エンドポイントのパブリック DNS アドレスで置き換えます。このアドレスを確認するには、AWS Glue コンソールで開発エンドポイントに移動して名前を選択し、[Endpoint details] (エンドポイントの詳細) ページに一覧表示されている [Public address] (パブリックアドレス) をコピーします。

```
ssh -i private-key-file-path -NTL 8998:169.254.76.1:8998 glue@dev-endpoint-public-dns
```

以下のような警告メッセージが表示されます。

```
The authenticity of host 'ec2-xx-xxx-xxx-xx.us-west-2.compute.amazonaws.com  
(xx.xxx.xxx.xx)'  
can't be established. ECDSA key fingerprint is SHA256:4e97875Brt+1wKzRko  
+Jf15np21X7aTP3BcFnHYLEts.  
Are you sure you want to continue connecting (yes/no)?
```

yes と入力し、JupyterLab を使用中はターミナルウィンドウを開いたままにしておきます。

3. SSH ポート転送が開発エンドポイントで正しく機能していることを確認します。

```
$ curl localhost:8998/sessions  
{"from":0,"total":0,"sessions":[]}
```

ステップ 4: ノートブックの段落でシンプルスクリプトフラグメントを実行する

JupyterLab のノートブックは、開発エンドポイントで動作するようになっています。次のスクリプトフラグメントをノートブックに入力して実行します。

1. Spark が正常に実行されていることを確認します。次のコマンドによって、Spark が 1 を計算して値を出力します。

```
spark.sql("select 1").show()
```

2. AWS Glue Data Catalog 統合が機能しているかどうか確認します。次のコマンドで、データカタログ内のテーブルが一覧表示されます。

```
spark.sql("show tables").show()
```

3. AWS Glue ライブラリを使用するシンプルスクリプトフラグメントが機能していることを確認します。

次のスクリプトでは、AWS Glue Data Catalog の `persons_json` テーブルメタデータを使用して、サンプルデータから `DynamicFrame` を作成します。次に、このデータの項目数およびスキーマが出力されます。

```
import sys
from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create a Glue context
glueContext = GlueContext(SparkContext.getOrCreate())

# Create a DynamicFrame using the 'persons_json' table
persons_DyF = glueContext.create_dynamic_frame.from_catalog(database="legislators",
    table_name="persons_json")

# Print out information about *this* data
print("Count: ", persons_DyF.count())
persons_DyF.printSchema()
```

スクリプトの出力は次のとおりです。

```
Count: 1961
root
|-- family_name: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
```

```
|-- gender: string
|-- image: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- name: string
|   |   |-- lang: string
|-- sort_name: string
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- death_date: string
```

トラブルシューティング

- JupyterLab のインストール中に、コンピュータが企業のプロキシまたはファイアウォールの内側にある場合、企業の IT 部門によって管理されるカスタムセキュリティプロファイルにより HTTP および SSL エラーが発生することがあります。

conda が独自のリポジトリに接続できないときに、よく発生するエラーの例を次に示します。

```
CondaHTTPError: HTTP 000 CONNECTION FAILED for url <https://repo.anaconda.com/pkg/main/win-64/current_repodata.json>
```

これは、会社が Python および JavaScript コミュニティで広く使用されているリポジトリへの接続をブロックすることがあるために発生した可能性があります。詳細については、JupyterLab ウェブサイトの「[Installation Problems](#)」を参照してください。

- 開発エンドポイントに接続しようとしたときに「connection refused (接続拒否)」エラーが発生した場合、古い開発エンドポイントを使用している可能性があります。新しい開発エンドポイントを作成して再接続してみます。

チュートリアル: 開発エンドポイントで Amazon SageMaker ノートブックを使用する

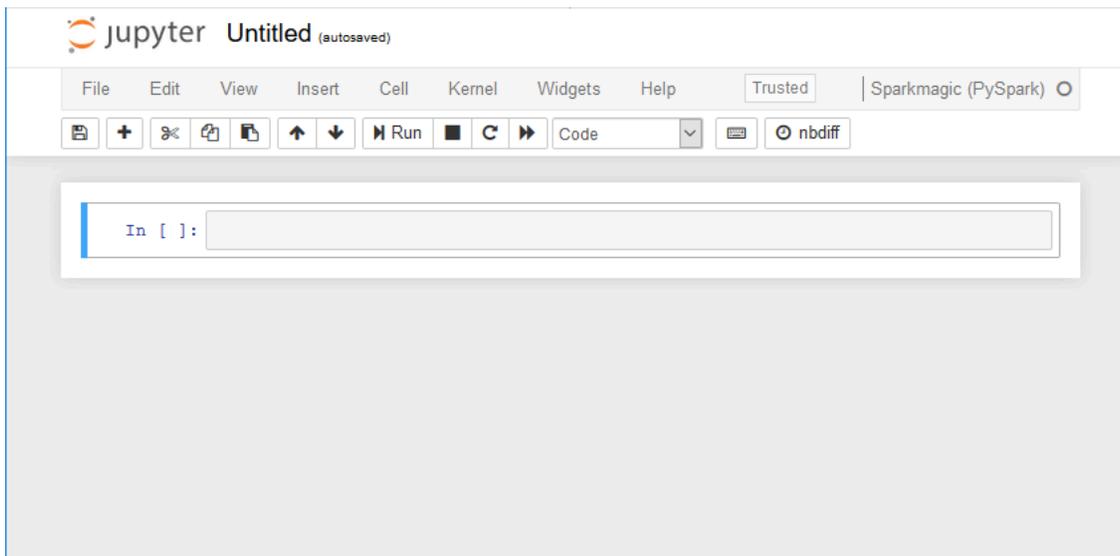
AWS Glue では、開発エンドポイントを作成し、ETL および機械学習スクリプトの開発に役立つ SageMaker ノートブックを作成できます。SageMaker ノートブックは、Jupyter Notebook アプリケーションを実行しているフルマネージド型機械学習コンピューティングインスタンスです。

1. AWS Glue コンソールで [Dev endpoints (開発エンドポイント)] を選択して、開発エンドポイントのリストに移動します。
2. 使用する開発エンドポイントの名前の横にあるチェックボックスをオンにし、[アクション] メニューで [Create SageMaker notebook (SageMaker ノートブックの作成)] を選択します。
3. [Create and configure a notebook (ノートブックの作成と設定)] ページに次のように入力します。
 - a. ノートブック名を入力します。
 - b. [Attach to development endpoint (開発エンドポイントにアタッチ)] で、開発エンドポイントを確認します。
 - c. AWS Identity and Access Management (IAM) ロールを作成または選択します。

ロールを作成することをお勧めします。既存のロールを使用する場合は、必要な権限があることを確認します。(詳細については、[the section called “ステップ 6: SageMaker ノートブック用に IAM ポリシーを作成する”](#) を参照してください)。

- d. (オプション) VPC、サブネット、および 1 つ以上のセキュリティグループを選択します。
 - e. (オプション) AWS Key Management Service 暗号化キーを選択します。
 - f. (オプション) ノートブックインスタンスのタグを追加します。
4. [Create notebook (ノートブックの作成)] を選択します。[ノートブック] ページで、右上にある更新アイコンを選択し、[ステータス] に Ready と表示されるまで続行します。
 5. 新しいノートブック名の横にあるチェックボックスをオンにし、[ノートブックを開く] を選択します。
 6. 新しいノートブックを作成する: jupyter ページで [新規] を選択し、[Sparkmagic (PySpark)] を選択します。

これで、画面は以下のようになります。



7. (オプション) ページの上部にある [Untitled (無題)] を選択し、ノートブックに名前を付けます。
8. Spark アプリケーションを起動するには、次のコマンドをノートブックに入力し、ツールバーで [実行] を選択します。

```
spark
```

しばらくすると、次のようなレスポンスが表示されます。

```
In [1]: spark
Starting Spark application

  ID      YARN Application ID  Kind  State  Spark UI  Driver log  Current session?
  --      -
  0      application_1576209965005_0001  pyspark  idle   Link      Link          ✓

SparkSession available as 'spark'.

<pyspark.sql.session.SparkSession object at 0x7f3d54913550>
```

9. 動的フレームを作成し、それに対してクエリを実行します。persons_json テーブルのカウントとスキーマを出力する次のコードをコピー、貼り付け、実行します。

```
import sys
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.transforms import *
glueContext = GlueContext(SparkContext.getOrCreate())
persons_DyF = glueContext.create_dynamic_frame.from_catalog(database="legislators",
    table_name="persons_json")
print ("Count: ", persons_DyF.count())
```

```
persons_DyF.printSchema()
```

チュートリアル: 開発エンドポイントで REPL シェルを使用する

AWS Glue では、開発エンドポイントを作成してから、REPL (Read-Evaluate-Print Loop) シェルを呼び出して PySpark コードを増分的に実行し、ETL スクリプトをデプロイする前にインタラクティブにデバッグできるようにします。

開発エンドポイントで REPL を使用するには、エンドポイントに SSH で接続するための承認が必要です。

1. ローカルコンピュータで、SSH コマンドを実行できるターミナルウィンドウを開き、編集した SSH コマンドを貼り付けます。コマンドを実行します。

開発エンドポイントで AWS Glue バージョン 1.0 (Python 3 を使用) を受け入れた場合、出力は次のようになります。

```
Python 3.6.8 (default, Aug  2 2019, 17:42:44)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-28)] on linux
Type "help", "copyright", "credits" or "license" for more information.
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/share/aws/glue/etl/jars/glue-assembly.jar!/
org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/lib/spark/jars/slf4j-log4j12-1.7.16.jar!/
org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use
setLogLevel(newLevel).
2019-09-23 22:12:23,071 WARN [Thread-5] yarn.Client (Logging.scala:logWarning(66))
- Neither spark.yarn.jars nor spark.yarn.archive is set, falling back to uploading
libraries under SPARK_HOME.
2019-09-23 22:12:26,562 WARN [Thread-5] yarn.Client (Logging.scala:logWarning(66))
- Same name resource file:/usr/lib/spark/python/lib/pyspark.zip added multiple
times to distributed cache
2019-09-23 22:12:26,580 WARN [Thread-5] yarn.Client (Logging.scala:logWarning(66))
- Same path resource file:///usr/share/aws/glue/etl/python/PyGlue.zip added
multiple times to distributed cache.
```

```

2019-09-23 22:12:26,581 WARN [Thread-5] yarn.Client (Logging.scala:logWarning(66))
  - Same path resource file:///usr/lib/spark/python/lib/py4j-src.zip added multiple
  times to distributed cache.
2019-09-23 22:12:26,581 WARN [Thread-5] yarn.Client (Logging.scala:logWarning(66))
  - Same path resource file:///usr/share/aws/glue/libs/pyspark.zip added multiple
  times to distributed cache.

```

Welcome to

```

  _ _ _ _
 /   \   \   /
 _\  \_  \_  /  \
 /_  /  \_  \  /  \  version 2.4.3
 /_  /  \_  \  /  \
 /_  /  \_  \

```

Using Python version 3.6.8 (default, Aug 2 2019 17:42:44)

SparkSession available as 'spark'.

>>>

2. ステートメント `print(spark.version)` を入力して REPL シェルが正常に動作しているかテストします。Spark のバージョンが表示されれば、REPL を使用する準備ができたこととなります。
3. シェルで次のシンプルなスクリプトを行単位で実行することができます。

```

import sys
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.transforms import *
glueContext = GlueContext(SparkContext.getOrCreate())
persons_DyF = glueContext.create_dynamic_frame.from_catalog(database="legislators",
  table_name="persons_json")
print ("Count: ", persons_DyF.count())
persons_DyF.printSchema()

```

チュートリアル: 開発エンドポイントで PyCharm Professional をセットアップする

このチュートリアルでは、ローカルマシンで実行中の [PyCharm Professional](#) Python IDE を開発エンドポイントに接続し、AWS Glue ETL (抽出、転送、およびロード) スクリプトをデプロイ前にインタラクティブに実行、デバッグ、およびテストします。チュートリアルの手順とスクリーンショットは、PyCharm Professional バージョン 2019.3 に基づいています。

開発エンドポイントをインタラクティブに接続するには、PyCharm Professional がインストールされている必要があります。無料版を使用してこれを行うことはできません。

Note

このチュートリアルでは、Amazon S3 をデータソースとして使用します。代わりに JDBC データソースを使用する場合は、Virtual Private Cloud (VPC) で開発エンドポイントを実行する必要があります。SSH を使用して VPC の開発エンドポイントに接続するには、SSH トンネルを作成する必要があります。このチュートリアルには、SSH トンネルを作成する手順は含まれていません。SSH を使用して VPC の開発エンドポイントに接続する方法については、AWS セキュリティブログの「[Securely Connect to Linux Instances Running in a Private Amazon VPC](#)」を参照してください。

トピック

- [PyCharm Professional を開発エンドポイントに接続する](#)
- [開発エンドポイントにスクリプトをデプロイする](#)
- [リモートインタープリタの設定](#)
- [開発エンドポイントでスクリプトを実行する](#)

PyCharm Professional を開発エンドポイントに接続する

1. PyCharm に legislators という名前の新しい純粋な Python プロジェクトを作成します。
2. プロジェクトに get_person_schema.py という名前のファイルを、次の内容で作成します。

```
from pyspark.context import SparkContext
from awsglue.context import GlueContext

def main():
    # Create a Glue context
    glueContext = GlueContext(SparkContext.getOrCreate())

    # Create a DynamicFrame using the 'persons_json' table
    persons_DyF =
    glueContext.create_dynamic_frame.from_catalog(database="legislators",
    table_name="persons_json")

    # Print out information about this data
    print("Count: ", persons_DyF.count())
    persons_DyF.printSchema()
```

```
if __name__ == "__main__":  
    main()
```

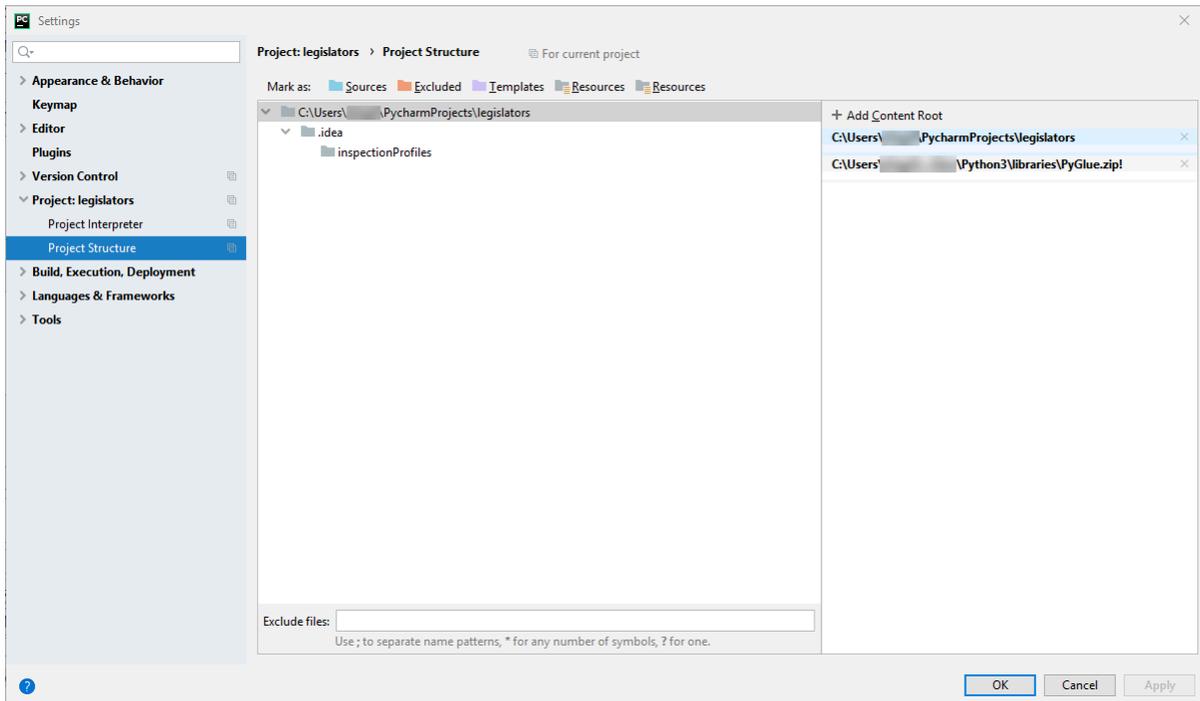
3. 次のいずれかを行います。

- AWS Glue バージョン 0.9 の場合、AWS Glue Python ライブラリファイル PyGlue.zip を <https://s3.amazonaws.com/aws-glue-jes-prod-us-east-1-assets/etl/python/PyGlue.zip> から、ローカルマシンの便利な場所にダウンロードします。
- AWS Glue バージョン 1.0 以降の場合、AWS Glue Python ライブラリファイル PyGlue.zip を <https://s3.amazonaws.com/aws-glue-jes-prod-us-east-1-assets/etl-1.0/python/PyGlue.zip> から、ローカルマシンの便利な場所にダウンロードします。

4. PyCharm のプロジェクトのコンテンツルートとして PyGlue.zip を追加します。

- PyCharm で、[File] (ファイル)、[Settings] (設定) の順に選択し、[Settings] (設定) ダイアログボックスを開きます。(Ctrl+Alt+S キーを押す方法もあります。)
- legislators プロジェクトを展開し、[Project Structure] (プロジェクト構造) を選択します。次に、右ペインで [+Add Content Root] (+ コンテンツルートの追加) を選択します。
- PyGlue.zip を保存した場所に移動して選択し、[Apply] (適用) を選択します。

[Settings] (設定) 画面は以下のようになります。



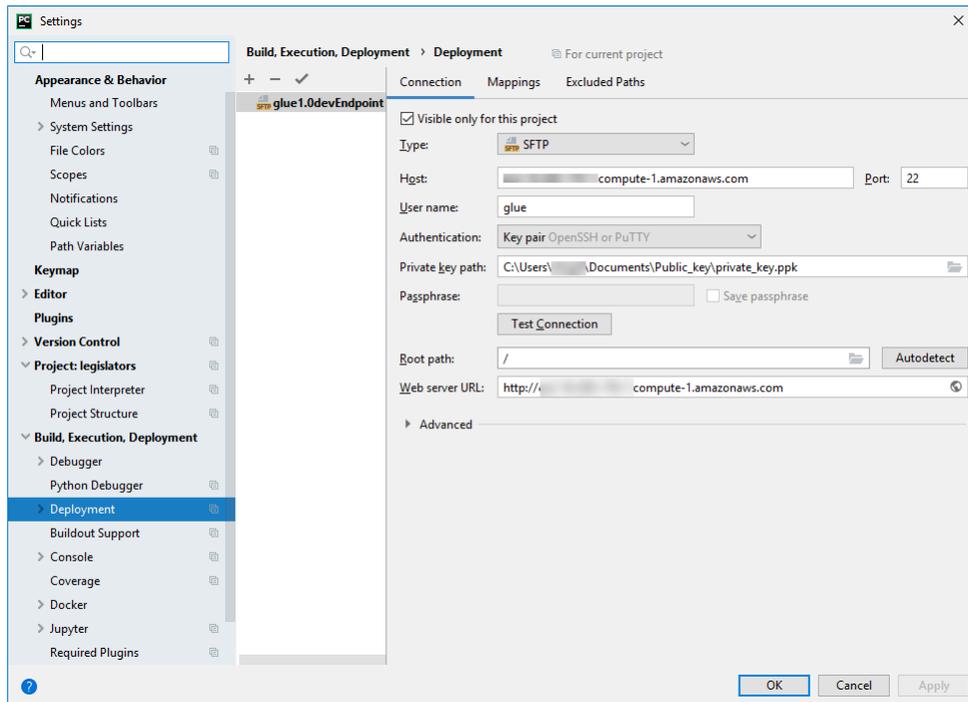
[Apply] (適用) を選択した後は、[Settings] (設定) ダイアログボックスを開いたままにします。

5. デプロイオプションを設定し、SFTP を使用してローカルスクリプトを開発エンドポイントにアップロードします (この機能は、PyCharm Professional でのみ使用できます)。
 - [Settings] (設定) ダイアログボックスで、[Build, Execution, Deployment] (ビルド、実行、デプロイ) セクションを展開します。[Deployment] (デプロイ) サブセクションを選択します。
 - 中央のペインの一番上にある [+] アイコンを選択し、新しいサーバーを追加します。[タイプ] を SFTP に設定し、名前を付けます。
 - 詳細ページにリストされているように、[SFTP ホスト] を開発エンドポイントの [パブリックアドレス] に設定します (詳細ページを表示するには、AWS Glue コンソールで開発エンドポイントの名前を選択します)。VPC で実行されている開発エンドポイントの場合、[SFTP ホスト] をホストアドレスに設定し、SSH トンネルのローカルポートを開発エンドポイントに設定します。
 - [User name] (ユーザー名) を glue に設定します。
 - [Auth type] (認証タイプ) を [Key pair (OpenSSH or Putty)] (キーペア、OpenSSH または Putty) に設定します。開発エンドポイントのプライベートキーファイルがある場所を参照し、プライベートキーファイルを設定します。PyCharm は DSA、RSA、ECDSA の OpenSSH キータイプのみをサポートし、Putty のプライベート形式のキーは受け入れないことに注意してください。ssh-keygen の最新バージョンを使用して、PyCharm が受け入れるキーペアタイプを以下のような構文で生成できます。

```
ssh-keygen -t rsa -f <key_file_name> -C "<your_email_address>"
```

- [Test connection (接続のテスト)] を選択し、接続をテストします。接続が成功したら、[Apply] (適用) を選択します。

[Settings] (設定) 画面は以下のようになります。

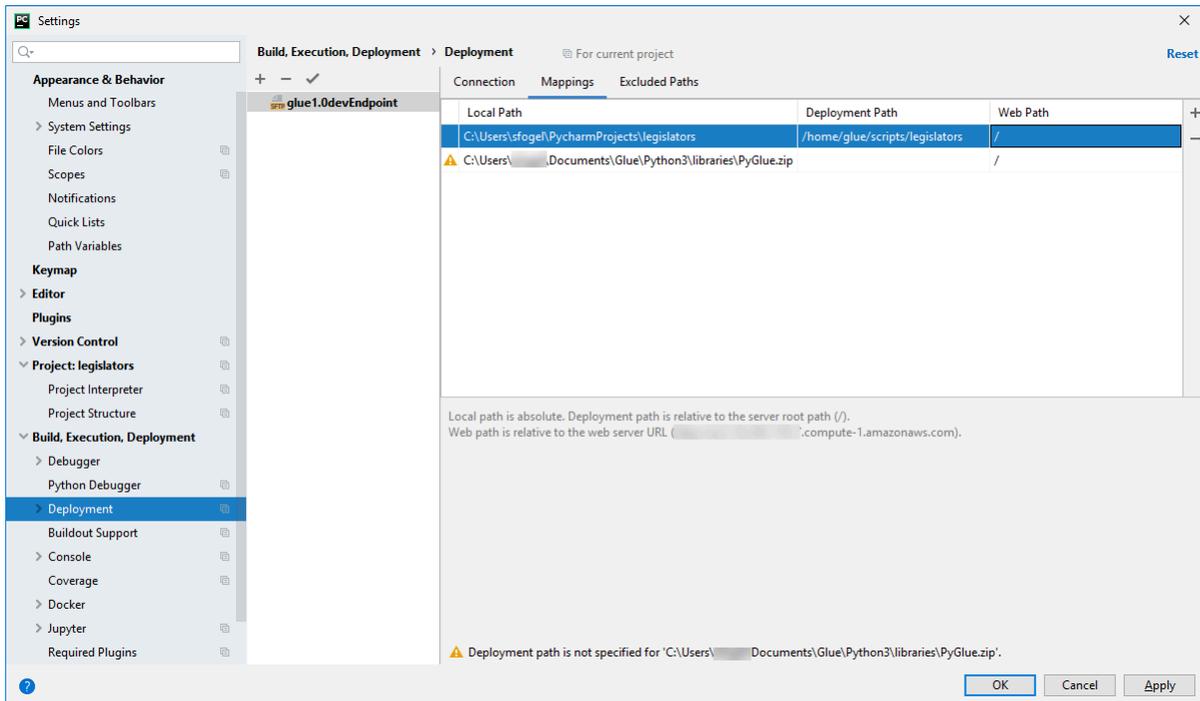


再び、[Apply] (適用) を選択した後は、[Settings] (設定) ダイアログボックスを開いたままにします。

6. ローカルディレクトリをデプロイ用のリモートディレクトリにマッピングします。

- 右のペインの [Deployment] (デプロイ) ページで、[Mappings] (マッピング) と書かれた中央上部のタブを選択します。
- [Deployment Path] (デプロイパス) 列で、プロジェクトパスのデプロイ用に /home/glue/scripts/ の下にパスを入力します。例: /home/glue/scripts/legislators。
- [Apply] を選択します。

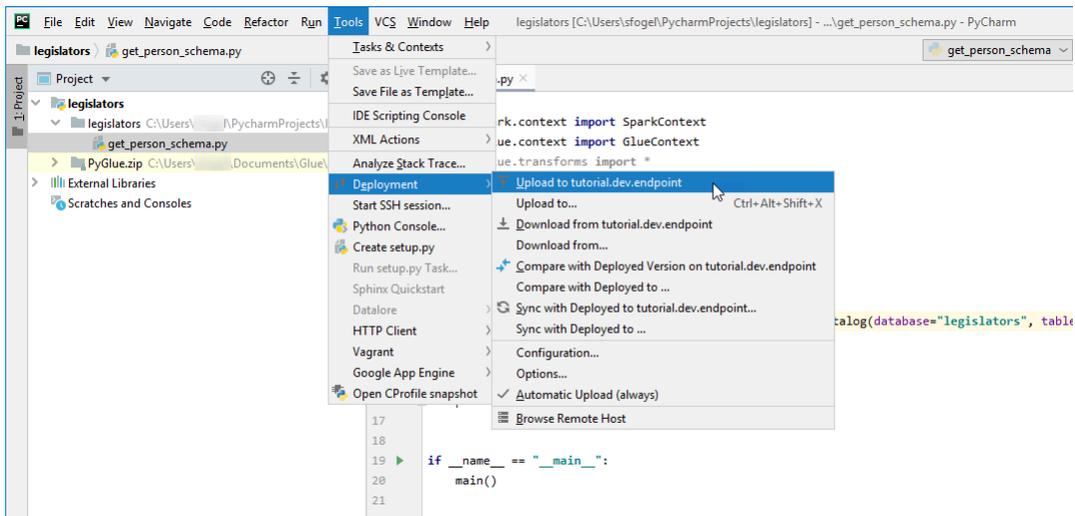
[Settings] (設定) 画面は以下のようになります。



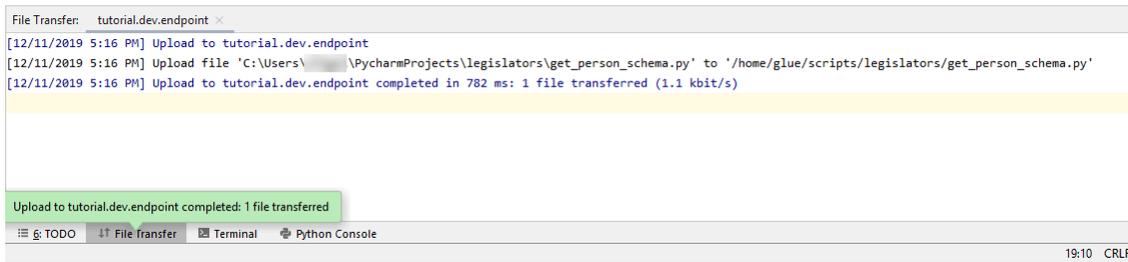
[OK] を選択し、[設定] ダイアログボックスを閉じます。

開発エンドポイントにスクリプトをデプロイする

1. [Tools (ツール)]、[Deployment (デプロイ)] の順に選択した後、次の図に示すように、開発エンドポイントをセットアップする名前を選択します。



スクリプトがデプロイされた後、画面の下部は、次のようになります。



```
File Transfer: tutorial.dev.endpoint x
[12/11/2019 5:16 PM] Upload to tutorial.dev.endpoint
[12/11/2019 5:16 PM] Upload file 'C:\Users\... \PycharmProjects\legislators\get_person_schema.py' to '/home/glue/scripts/legislators/get_person_schema.py'
[12/11/2019 5:16 PM] Upload to tutorial.dev.endpoint completed in 782 ms: 1 file transferred (1.1 kbit/s)

Upload to tutorial.dev.endpoint completed: 1 file transferred
```

2. メニューバーで、[Tools (ツール)]、[Deployment (デプロイ)]、[Automatic Upload (always) (自動アップロード (常時))] の順に選択します。[Automatic Upload (always) (自動アップロード (常時))] の横にチェックマークが表示されていることを確認します。

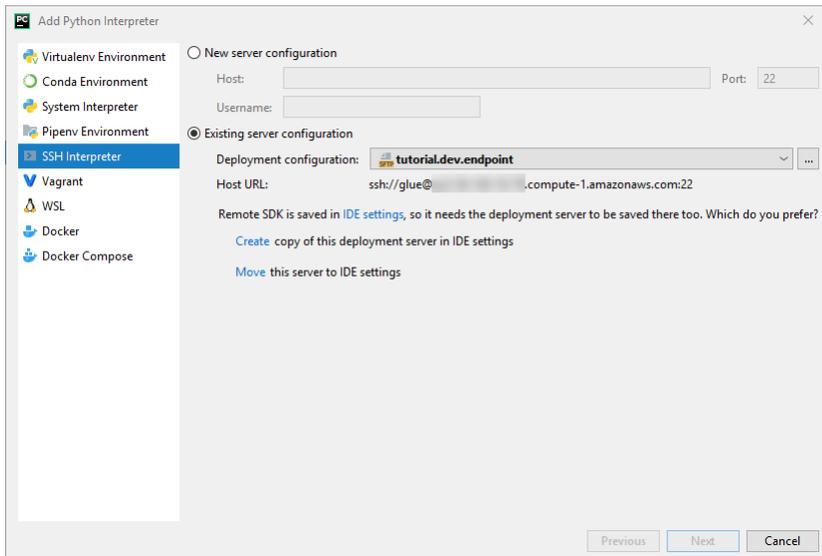
このオプションを有効にすると、PyCharm は変更されたファイルを開発エンドポイントに自動的にアップロードします。

リモートインタープリタの設定

開発エンドポイントで Python インタプリタを使用するように PyCharm を設定します。

1. [ファイル] メニューの [設定] を選択します。
2. プロジェクトの [legislators (立法者)] を展開し、[Project Interpreter (プロジェクトインタプリタ)] を選択します。
3. [Project Interpreter (プロジェクトインタプリタ)] リストの横にある歯車アイコンを選択し、[追加] を選択します。
4. [Add Python Interpreter (Python インタプリタを追加)] ダイアログボックスの左ペインで、[SSH Interpreter (SSH インタプリタ)] を選択します。
5. [Existing server configuration (既存のサーバー設定)] を選択し、[Deployment configuration (展開設定)] リストで設定を選択します。

これで、画面は以下の図のようになります。



6. [Move this server to IDE settings (このサーバーを IDE 設定に移動する)] を選択し、[次へ] を選択します。
7. [Interpreter (インタプリタ)] フィールドで、Python 2 を使用している場合は、パスを `/usr/bin/gluepython` に変更し、Python 3 を使用している場合は `/usr/bin/gluepython3` に変更します。次に、[Finish] (終了) を選択します。

開発エンドポイントでスクリプトを実行する

スクリプティングを実行するには:

- 左ペインでファイル名を右クリックし、[<filename> を実行] を選択します。
一連のメッセージの後、最終的な出力にはカウントとスキーマが表示されます。

```
Count: 1961
root
|-- family_name: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- gender: string
|-- image: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
```

```
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- sort_name: string
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- death_date: string
```

Process finished with exit code 0

これで開発エンドポイントでスクリプトをリモートでデバッグするためのセットアップができました。

高度な設定: 複数のユーザー間で開発エンドポイントを共有する

このセクションでは、一般的なユースケースで SageMaker ノートブックを使用して開発エンドポイントを活用し、複数のユーザー間で開発エンドポイントを共有する方法について説明します。

シングルテナンシーの設定

シングルテナントのユースケースでは、デベロッパーの作業をシンプルにし、リソースの競合を避けるために、各デベロッパーが、作業しているプロジェクトに適したサイズの独自の開発エンドポイントを使用することをお勧めします。これにより、ワーカータイプと DPU 数に関する決定がシンプルになり、デベロッパーおよび作業しているプロジェクトの裁量に任せられます。

複数のノートブックファイルを同時に実行しない限り、リソースの割り当てを検討する必要はありません。同時に複数のノートブックファイルでコードを実行すると、複数の Livy セッションが同時に起動します。複数の Livy セッションを同時に実行するために、Spark クラスタ設定を分離するには、マルチテナントのユースケースで導入された手順に従います。

例えば、開発エンドポイントに 10 ワーカーあり、ワーカータイプが G.1X の場合、Spark エグゼキュターは 9 つになり、各エグゼキュターは 10G のメモリを持つため、クラスター全体のエグゼキュターメモリは 90G になります。

指定されたワーカータイプに関係なく、Spark 動的リソース割り当てが有効になります。データセットが十分に大きい場合、`spark.dynamicAllocation.maxExecutors` はデフォルトで設定されていないため、Spark はすべてのエグゼキュターを単一の Livy セッションに割り当てることができます。つまり、同じ開発エンドポイント上の他の Livy セッションは、新しいエグゼキュターの起動を待つことになります。データセットが小さい場合、Spark は同時に複数の Livy セッションにエグゼキュターを割り当てることができます。

Note

さまざまなユースケースでリソースがどのように割り当てられるか、および動作を変更するための設定方法の詳細については、「[高度な設定: 複数のユーザー間で開発エンドポイントを共有する](#)」を参照してください。

マルチテナンシーの設定

Note

開発エンドポイントは、シングルテナント環境として AWS Glue ETL 環境をエミュレートするためのものであることに留意してください。マルチテナントの使用は可能ですが、これは高度なユースケースであり、ほとんどのユーザーには、開発エンドポイントごとにシングルテナンシーのパターンを維持することをお勧めします。

マルチテナントのユースケースでは、リソース割り当てを検討する必要があることがあります。重要な要因は、同時に Jupyter Notebook を使用するユーザーの数です。チームが「フォローザサン」ワークフローで作業し、各タイムゾーンに Jupyter ユーザーが 1 人しかいない場合、同時ユーザーの数は 1 人であるため、リソースの割り当てについて心配する必要はありません。ただし、ノートブックが複数のユーザー間で共有され、各ユーザーがアドホックベースでコードを送信する場合は、以下の点を考慮する必要があります。

Spark クラスターリソースを複数のユーザー間でパーティション化するには、SparkMagic の設定を使用します。Spark Magic の設定には、2 つの異なる方法があります。

(A) %%configure -f デイレクティブを使用する

ノートブックから Livy セッションごとの設定を変更したい場合は、%%configure -f デイレクティブをノートブックの段落で実行します。

例えば、5 つのエグゼキュターで Spark アプリケーションを実行する場合は、ノートブックの段落で次のコマンドを実行します。

```
%%configure -f  
{"numExecutors":5}
```

そうすると、Spark UI 上でジョブに対して実行されているエグゼキュターが 5 つだけ表示されます。

動的リソース割り当てでは、エグゼキュターの最大数を制限することをお勧めします。

```
%%configure -f  
{"conf":{"spark.dynamicAllocation.maxExecutors":"5"}}
```

(B) SparkMagic Config ファイルを変更する

SparkMagic は、[Livy API](#) に基づいて動作します。SparkMagic は、driverMemory、driverCores、executorMemory、executorCores、numExecutors、conf などの設定で Livy セッションを作成します。これらは、Spark クラスター全体で消費されるリソースの量を決定する重要な要素です。SparkMagic を使用すると、Livy に送信されるパラメータを指定するための設定ファイルを指定することができます。この [GitHub リポジトリ](#) にサンプルの設定ファイルがあります。

ノートブックからすべての Livy セッションの設定を変更したい場合は、/home/ec2-user/.sparkmagic/config.json を変更して session_config を追加します。

SageMaker ノートブックインスタンスの設定ファイルを変更するには、以下の手順に従います。

1. SageMaker ノートブックを開きます。
2. ターミナルカーネルを開きます。
3. 以下のコマンドを実行します。

```
sh-4.2$ cd .sparkmagic  
sh-4.2$ ls  
config.json logs
```

```
sh-4.2$ sudo vim config.json
```

例えば、これらの行を `/home/ec2-user/.sparkmagic/config.json` に追加し、ノートブックから Jupyter カーネルを再起動します。

```
"session_configs": {
  "conf": {
    "spark.dynamicAllocation.maxExecutors": "5"
  }
},
```

ガイドラインとベストプラクティス

このようなリソースの競合を避けるために、次のような基本的なアプローチを使用できます。

- `NumberOfWorkers` (水平方向のスケーリング) を増やし、`workerType` (垂直スケーリング) をアップグレードして、Spark クラスターを大きくする
- ユーザーあたりの割り当てリソースを削減 (Livy セッションあたりのリソースの削減)

アプローチはユースケースによって異なります。開発エンドポイントが大きく、大量のデータがない場合、Spark は動的割り当て戦略に基づいてリソースを割り当てることができるため、リソースの競合の可能性が大幅に低下します。

上記のように、Spark エグゼキュターの数、DPU (または `NumberOfWorkers`) とワーカータイプの組み合わせに基づいて自動的に計算できます。各 Spark アプリケーションは、1つのドライバーと複数のエグゼキュターを起動します。計算するには、`NumberOfWorkers = NumberOfExecutors + 1` が必要です。以下のマトリックスは、同時ユーザー数に基づいて、開発エンドポイントに必要な容量を示しています。

同時に実行するノートブックユーザーの数	ユーザーごとに割り当てる Spark エグゼキュターの数	開発エンドポイントの <code>NumberOfWorkers</code> 合計
3	5	18
10	5	60
50	5	300

ユーザーごとに割り当てるリソースを少なくしたい場合、`spark.dynamicAllocation.maxExecutors` (または `numExecutors`) は、Livy セッションパラメータとして設定する最も簡単なパラメータです。`/home/ec2-user/.sparkmagic/config.json` で以下のように設定した場合、SparkMagic は Livy セッションごとに最大 5 つのエグゼキュターを割り当てます。これは、Livy セッションごとにリソースを分離するのに役立ちます。

```
"session_configs": {
  "conf": {
    "spark.dynamicAllocation.maxExecutors": "5"
  }
},
```

18 のワーカー (G.1X) が使用されている開発エンドポイントで、同時実行するノートブックユーザー数が 3 とします。セッション設定に `spark.dynamicAllocation.maxExecutors=5` が含まれていると、各ユーザーは 1 つのドライバーと 5 つのエグゼキュターを使用できます。複数のノートブックの段落を同時に実行しても、リソースの競合は発生しません。

トレードオフ

この `"spark.dynamicAllocation.maxExecutors": "5"` というセッション設定で、リソースの競合エラーを回避することができ、同時ユーザーアクセスがあるときにリソース割り当てを待つ必要はありません。ただし、多くの空きリソースがある場合 (例えば、他の同時ユーザーがない場合) でも、Spark は Livy セッションに 5 つまでしかエグゼキュターを割り当てることができません。

その他の注意事項

ノートブックの使用を停止するときは、Jupyter カーネルを停止することをお勧めします。これにより、リソースが解放され、他のノートブックユーザーはカーネルの有効期限 (自動シャットダウン) を待たずにすぐにこれらのリソースを使用できます。

一般的な問題

ガイドラインに従っている場合でも、特定の問題が発生する可能性があります。

セッションが見つからない

Livy セッションがすでに終了しているにもかかわらず、ノートブック段落を実行しようとする、以下のメッセージが表示されます。Livy セッションをアクティブにするには、Jupyter カーネルを再起動する必要があります。Jupyter メニューで [Kernel] (カーネル) > [Restart] (再起動) を選択し、ノートブックの段落をもう一度実行します。

```
An error was encountered:  
Invalid status code '404' from http://localhost:8998/sessions/13 with error payload:  
"Session '13' not found."
```

YARN リソースが不足

Spark クラスターに新しい Livy セッションを開始するのに十分なリソースがないにもかかわらず、ノートブックの段落を実行しようとする、以下のメッセージが表示されます。多くの場合、ガイドラインに従うことでこの問題を回避できますが、この問題に直面する可能性があります。この問題を回避するには、不要なアクティブな Livy セッションがあるかどうかを確認します。不要な Livy セッションがある場合、クラスターリソースを解放するためにセッションを終了する必要があります。詳細については、次のセクションを参照ください。

```
Warning: The Spark session does not have enough YARN resources to start.  
The code failed because of a fatal error:  
    Session 16 did not start up in 60 seconds..
```

Some things to try:

- Make sure Spark has enough available resources for Jupyter to create a Spark context.
- Contact your Jupyter administrator to make sure the Spark magics library is configured correctly.
- Restart the kernel.

モニタリングとデバッグ

このセクションでは、リソースとセッションをモニタリングするための手法について説明します。

クラスターリソース割り当てのモニタリングとデバッグ

Spark UI を見て、Livy セッションごとに割り当てられているリソースの数と、ジョブで有効な Spark 設定をモニタリングできます。Spark UI をアクティブ化するには、「[Enabling the Apache Spark Web UI for Development Endpoints](#)」を参照してください。

(オプション) Spark UI のリアルタイムビューが必要な場合は、Spark クラスターで実行されている Spark 履歴サーバーに対して SSH トンネルを設定できます。

```
ssh -i <private-key.pem> -N -L 8157:<development endpoint public address>:18080  
glue@<development endpoint public address>
```

ブラウザで <http://localhost:8157> を開き、Spark UI を表示できます。

不要な Livy セッションの解放

ノートブックまたは Spark クラスターから不要な Livy セッションをすべてシャットダウンするには、以下の手順を確認してください。

(a). ノートブックから Livy セッションを終了する

Jupyter Notebook でカーネルをシャットダウンして、不要な Livy セッションを終了できます。

(b). Spark クラスターから Livy セッションを終了する

不要な Livy セッションがまだ実行されている場合は、Spark クラスターで Livy セッションをシャットダウンできます。

この手順を実行するための前提条件として、開発エンドポイントの SSH パブリックキーを設定する必要があります。

Spark クラスターにログインするには、次のコマンドを実行します。

```
$ ssh -i <private-key.pem> glue@<development endpoint public address>
```

次のコマンドを実行して、アクティブな Livy セッションを表示できます。

```
$ yarn application -list
20/09/25 06:22:21 INFO client.RMPProxy: Connecting to ResourceManager at
ip-255-1-106-206.ec2.internal/172.38.106.206:8032
Total number of applications (application-types: [] and states: [SUBMITTED, ACCEPTED,
RUNNING]):2
Application-Id Application-Name Application-Type User Queue State Final-State Progress
Tracking-URL
application_1601003432160_0005 livy-session-4 SPARK livy default RUNNING UNDEFINED 10%
http://ip-255-1-4-130.ec2.internal:41867
application_1601003432160_0004 livy-session-3 SPARK livy default RUNNING UNDEFINED 10%
http://ip-255-1-179-185.ec2.internal:33727
```

その後、次のコマンドを使用して Livy セッションをシャットダウンできます。

```
$ yarn application -kill application_1601003432160_0005
20/09/25 06:23:38 INFO client.RMPProxy: Connecting to ResourceManager at
ip-255-1-106-206.ec2.internal/255.1.106.206:8032
Killing application application_1601003432160_0005
```

```
20/09/25 06:23:39 INFO impl.YarnClientImpl: Killed application
application_1601003432160_0005
```

ノートブックの管理

Note

開発エンドポイントは、AWS Glue 2.0 より前のバージョンでのみサポートされます。ETL スクリプトを作成およびテストできるインタラクティブな環境の場合は、[AWS Glue Studio](#) でノートブックを使用します。

ノートブックを使用すると、開発エンドポイントにおける ETL (抽出、変換、ロード) スクリプトのインタラクティブな開発とテストを行えます。AWS Glue には、SageMaker Jupyter Notebook に対するインターフェイスが用意されています。AWS Glue を使用して、SageMaker ノートブックを作成および管理します。AWS Glue コンソールから SageMaker ノートブックを開くこともできます。

また、SageMaker (AWS Glue ETL ジョブではない) をサポートする AWS Glue 開発エンドポイント上の SageMaker で Apache Spark を使用できます。SageMaker Spark は、オープンソースの SageMaker 用 Apache Spark ライブラリです。詳細については、「[Amazon SageMaker で Apache Spark を使う](#)」を参照してください。

Important

AWS Glue 開発エンドポイントでの SageMaker ノートブックの管理は、以下の AWS リージョンで利用できます。

リージョン	Code
米国東部 (オハイオ)	us-east-2
米国東部 (バージニア北部)	us-east-1
米国西部 (北カリフォルニア)	us-west-1
米国西部 (オレゴン)	us-west-2
アジアパシフィック (東京)	ap-northeast-1

リージョン	Code
アジアパシフィック (ソウル)	ap-northeast-2
アジアパシフィック (ムンバイ)	ap-south-1
アジアパシフィック (シンガポール)	ap-southeast-1
アジアパシフィック (シドニー)	ap-southeast-2
カナダ (中部)	ca-central-1
欧州 (フランクフルト)	eu-central-1
欧州 (アイルランド)	eu-west-1
欧州 (ロンドン)	eu-west-2

AWS Glue Studio でビジュアル ETL ジョブを作成する

AWS Glue ジョブには、ソースデータに接続して処理し、データターゲットに書き出すスクリプトがカプセル化されています。通常、ジョブは、抽出、変換、ロード (ETL) スクリプトを実行します。ジョブは、Apache Spark および Ray ランタイム環境向けに設計されたスクリプトを実行できます。ジョブでは、汎用 Python スクリプト (Python シェルジョブ) を実行することもできます。AWS Glue トリガーでは、スケジュールまたはイベントに基づいて、またはオンデマンドでジョブを開始できます。ジョブ実行をモニタリングすると、完了ステータス、継続時間、開始時間などのランタイムメトリクスを知ることができます。

AWS Glue で生成されたスクリプトを使用することも、独自のスクリプトを使用することもできます。ソーススキーマとターゲット位置またはスキーマを指定すると、AWS Glue Studio コードジェネレーターで Apache Spark API (PySpark) スクリプトを自動的に作成できます。このスクリプトを出発点として使用し、目標に合わせて編集できます。

AWS Glue では出力ファイルを複数のデータ形式で書き込むことができます。各ジョブタイプは異なる出力形式をサポートしている場合があります。一部のデータ形式では、一般的な圧縮形式を記述できます。

AWS Glue コンソールにサインインする

AWS Glue でのジョブは、抽出、変換、ロード (ETL) 作業を実行するビジネスロジックで構成されます。コンソールの [ETL AWS Glue] セクションでジョブを作成できます。

既存のジョブを表示するには、AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。その後、の [ジョブ AWS Glue] タブを選択します。[Jobs] (ジョブ) リストは、ジョブが最後に変更されたとき各ジョブに関連付けられたスクリプトの場所、および現在のジョブのブックマークオプションを表示します。

新しいジョブの作成中またはジョブの保存後、can AWS Glue Studio を使用して、ETL ジョブを変更できます。これを行うには、ビジュアルエディタでノードを編集するか、デベロッパーモードでジョブスクリプトを編集します。ビジュアルエディタでノードを追加および削除して、より複雑な ETL ジョブを作成することもできます。

AWS Glue Studio でジョブを作成するための次のステップ

ジョブのノードを設定するには、ビジュアルジョブエディタを使用します。各ノードは、ソース位置からのデータの読み取り、データへの変換の適用などのアクションに対応しています。ジョブに追加する各ノードには、データの場所または変換に関する情報を指定するためのプロパティがあります。

ジョブを作成および管理するための次のステップは、以下のとおりです。

- [ビジュアル ETL と AWS Glue Studio](#)
- [ジョブスクリプトの表示](#)
- [ジョブのプロパティを変更する](#)
- [ジョブの保存](#)
- [ジョブの実行の開始](#)
- [最近のジョブの実行の情報を表示する](#)
- [ジョブモニタリングダッシュボードにアクセスする](#)

ビジュアル ETL と AWS Glue Studio

AWS Glue Studio ではシンプルなビジュアルインターフェイスを使用して、ETL ジョブを作成できます。新しいジョブは、[Jobs] (ジョブ) ページを使用して作成します。また、スクリプトエディタを使用して、AWS Glue Studio ETL ジョブスクリプト内のコードを直接操作することも可能です。

AWS Glue Studio または AWS Glue で作成したすべてのジョブは、[Jobs] (ジョブ) ページで確認することができます。このページでは、ジョブを表示、管理、および実行できます。

AWS Glue Studio で ETL ジョブを作成する方法の別の例については、[ブログチュートリアル](#)も参照してください。

AWS Glue Studio でのジョブの開始

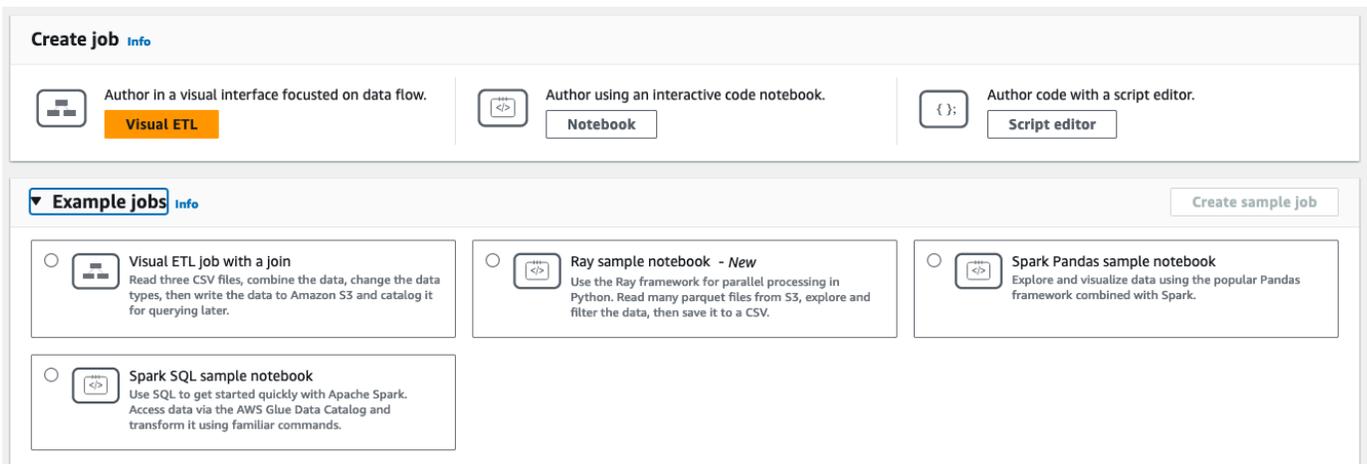
AWS Glue では、ビジュアルインターフェイスやインタラクティブなコードの Notebook を通じて、またはスクリプトエディタを使用してジョブを作成できます。オプションのいずれかをクリックしてジョブを開始することも、サンプルジョブに基づいて新しいジョブを作成することもできます。

サンプルジョブでは、選択したツールでジョブが作成されます。例えば、サンプルジョブを使用すると、CSV ファイルをカタログテーブルに結合するビジュアル ETL ジョブを作成したり、Pandas を操作するときに AWS Glue for Ray または AWS Glue for Spark を使用してインタラクティブな

コードの Notebook でジョブを作成したり、SparkSQL を使用してインタラクティブなコードの Notebook でジョブを作成したりできます。

AWS Glue Studio でゼロからのジョブの作成

1. AWS Management Console にサインインし、<https://console.aws.amazon.com/gluestudio/> **AWS Glue Studio** でコンソールを開きます。
2. ナビゲーションペインで、[ETL ジョブ] を選択します。
3. [ジョブを作成する] セクションで、ジョブの設定オプションを選択します。



ジョブを最初から作成するためのオプション:

- ビジュアル ETL – データフローに重点を置いたビジュアルインターフェイスでの作成
- インタラクティブなコードの Notebook を使用して作成 – Jupyter Notebook に基づいて、Notebook インターフェイスでジョブをインタラクティブに作成

このオプションを選択した場合、Notebook 作成セッションを作成する前に、追加情報を指定する必要があります。この情報の指定方法の詳細については、[AWS Glue Studio 中でのノートブックの使用開始](#) を参照してください。

- スクリプトエディタを使用してコードを作成 – プログラミングと ETL スクリプトの記述をよく知っている場合には、このオプションを選択して、新しい Spark ETL ジョブを作成します。エンジン (Python シェル、Ray、Spark (Python)、または Spark (Scala)) を選択します。次に、[新規に開始] または [スクリプトをアップロード] を選択し、ローカルファイルから既存のスクリプトをアップロードします。スクリプトエディタの使用を選択した場合は、ビジュアルジョブエディタを使用してジョブを設計または編集することはできません。

Spark ジョブは、AWS Glue によって管理される Apache Spark 環境で実行されます。デフォルトでは、新しいスクリプトは Python でコーディングします。新しい Scala スクリプトを作成する場合は、「[AWS Glue Studio 中の Scala スクリプトの作成および編集](#)」を参照してください。

AWS Glue Studio で、サンプルジョブからジョブの作成

サンプルジョブからジョブを作成できます。[サンプルジョブ] セクションでオプションを選択し、[サンプルジョブの作成] を選択します。いずれかのオプションから作成したサンプルジョブを手軽なテンプレートとして作業を開始できます。

1. AWS Management Console にサインインし、<https://console.aws.amazon.com/gluestudio/> [AWS Glue Studio](#) でコンソールを開きます。
2. ナビゲーションペインで、[ETL ジョブ] を選択します。
3. サンプルジョブからジョブを作成するためのオプションを次の中から選択します。
 - 複数のソースを結合するビジュアル ETL ジョブ – 3 つの CSV ファイルを読み取り、データを結合し、データ型を変更し、データを Amazon S3 に書き込み、後でクエリできるようにカタログ化します。
 - Pandas を使用した Spark Notebook – 一般的な Pandas フレームワークと Spark を組み合わせてデータを探索および可視化します。
 - SQL を使用した Spark Notebook – SQL を使用して Apache Spark をすぐに開始できます。AWS Glue データカタログからデータにアクセスし、使い慣れたコマンドを使用してデータを変換します。
4. [サンプルジョブの作成] を選択します。

ジョブエディタの機能

ジョブエディタでは、ジョブの作成と編集のための次の機能が提供されています。

- ジョブの視覚的な図。各ジョブタスクのノード (データを読み取るデータソースノード、データを修正する変換ノード、データを書き込むデータターゲットノード) が示されます。

ジョブ図では、各ノードのプロパティを表示および設定できます。また、各ノードのスキーマとサンプルデータを表示することもできます。これらの機能により、ジョブを実行することなく、ジョブでデータが正しい方法で変更および変換されていることを確認できます。

- [Script viewing and editing] (スクリプトの表示と編集) タブ。ジョブ用に生成されたコードを変更できます。
- [Job details] (ジョブの詳細) タブ。さまざまな設定を行い、AWS Glue ETL ジョブの実行環境をカスタマイズできます。
- [Runs] (実行) タブ。現在実行中のジョブや以前実行したジョブの表示、ジョブの実行のステータスの表示、ジョブの実行用のログへのアクセスを行えます。
- [データ品質] タブでは、データ品質ルールをジョブに適用できます。
- [Schedules] (スケジュール) タブ。ジョブの開始時間を設定したり、ジョブの実行を定期的に設定できます。
- バージョン管理タブ。ジョブで使用する Git サービスを設定できます。

ビジュアルジョブエディタでスキーマのプレビューを使用する

ジョブの作成中または編集中に、[Output schema] (出力スキーマ) タブをクリックして、データのスキーマを表示できます。

スキーマを表示するには、ジョブエディタにデータソースへのアクセス許可が必要です。エディタの [Job details] (ジョブの詳細) タブ、またはノードの [Output schema] (出力スキーマ) タブで IAM ロールを指定できます。データソースへのアクセスに必要なすべてのアクセス許可が IAM ロールにある場合は、ノードの [Output schema] (出力スキーマ) タブのスキーマを表示できます。

ビジュアルジョブエディタでデータのプレビューを使用する

データのプレビューを使用することで、ジョブを繰り返し実行することなく、データのサンプルを使用してジョブを作成およびテストできます。データプレビューを使用すると、次を実行できます。

- IAM ロールをテストして、データソースまたはデータターゲットにアクセスできることを確認します。
- 変換により、意図した方法でデータが変更されていることを確認します。例えば、フィルター変換を使用する場合、フィルターでデータの適切なサブセットが選択されていることを確認できます。
- データを確認する。データセットに複数のタイプの値を持つ列が含まれている場合、データのプレビューには、これらの列のタプルのリストが表示されます。各タプルには、データ型とその値が含まれています。

ジョブの作成中または編集中に、[Data preview] タブをクリックして、データのサンプルを表示できます。ジョブでロールが既に設定されている場合、またはアカウントでデフォルトの IAM ロールが設定されている場合、新しいデータプレビューセッションが自動的に開始されます。ロールが事前に設定されていない場合は、ロールを選択してセッションを開始できます。

Data preview | **Output schema** 🗨️ 🗨️

Start a data preview session

IAM role
To start a data preview session, choose an IAM role for this job. Changing the role will end an existing data preview session.

Admin ▼
No description available.

[Create IAM role.](#)

▶ **Additional Settings**

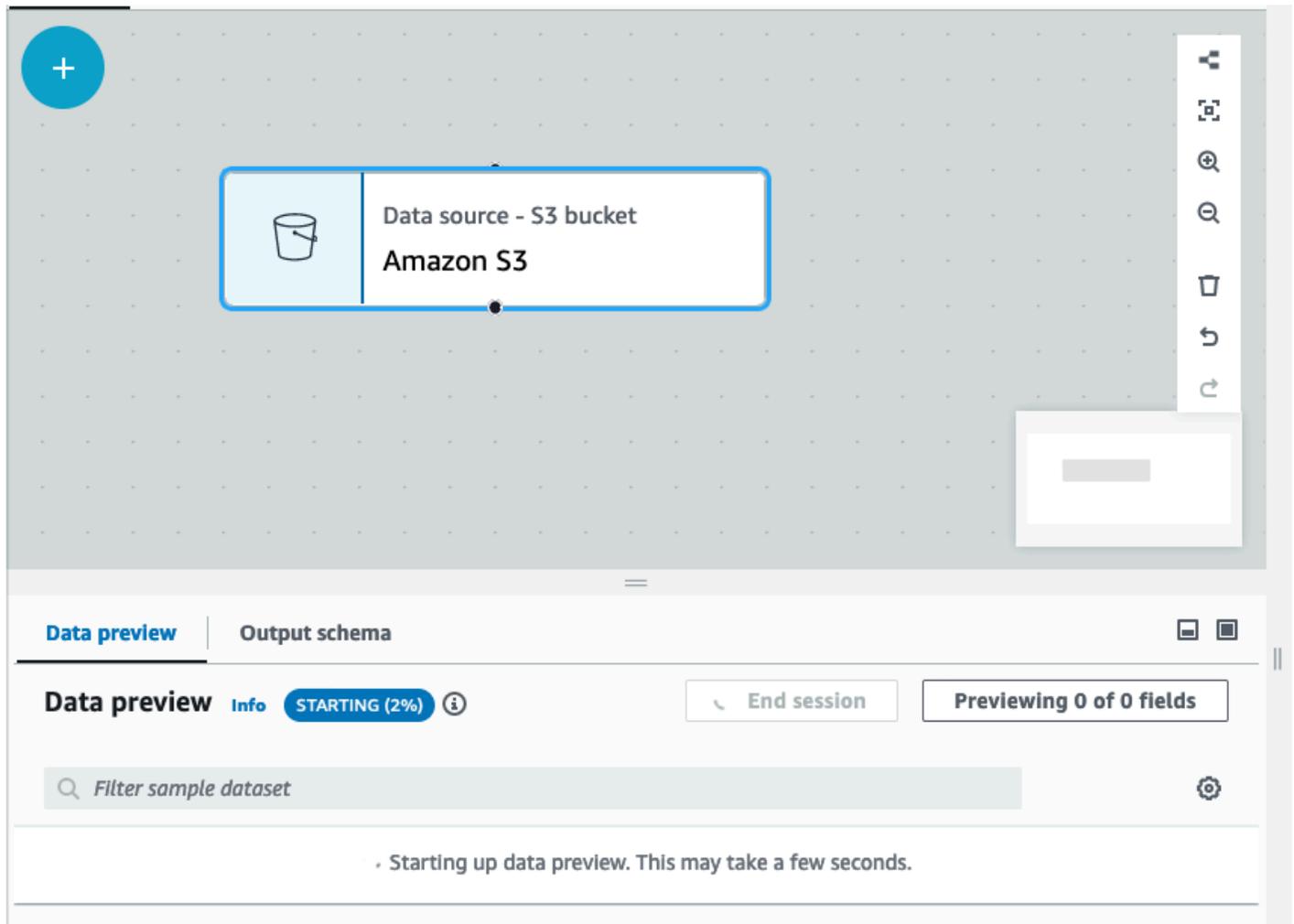
Start session

Note

データプレビューセッション用に選択したロールは、ジョブにも使用されます。

情報アイコンをクリックすると、セッションのステータスと進行状況、およびセッションの詳細を確認できます。

セッションの準備が完了すると、AWS Glue Studio は選択したノードのためにデータをロードします。進捗に伴う [完了率 (%)] を表示できます。



ビジュアルジョブを作成する際に、[出力スキーマ] タブで [セッションからスキーマを推測] を切り替えると、AWS Glue Studio は選択したノードのスキーマを自動的に更新します。

The screenshot displays the AWS Glue console interface for configuring a job. A job graph at the top shows a 'Transform - SQL Query' node selected. The configuration panel on the right includes a dropdown for 'Choose one or more parent node', an 'Amazon S3' input source, and an 'SQL alias' of 'myDataSource'. Below this, the 'SQL query' field contains the statement: `select firstname, lastname, title from myDataSource`. At the bottom, the 'Output schema' section shows a table with columns: Key, Data type, firstname (string), lastname (string), and title (string).

データプレビューの設定を行うには:

設定アイコン (歯車の記号) を選択して、データのプレビューの設定を行えます。これらの設定は、ジョブ図のすべてのノードに適用されます。次のようにできます。

- テキストを 1 行から次の行に折り返す場合に選択します。デフォルトでは、このオプションは有効になっています
- 行数を変更する (デフォルトは 200)
- IAM ロールを選択するか、または、必要であれば新しいロールを作成
- ジョブの作成時に新しいセッションを自動的に開始することを選択します。これにより、ジョブの作成時に新しいインタラクティブセッションがプロビジョニングされます。この設定はアカウントレベルで適用されます。一度設定すると、ジョブの編集時にアカウント内のすべてのユーザーに適用されます。
- スキーマを自動的に推論する場合に選択します。出力スキーマは、選択したノードのために自動的に推論されます
- AWS Glue ライブラリを自動的にインポートする場合に選択します。これは、セッションの再起動を必要とする新しい変換を追加する際に、データプレビューが新しいセッションを再起動するのを防ぐため、有益です。

Preferences ✕

Wrap lines
Enable to wrap lines of table cell content, disable to truncate text.

Number of rows
Enter the amount of entries to sample from the dataset.

IAM role
To start a data preview session, choose an IAM role for this job. Changing the role will end an existing data preview session.

Admin
No description available. ▼

[Create IAM role.](#) 

Automatically start data preview sessions
Data preview will automatically start new interactive sessions when entering the visual job editor enabling you to preview data more efficiently.
⚠ This setting applies to all users in your account.

Infer schema from session
Output schemas will be automatically inferred based on the result of the datapreview execution

Automatically import glue libraries
Some ETL transform require extra libraries to be imported in the datapreview session, enabling this option will automatically import them to your sessions in order to prevent session from restarting during your job authoring. Note: the IAM role require read permission to Glue S3 bucket to prevent failures.

Cancel **Confirm**

追加機能には次の機能が含まれます。

- [Previewing x of y fields] (X/Y フィールドのプレビュー) ボタンをクリックして、プレビューする列 (フィールド) を選択できます。デフォルトの設定を使用してデータをプレビューすると、ジョブ

エディタにはデータセットの最初の 5 列が表示されます。これを、[show all] (すべて表示) または [none] (なし) (非推奨) に変更できます。

- データのプレビューウィンドウは、水平方向と垂直方向の両方にスクロールします。
- 最大化ボタンを使用して、[Data preview] タブをジョブグラフ画面全体に展開し、データおよびデータ構造をより詳しく表示できます。同様に、最小化ボタンを使用して [Data preview] タブを最小化します。ハンドルペインをつかんで上にドラッグして [Data preview] タブを展開することもできます。

The screenshot shows the AWS Glue console interface. At the top, there are tabs for 'Visual', 'Script', 'Job details', 'Runs', 'Data quality New', 'Schedules', and 'Version Control'. The main area displays a job graph with a 'Data source - S3 bucket Amazon S3' and a 'Data target - Snowflake'. Below the graph, the 'Data preview' tab is active, showing a table of venue data. A red box highlights the maximize and close buttons in the top right corner of the preview window.

Data preview (200) Info **READY** ⓘ

Filter sample dataset

venueid	venue name	venue city	venue state	venue seats
1	Toyota Park	Bridgeview	IL	0
2	Columbus Crew Stadium	Columbus	OH	0
3	RFK Stadium	Washington	DC	0

- [セッションの終了] を使用してデータプレビューを停止します。セッションを停止すると、新しい IAM ロールを選択し、新しいセッションを自動的に開始したり、スキーマを推測したり、AWS Glue ライブラリをインポートしたりするための追加設定 (オン/オフなど) を設定して、セッションを再開できます。

データのプレビューを使用する際の制限事項

データのプレビューを使用する際、次の制約または制限事項が発生する場合があります。

- 初めて [Data preview] (データのプレビュー) タブを選択する際、IAM ロールを選択する必要があります。このロールには、そのデータおよびデータのプレビューの作成に必要なその他のリソースへのアクセス許可が必要です。
- IAM ロールを指定した後、データを表示できるようになるまでには時間がかかります。データが 1 GB 未満のデータセットの場合、最大で 1 分かかることがあります。大きなデータセットがある場合は、パーティションを使用してロード時間を改善する必要があります。Amazon S3 からデータを直接ロードすると、パフォーマンスが最も良くなります。
- 非常に大きなデータセットがあり、データのプレビュー用のデータのクエリに 15 分以上かかる場合、リクエストはタイムアウトします。データプレビューには 30 分の IDLE タイムアウトがあります。これを軽減するために、データセットのサイズを小さくして、データプレビューを使用できます。
- デフォルトでは、[Data preview] タブには最初の 50 列が表示されます。列にデータの値がない場合は、表示するデータはありませんというメッセージが表示されます。サンプリングされる行数を増やすか、別の列を選択してデータの値を表示できます。
- 現在、データのプレビューは、ストリーミングデータソースまたはカスタムコネクタを使用するデータソースではサポートされていません。
- 1 つのノードでのエラーは、ジョブ全体に影響します。いずれかのノードでデータのプレビューのエラーがある場合、エラーは修正されるまですべてのノードに表示されます。
- ジョブのデータソースを変更する場合、そのデータソースの子ノードは、新しいスキーマに合わせて更新する必要があります。例えば、列を変更する ApplyMapping ノードがあり、その列が代替データソースに存在しない場合、ApplyMapping 変換ノードを更新する必要があります。
- SQL クエリ変換ノードの [Data preview] (データのプレビュー) タブを表示する際に SQL クエリで誤ったフィールド名が使用されている場合、[Data preview] (データのプレビュー) タブにエラーが表示されます。

スクリプトコードの生成

ビジュアルエディタを使用してジョブを作成する場合、ETL コードが自動的に生成されます。AWS Glue Studio は機能的で完全なジョブスクリプトを作成し、Amazon S3 の場所に保存します。

AWS Glue Studio によって生成されるコードには 2 つの形式があります。: オリジナルバージョン、またはクラシックバージョン、およびより新しいストリームラインされたバージョンです。デフォル

トでは、新しいコードジェネレーターがジョブスクリプトの作成に使用されます。ジョブスクリプトは、[Generate classic script] トグルボタンを選択することにより、[Script] タブ上のクラシックコードジェネレーターを使用して生成できます。

生成されたコードの新しいバージョンでは、いくつかの違いを含みます。:

- 大きなコメントブロックは、もうスクリプトに追加されなくなりました。
- コード内の出力のデータ構造は、ビジュアルエディタで指定したノード名を使用します。クラススクリプトでは、出力のデータ構造は単に DataSource0, DataSource1, Transform0, Transform1, DataSink0, DataSink1 などと名前が付けられます。
- 長いコマンドは複数の行に分割されるため、ページ全体をスクロールしてコマンド全体を表示する必要がなくなります。

AWS Glue Studio 中の新機能には新しいバージョンのコード生成が必要で、クラシックコードスクリプトでは動作しません。これらのジョブを実行しようとする、これらのジョブを更新するように求められます。

AWS Glue マネージドデータ変換ノードの編集

AWS Glue Studio には、次の 2 種類の変換が用意されています。

- AWS Glue ネイティブ変換 - すべてのユーザーが利用でき、AWS Glue によって管理されます。
- カスタムビジュアル変換 - 自分のカスタムビジュアル変換をアップロードして、AWS Glue Studio で使用できます。

AWS Glue マネージドデータ変換ノード

AWS Glue Studio には、データを処理するために使用できる一式の組み込み変換が用意されています。データは、ジョブ図内のあるノードから DynamicFrame と呼ばれるデータ構造の別のノードに渡されます。これは、Apache Spark SQL DataFrame の拡張機能です。

事前設定されたジョブ図では、データソースとデータターゲットノード間で、スキーマ変更変換ノードを使用できます。この変換ノードを設定すると、データを修正したり、他の変換を使用できます。

AWS Glue Studio を用いて、次の built-in 変換が利用可能です。

- [ChangeSchema](#): データソースのデータプロパティキーを、データターゲットのデータプロパティキーにマッピングします。キーの名前を変更したり、データ型を変更したり、データセットから削除するキーを選択できます。
- [SelectFields](#): 保持したいデータプロパティキーを選択します。
- [DropFields](#): 削除したいデータプロパティキーを選択します。
- [RenameField](#): 単一のデータプロパティキーの名前を変更します。
- [Spigot](#): Amazon S3 バケットにデータのサンプルを書き込みます。
- [Join](#): 指定したデータプロパティキー上の比較フレーズを使用して、2つのデータセットを1つに結合します。結合タイプは、内部結合、外部結合、左結合、右結合、左半結合、左反結合を使用できます。
- [Union](#): 同じスキーマを持つ複数のデータソースの行を結合します。
- [SplitFields](#): データプロパティキーを2つの DynamicFrames に分割します。出力は DynamicFrames のコレクションです。一方は選択したデータプロパティキー、他方は残っている方のデータプロパティキーを持ちます。
- [SelectFromCollection](#): DynamicFrame のコレクションから1つの DynamicFrames を選択します。出力は選択された DynamicFrame です。
- [\[FillMissingValues\]](#): 欠落値があるデータセットの記録を配置し、代入によって決定された提案値を用いた新しいフィールドを追加します。
- [Filter](#): フィルター条件に基づいて、データセットを2つに分割します。
- [\[Drop Null Fields\]](#): 列のすべての値が「null」の場合、データセットから列を削除します。
- [\[Drop Duplicates\]](#): 行全体を照合するか、キーを指定するかを選択して、データソースから行を削除します。
- [SQL](#): SparkSQL コードをテキスト入力フィールドに入力して、データを変換するために SQL クエリを使用します。出力は、単一の DynamicFrame です。
- [Aggregate](#): 選択したフィールドと行で計算 (平均、合計、最小、最大など) を実行し、計算された値を使って新しいフィールドを作成します。
- [Flatten](#): Struct 内のフィールドを最上位のフィールドに抽出します。
- [UUID](#): 各行にユニバーサル意識別子を含む列を追加します。
- [Identifier](#): 各行に数値識別子を含む列を追加します。
- [To timestamp](#): 列をタイムスタンプタイプに変換します。
- [Format timestamp](#): タイムスタンプ列をフォーマットされた文字列に変換します。

- [Conditional Router transform](#): 受信データに複数の条件を適用します。受信データの行はそれぞれグループフィルター条件により評価され、対応するグループに取り入れて処理されます。
- [\[列の連結\] 変換](#): オプションのスペーサーを使用して、他の列の値を使った新しい文字列の列を作成します。
- [\[文字列の分割\] 変換](#): 正規表現を使用して文字列をトークンの配列に分割し、分割方法を定義します。
- [\[配列から列へ\] 変換](#): 配列タイプの列の一部またはすべての要素を抽出し、新しい列に追加します。
- [\[現在のタイムスタンプを追加\] 変換](#): データが処理された時刻で行をマークします。これは、監査目的やデータパイプラインでのレイテンシーの追跡に役立ちます。
- [\[行から列へのピボット\] 変換](#): 選択した列の固有の値を回転させて数値列を集約し、新しい列にします。複数の列を選択した場合、値が連結されて新しい列に名前が付けられます。
- [\[列から行へのピボット解除\] 変換](#): 列を新しい列の値に変換し、固有の値ごとに行を生成します。
- [\[オートバランス処理\] 変換](#): データをワーカー間でより適切に再配分します。これは、データのバランスが取れていない場合や、ソースから取得したデータでは十分に並行処理ができない場合に役立ちます。
- [\[派生列\] 変換](#): 数式または SQL 式に基づいて新しい列を定義します。これらの式では、データ内の他の列および定数やリテラルを使用できます。
- [\[ルックアップ\] 変換](#): キーがデータ内の定義済みのルックアップ列と一致する場合、定義済みのカタログテーブルから列を追加します。
- [\[配列またはマップを行に分解\] 変換](#): ネストされたデータ構造から値を抽出し、操作しやすい個々の行に値を追加します。
- [レコードマッチング変換](#): 既存のレコードマッチング機械学習データ分類変換を呼び出します。
- [null 行変換を削除](#): すべての列が null または空の行をデータセットから削除します。
- [JSON 列変換の解析](#): JSON データを含む文字列の列を解析し、JSON がオブジェクトか配列かに応じて、それぞれを構造体または配列の列に変換します。
- [JSON パス変換の抽出](#): JSON 文字列の列から新しい列を抽出します。
- [正規表現から文字列の断片を抽出](#): 正規表現を使用して文字列の断片を抽出し、そこから新しい列を作成します。正規表現グループを使用する場合は、複数の列を作成します。
- [Custom transform](#): カスタム transforms を使用するために、テキスト入力フィールドにコードを入力します。出力は、DynamicFrames のコレクションです。

AWS Glue Studio でのデータ準備レシピの使用

AWS Glue Studio では、ビジュアルワークフローの AWS Glue DataBrew レシピを使用できます。これにより、お客様の AWS Glue DataBrew レシピをその他の AWS Glue Studio ノードと共に AWS Glue ジョブで実行できます。

DataBrew では、一連のデータ変換ステップをレシピといいます。DataBrew レシピには読み取り済みデータの変換方法は規定されていますが、データを読み取る場所と方法、およびデータを書き込む場所と方法は記載されていません。これは、AWS Glue Studio のソースノードとターゲットノードで設定します。レシピの詳細については、「[Creating and using AWS Glue DataBrew recipes](#)」を参照してください。

データ準備レシピノードは、リソースパネルから使用できます。データ準備レシピノードは、データソースノードであるか別の変換ノードであるかに関らず、ビジュアルワークフロー内の別のノードに接続できます。AWS Glue DataBrew レシピとバージョンを選択すると、レシピに適用されたステップがノードのプロパティタブに表示されます。

前提条件

- AWS Glue DataBrew で作成された AWS Glue DataBrew レシピがあります。
- 以下のセクションで説明されているような IAM アクセス許可が必要です。

AWS Glue DataBrew に対する IAM アクセス許可

このトピックでは、IAM 管理者がデータ準備レシピ変換の AWS Identity and Access Management (IAM) ポリシーで使用できるアクションとリソースを理解するための情報について説明します。

AWS Glue のセキュリティに関する追加情報については、「[Access Management](#)」を参照してください。

次の表は、データ準備レシピ変換を使用するために特定の操作を実行するユーザーが必要とするアクセス許可を一覧表示しています。

データ準備レシピ変換アクション

[アクション]	説明
<code>databrew:ListRecipes</code>	AWS Glue DataBrew レシピを取得するアクセス許可を付与します。

[アクション]	説明
databrew:ListRecipeVersions	AWS Glue DataBrew レシピバージョンを取得するアクセス許可を付与します。
databrew:DescribeRecipe	AWS Glue DataBrew レシピの説明を取得するアクセス許可を付与します。

この機能にアクセスするために使用するロールには、複数の AWS Glue DataBrew を許可するポリシーが必要です。このためには、必要なアクションを含む AWSGlueConsoleFullAccess ポリシーを使用するか、次のインラインポリシーをロールに追加します。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "databrew:ListRecipes",
        "databrew:ListRecipeVersions",
        "databrew:DescribeRecipe"
      ],
      "Resource": [
        "*"
      ]
    }
  ]
}
```

データ準備レシピ変換を使用するには、アクセス許可ポリシーに IAM:PassRole アクションを追加する必要があります。

追加で必要な許可

[アクション]	説明
iam:PassRole	承認済みのロールをユーザーが渡すことを許可するアクセス許可を、IAM に付与します。

これらのアクセス許可がないと、次のエラーが発生します。

```
"errorCode": "AccessDenied"
"errorMessage": "User: arn:aws:sts::account_id:assumed-role/AWSGlueServiceRole is not
authorized to perform: iam:PassRole on resource: arn:aws:iam::account_id:role/service-
role/AWSGlueServiceRole
because no identity-based policy allows the iam:PassRole action"
```

制限事項

- すべての AWS Glue DataBrew レシピが AWS Glue でサポートされているわけではありません。一部のレシピは AWS Glue Studio では実行できません。
- Union および Join 変換を使用するレシピはサポートされていませんが、AWS Glue Studio には、データ準備レシピノードの前後で使用できる「Union」および「Join」変換ノードが既に存在します。
- データ準備レシピノードは AWS Glue バージョン 4.0 以降のジョブでサポートされています。このバージョンは、データ準備レシピノードがジョブに追加されると自動選択されます。
- データ準備レシピノードには Python が必要です。Python は、データ準備レシピノードがジョブに追加されると自動設定されます。
- [データのプレビュー] を使用する場合、[データ準備レシピ] ノードをジョブに追加した後、データのプレビュー セッションを再起動する必要があります。

AWS Glue Studio で AWS Glue DataBrew レシピを使用する方法

AWS Glue Studio で AWS Glue DataBrew レシピを使用するには、最初に AWS Glue DataBrew でレシピを作成します。使用するレシピがある場合は、このステップを省略できます。

AWS Glue DataBrew で AWS Glue DataBrew レシピを作成するには

1. AWS Glue DataBrew でレシピを作成します。詳細については、「[Getting started with AWS Glue DataBrew](#)」を参照してください。
2. レシピを保存します。
3. レシピを公開します。これにより、レシピがバージョン 1.0 として公開されます。

AWS Glue Studio でデータ準備レシピノードを使用するには

ビジュアル ETL ジョブで複数のデータ準備レシピノードを使用できます。このためには、以下のステップに従って別のデータ準備レシピノードをジョブに追加し、データ準備レシピノードを追加します。例えば、ワークフローは次のパターンに従う場合があります。

- データソース 1 > レシピ 1 > 出力 1
- データソース 2 > レシピ 2 > 出力 2
- 出力 1、出力 2 > Join

1. データソースを使用して、AWS Glue Studio で AWS Glue ジョブを開始します。
2. データ準備レシピノードをデータソースに追加します。
3. 検索フィールドにレシピ名を入力して、レシピを名前で絞り込みます。
4. 公開済みバージョンを選択します。公開済みバージョンのみを使用可能です。
5. 必要に応じてその他の変換ノードを追加してジョブ作成を終了し、データターゲットノードを追加してジョブ出力を保存します。
6. 必要に応じてジョブに名前を付けたり、割り当てられた容量を調整したりして、[ジョブの詳細] タブで必要な設定変更を行い、ジョブを保存します。
7. [アクション] ドロップダウン メニューで [実行] を選択し、ジョブを実行します。

データソースが Amazon S3 でデータ形式が CSV の場合にスキーマを変更するには

CSV ファイル内のすべての列が AWS Glue Studio で最初に文字列データ型としてロードされる場合は、列データ型が AWS Glue DataBrew レシピの残りのステップと互換性があることを確認する必要があります。

AWS Glue DataBrew レシピは、読み取り済みデータの変換方法のみを規定しています。データ読み取りの場所と方法については説明していません。

1. マルチステップレシピノードの前にスキーマ変更ノードを追加します。
2. スキーマ変更ノードをクリックし、必要に応じて「列の変換」で新しいデータ型を選択して、AWS Glue DataBrew の列のデータ型と同じになるようにスキーマを変更します。

Transform

Name
Change Schema

Node parents
Choose which nodes will provide inputs for this one.
Choose one or more parent node

S3 bucket X
S3 - DataSource

Change Schema (Apply mapping)

Source key	Target key	Data type	Drop
col0	col0	string ▼	<input type="checkbox"/>
col1	col1	string ▼	<input type="checkbox"/>
col2	col2	string ▼	<input type="checkbox"/>
col3	col3	string ▼	<input type="checkbox"/>
col4	col4	string ▼	<input type="checkbox"/>
col5	col5	string ▼	<input type="checkbox"/>
col6	col6	string ▼	<input type="checkbox"/>
col7	col7	string ▼	<input type="checkbox"/>
col8	col8	string ▼	<input type="checkbox"/>

データソースがヘッダーレスの場合にスキーマを変更するには

AWS Glue DataBrew レシピは、読み取り済みデータの変換方法のみを規定しています。データ読み取りの場所と方法については説明していません。

AWS Glue Studio でヘッダーレスデータセットをロードする場合、デフォルトのヘッダー名は AWS Glue DataBrew でロードされるものとは異なります。

1. ETL ジョブで、データ準備レシピノードの前にスキーマ変更ノードを追加します。
2. スキーマ変更ノードを選択し、列名を AWS Glue DataBrew レシピで使用されているものと同じ名前に変更します。

スキーマの変更によるデータプロパティキーの再マッピング

スキーマ変更変換を使用して、ソースデータのプロパティキーをターゲットデータに適切な設定に再マッピングします。スキーマ変更変換ノードでは、次のことができます。

- 複数のデータプロパティキーの名前の変更。
- 新しいデータ型がサポートされており、2つのデータ型間の変換パスがある場合、データプロパティキーのデータ型の変更。
- 削除するデータプロパティキーの指定によるデータプロパティキーのサブセットの選択。

必要に応じて、ジョブ図にスキーマ変更ノードを追加することもできます。例えば、追加のデータソースを変更したり、Join 変換を実行したりできます。

10 進データ型でスキーマの変更を使用する

10 進データ型でスキーマ変更変換を使用する場合、スキーマ変更変換は精度をデフォルト値 (10,2) に変更します。これを変更してユースケースの精度を設定するには、SQL クエリ変換を使用して特定の精度で列をキャストできます。

例えば、「」という名前 DecimalCol の入力列が 10 進数で、特定の精度 (18,6) OutputDecimalCol で「」という名前の出力列に再マップする場合、次のようになります。

1. スキーマ変更変換後に後続の SQL クエリ変換を追加します。
2. SQL クエリ変換では、SQL クエリを使用して、再マッピングされた列を希望の精度にキャストします。SQL クエリは次のようになります。

```
SELECT col1, col2, CAST(DecimalCol AS DECIMAL(18,6)) AS OutputDecimalCol
FROM __THIS__
```

上記の SQL クエリでは、次の操作を行います。

- `col1` と `col2` は、データ内の他の列であり、変更なしでパススルーします。
- `DecimalCol` は入力データの元の列名です。

- `CAST(DecimalCol AS DECIMAL(18,6))` は、DecimalCol` を 18 桁、小数点以下 6 桁の精度の 10 進数型にキャストします。
- `AS OutputDecimalCol` は、キャストされた列の名前を OutputDecimalCol` に変更します。

SQL クエリ変換を使用すると、スキーマ変更変換で設定されたデフォルトの精度を上書きし、10 進数の列を希望の精度に明示的にキャストできます。このアプローチにより、変更スキーマ変換を活用してデータの名前を変更および再構築しながら、その後の SQL クエリ変換を通じて 10 進数列の精度要件を処理できます。

変更スキーマ変換をジョブに追加する

Note

スキーマ変更変換は、大文字と小文字を区別しません。

スキーマ変更変換ノードをジョブ図に追加するには

1. (オプション) リソースパネルを開いて、[スキーマを変更する] を選択し、必要に応じてジョブ図に新しい変換を追加します。
2. ノードのプロパティパネルで、ジョブ図にノード名を入力します。ノードの親がまだ選択されていない場合、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. ノードのプロパティパネルで、[変換] タブを選択します。
4. その後、入カスキーマを変更します。
 - データプロパティキーの名前を変更するには、キーの新しい名前を [Target key] (ターゲットキー) フィールドに入力します。
 - データプロパティキーのデータ型を変更するには、[Data type] (データ型) リストから、新しいデータ型を選択します。
 - ターゲットスキーマからデータプロパティキーを削除するには、そのキーの [Drop] (削除) チェックボックスをオンにします。
5. (オプション) 変換ノードのプロパティを設定した後、ノードの詳細パネルの [Output schema] (出カスキーマ) タブを選択して、データ用に変更されたスキーマを表示できます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定

するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。

- (オプション) ノードおよび変換のプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、変更されたデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

Drop Duplicates を使用する

Drop Duplicates 変換では、2 つのオプションによってデータソースから行を削除します。完全に一致した重複する行を削除するか、照合するフィールドを選択して、選択したフィールドに基づいて重複する行のみを削除するかのどちらかを選択できます。

例えば、次のデータセットには重複する行があり、行内の一部の値が一致していたり異なっていたりします。また、いくつかの行は、別の行とすべての値が完全に一致しています。

行	名前	Email(メール)	年齢	都道府県	注記
1	Joy	joy@gmail	33	NY	
2	Tim	tim@gmail	45	OH	
3	Rose	rose@gmail	23	NJ	
4	Tim	tim@gmail	42	OH	
5	Rose	rose@gmail	23	NJ	
6	Tim	tim@gmail	42	OH	これは重複した行であり、行 4 とすべての値が完全に一致しています
7	Rose	rose@gmail	23	NJ	これは重複した行であり、行 5 とすべて

行	名前	Email(メール)	年齢	都道府県	注記
					の値が完全に一致していません

行全体で照合することを選択した場合、行 6 と 7 はデータセットから削除されます。データセットは次のようになります。

行	名前	Email(メール)	年齢	都道府県
1	Joy	joy@gmail	33	NY
2	Tim	tim@gmail	45	OH
3	Rose	rose@gmail	23	NJ
4	Tim	tim@gmail	42	OH
5	Rose	rose@gmail	23	NJ

キーを指定することを選択した場合、「名前」と「Email(メール)」が一致する行を削除するように選択できます。これにより、データセットの「重複する行」をより細かく制御できます。「名前」と「Email(メール)」を指定すると、データセットは次のようになります。

行	名前	Email(メール)	年齢	都道府県
1	Joy	joy@gmail	33	NY
2	Tim	tim@gmail	45	OH
3	Rose	rose@gmail	23	NJ

留意すべき点

- 行の重複を認識するため、値では大文字と小文字が区別されます。行のすべての値は、大文字と小文字も一致する必要があります。これは、選択したオプション (行全体での一致/指定キーでの一致) のどちらにも当てはまりません。
- すべての値は文字列として読み込まれます。
- Drop Duplicates 変換は、Spark の dropDuplicates コマンドを使用します。
- Drop Duplicates 変換を使用すると、最初の行は保持され、他の行は削除されます。
- Drop Duplicates 変換は、データフレームのスキーマを変更しません。キーの指定を選択した場合、生成されたデータフレームにはすべてのフィールドが保持されます。

SelectFields を使用して大多数のデータプロパティキーを削除する

SelectFields 変換を使用して、データセットからデータプロパティキーのサブセットを作成できます。保持するデータプロパティキーを指定します。残りはデータセットから削除されます。

Note

SelectFields 変換では、大文字と小文字が区別されます。大文字と小文字を区別しない方法でフィールドを選択する必要がある場合は、ApplyMapping を使用します。

SelectFields 変換ノードをジョブ図に追加するには

1. (オプション) リソースパネルを開いて、[SelectFields] を選択し、必要に応じてジョブ図に新しい変換を追加します。
2. [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入力ソースとして使用するノードを選択します。
3. ノードの詳細パネルで [Transform] (変換) タブを選択します。
4. 見出し [SelectFields] の下にある、保持するデータセット内のデータプロパティキーを選択します。選択されていないデータプロパティキーは、データセットから削除されます。

列見出し [Field] (フィールド) の横にあるチェックボックスをオンにして、データセット内のすべてのデータプロパティキーを自動的に選択することもできます。その後、個々のデータプロパティキーの選択を解除して、それらをデータセットから削除できます。

5. (オプション) 変換ノードのプロパティを設定した後、ノードの詳細パネルの [Output schema] (出力スキーマ) タブを選択して、データ用に変更されたスキーマを表示できます。ジョブ内の

任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。

6. (オプション) ノードおよび変換のプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、変更されたデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

DropFields を使用して大多数のデータプロパティキーを保持する

DropFields 変換を使用して、データセットからデータプロパティキーのサブセットを作成できます。データセットから削除するデータプロパティキーを指定します。残りのキーは保持されます。

Note

DropFields 変換では、大文字と小文字が区別されます。大文字と小文字を区別しない方法でフィールドを選択する必要がある場合は、スキーマ変更を使用します。

ジョブ図に DropFields 変換ノードを追加するには

1. (オプション) リソースパネルを開いて、[DropFields] を選択し、必要に応じてジョブ図に新しい変換を追加します。
2. [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. ノードの詳細パネルで [Transform] (変換) タブを選択します。
4. 見出し [DropFields] の下で、データソースから削除するデータプロパティキーを選択します。

列見出し [Field] (フィールド) の横にあるチェックボックスをオンにして、データセット内のすべてのデータプロパティキーを自動的に選択することもできます。その後、個々のデータプロパティキーの選択を解除して、それらをデータセット内に保持できます。

5. (オプション) 変換ノードのプロパティを設定した後、ノードの詳細パネルの [Output schema] (出力スキーマ) タブを選択して、データ用に変更されたスキーマを表示できます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定

するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。

- (オプション) ノードおよび変換のプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、変更されたデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

データセット内のフィールドの名前を変更する

RenameField 変換を使用して、データセット内の個々のプロパティキーの名前を変更します。

Note

RenameField 変換では、大文字と小文字が区別されます。大文字と小文字を区別しないで変換する必要がある場合は、ApplyMapping を使用します。

Tip

スキーマ変更変換では、1 回の変換でデータセット内に存在する複数のデータプロパティキーの名前を変更できます。

RenameField 変換ノードをジョブ図に追加するには

- (オプション) リソースパネルを開いて、[RenameField] を選択し、必要に応じてジョブ図に新しい変換を追加します。
- [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入力ソースとして使用するノードを選択します。
- [Transform] (変換) タブを選択します。
- 見出し [Data field] (データフィールド) の下で、ソースデータからプロパティキーを選択し、その後 [New field name] (新しいフィールド名) フィールドに新しい名前を入力します。
- (オプション) 変換ノードのプロパティを設定した後、ノードの詳細パネルの [Output schema] (出力スキーマ) タブを選択して、データ用に変更されたスキーマを表示できます。ジョブ内の

任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。

6. (オプション) ノードおよび変換のプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、変更されたデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

Spigot を使用してデータセットをサンプリングする

ジョブで実行される変換をテストするには、データのサンプルを取得して、変換が意図したとおりに機能することを確認します。Spigot 変換では、データセットから Amazon S3 バケットの JSON ファイルにレコードのサブセットが書き出されます。データのサンプリングには、ファイルの最初からの特定のレコード数、またはレコードの選択に使用される確率係数を使用します。

Spigot 変換ノードをジョブ図に追加するには

1. (オプション) リソースパネルを開いて、[スピゴット] を選択し、必要に応じてジョブ図に新しい変換を追加します。
2. [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. ノードの詳細パネルで [Transform] (変換) タブを選択します。
4. Amazon S3 パスを入力するか、[Browse S3] (S3 をブラウズ) を選択して、Amazon S3 内の場所を選択します。これは、ジョブによりデータサンプルを含む JSON ファイルが書き込まれる場所です。
5. サンプリング方法に関する情報を入力します。データセットの最初から書き込むレコード数の値、および任意のレコードを選択する確率のしきい値 (最大値が 1 の 10 進値として入力) を指定できます。

例えば、データセットから最初の 50 レコードを書き込むには、レコード数を 50、確率のしきい値を 1 (100%) に設定します。

データセットの結合

結合 変換を使用すると、2つのデータセットを1つに結合できます。比較する各データセットのスキーマのキー名を指定します。出力 `DynamicFrame` には、キーが結合条件を満たす行が含まれます。結合条件を満たす各データセットの行は、いずれかのデータセット内で見つかったすべての列を含む出力 `DynamicFrame` の単一の行に結合されます。

Join (結合) 変換ノードをジョブ図に追加するには

1. 使用可能なデータソースが1つしかない場合は、新しいデータソースノードをジョブ図に追加する必要があります。
2. 結合するソースノードを1つ選択します。リソースパネルを開いて、[Join] を選択し、ジョブ図に新しい変換を追加します。
3. [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。
4. [Node properties] (ノードのプロパティ) タブの、見出し [Node parents] (ノードの親) の下で、結合の入力を提供するデータセットが2つ存在するように親ノードを追加します。親ノードは、データソースノードまたは変換ノードです。

Note

結合すると、親ノードを2つだけ持つことができます。

5. [Transform] (変換) タブを選択します。

競合するキー名があることを示すメッセージが表示された場合は、次のいずれかの操作を行います。

- [Resolve it] (解決する) を選択して、ApplyMapping 変換ノードをジョブ図に自動的に追加します。ApplyMapping ノードにより、他のデータセットのキーと同じ名前を持つデータセット内のキーにプレフィックスが追加されます。例えば、デフォルト値 **right** を使用している場合、左側のデータセットと同じキー名を持つ右側のデータセットのキー名は `(right)key name` に変更されます。
 - ジョブ図に変換ノードを早期に手動で追加して、競合するキーを削除するか、名前を変更できます。
6. [Join type] (結合タイプ) のリストで、結合のタイプを選択します。
 - Inner join (内部結合): 結合条件に基づくすべての一致に対して、両方のデータセットの列を含む行を返します。結合条件を満たさない行は返されません。

- Left join (左結合): 左側のデータセットのすべての行と、結合条件を満たす右側のデータセットの行の結合です。
 - Right join (右結合): 右側のデータセットのすべての行と、結合条件を満たす左側のデータセットの行の結合です。
 - Outer join (外部結合): 両方のデータセットのすべての行の結合です。
 - Left semi join (左半結合): 右側のデータセットで結合条件に基づく一致がある、左側のデータセットのすべての行の結合です。
 - Left anti join (左反結合): 右側のデータセットで結合条件に基づく一致がない、左側のデータセットのすべての行の結合です。
7. [Transform] (変換) タブで、見出し [Join conditions] (結合条件) の下にある、[Add condition] (条件の追加) を選択します。比較する各データセットからプロパティキーを選択します。比較演算子の左側にあるプロパティキーを左側のデータセット、右側にあるプロパティキーを右側のデータセットと呼びます。

より複雑な結合条件の場合、[Add condition] (条件の追加) を複数回クリックして、他の一致するキーを追加できます。誤って条件を追加した場合は、削除アイコン () をクリックして削除できます。

8. (オプション) 変換ノードのプロパティを設定した後、ノードの詳細パネルの [Output schema] (出力スキーマ) タブを選択して、データ用に変更されたスキーマを表示できます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。
9. (オプション) ノードおよび変換のプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、変更されたデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

結合出力スキーマの例では、次のプロパティキーを持つ 2 つのデータセット間の結合を検討します。

```
Left: {id, dept, hire_date, salary, employment_status}
Right: {id, first_name, last_name, hire_date, title}
```

結合は、= 比較演算子を使用して、id および hire_date キーで一致するように設定されます。

両方のデータセットに id および hire_date キーが含まれているため、[Resolve it] (解決する) を選択して、右側のデータセットのキーにプレフィックス **right** が自動で追加されるようにします。

出力スキーマのキーは次のようになります。

```
{id, dept, hire_date, salary, employment_status,  
(right)id, first_name, last_name, (right)hire_date, title}
```

Union を使用して行を結合する

Union 変換ノードは、同じスキーマを持つ複数のデータソースの行を結合する場合に使用します。

Union 変換には次の 2 種類があります。

1. ALL – ALL を適用すると、Union の結果として、重複する行は削除されません。
2. DISTINCT – DISTINCT を適用すると、Union の結果として、重複する行は削除されます。

Union と Join

行を結合するには Union を使用します。列を結合するには Join を使用します。

Visual ETL キャンバスでの Union 変換の使用

1. Union 変換を実行するには、複数のデータソースを追加します。データソースを追加するには、リソースパネルを開き、[ソース] タブからデータソースを選択します。Union 変換を使用する前に、Union の対象となるすべてのデータソースのスキーマと構造が同じであることを確認する必要があります。
2. Union 変換を使用して結合したいデータソースが 2 つ以上ある場合は、キャンバスに追加して Union 変換を作成します。キャンバス上のリソースパネルを開き、「Union」を検索します。リソースパネルの [変換] タブを選択し、Union 変換が見つかるまで下にスクロールして、[Union] を選択することもできます。
3. ジョブキャンバスの Union ノードを選択します。ノードプロパティウィンドウで、Union 変換に接続する親ノードを選択します。
4. AWS Glue は互換性をチェックして、Union 変換がすべてのデータソースに適用できることを確認します。データソースのスキーマが同じ場合、実行できます。データソースのスキーマが同じでない場合は、無効であるというエラーメッセージが表示されます。「この Union の入力ス

スキーマは一致していません。ApplyMapping を使用してスキーマを一致させることを検討してください」これを修正するには、[ApplyMapping を使用] を選択します。

5. Union の種類を選択します。

1. All – デフォルトでは、Union の種類として All が選択されています。これにより、データを結合したときに、重複する行は重複したままになります。
2. Distinct – データを結合したときに、重複する行を削除したい場合は、Distinct を選択します。

SplitFields を使用してデータセットを 2 つに分割する

SplitFields 変換を使用すると、入力データセット内のデータプロパティキーをいくつか選択し、それらを 1 つのデータセットに配置したり、選択されていないキーを別のデータセットに配置できます。この変換からの出力は、DynamicFrames のコレクションです。

Note

出力をターゲットの場所へ送信するには、SelectFromCollection 変換を使用して、DynamicFrames のコレクションを単一の DynamicFrame に変換する必要があります。

SplitFields 変換では、大文字と小文字が区別されます。大文字と小文字を区別しないプロパティキー名が必要な場合は、ApplyMapping 変換を親ノードとして追加します。

SplitFields 変換ノードをジョブ図に追加するには

1. (オプション) リソースパネルを開いて、[SplitFields] を選択し、必要に応じてジョブ図に新しい変換を追加します。
2. [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入力ソースとして使用するノードを選択します。
3. [Transform] (変換) タブを選択します。
4. 最初のデータセットに入れるプロパティキーを選択します。選択しなかったキーは、2 番目のデータセットに配置されます。
5. (オプション) 変換ノードのプロパティを設定した後、ノードの詳細パネルの [Output schema] (出力スキーマ) タブを選択して、データ用に変更されたスキーマを表示できます。ジョブ内の

任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。

6. (オプション) ノードおよび変換のプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、変更されたデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。
7. SelectFromCollection 変換ノードを設定して、結果のデータセットを処理します。

SelectFromCollection 変換の概要

SplitFields のように、単一ではなく、複数のデータセットが出力される変換もあります。SelectFromCollection 変換では、データセットのコレクション (DynamicFrames の配列) から、1 つのデータセット (DynamicFrame) が選択されます。変換の出力は、選択された DynamicFrame です。

この変換は、次のような DynamicFrames のコレクションを作成する変換を使用した後に使用する必要があります。

- カスタムコード変換
- SplitFields

変換後に SelectFromCollection 変換ノードをジョブ図に追加しないと、ジョブでエラーが発生します。

この変換の親ノードは、DynamicFrames のコレクションを返すノードである必要があります。この変換ノードに、Join (結合) 変換などの単一の DynamicFrame を返す親ノードを選択すると、ジョブによりエラーが返されます。

同様に、ジョブ図の SelectFromCollection ノードを、入力として単一の DynamicFrame を想定している変換の親ノードとして使用する場合、ジョブによりエラーが返されます。

Node parents

Select which node(s) will provide inputs for this one

Select parents

Split Fields
SplitFields - Transform Parent node Split Fields outputs a collection, but node Drop Fields does not accept a collection.

SelectFromCollection を使用して保持するデータセットを選択する

SelectFromCollection 変換を使用して、DynamicFrames のコレクションを単一の DynamicFrame に変換できます。

SelectFromCollection 変換ノードをジョブ図に追加するには

1. (オプション) リソースパネルを開いて、[SelectFromCollection] を選択し、必要に応じてジョブ図に新しい変換を追加します。
2. [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. [Transform] (変換) タブを選択します。
4. 見出し [Frame index] (フレームインデックス)の下で、DynamicFrames のコレクションから選択する DynamicFrame に対応する配列のインデックス番号を選択します。

例えば、この変換の親ノードが SplitFields 変換である場合、そのノードの [Output schema] (出力スキーマ) タブで各 DynamicFrame のスキーマを表示できます。[出力 2] のスキーマに関連付けられている DynamicFrame を保持する場合、フレームインデックスの値は **1** を選択します。これは、リストの 2 番目の値です。

選択した DynamicFrame だけが出力に含まれます。

5. (オプション) 変換ノードのプロパティを設定した後、ノードの詳細パネルの [Output schema] (出力スキーマ) タブを選択して、データ用に変更されたスキーマを表示できます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。
6. (オプション) ノードおよび変換のプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、変更されたデータセットをプレビューできま

す。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

データセット内の欠落値を検索して埋める

FillMissingValues 変換を使用して、データセット内に欠落値があるレコードを検索し、補完により決定する値を持つ新しいフィールドを追加します。入力データセットは、欠落値を決定する機械学習 (ML) モデルのトレーニングに使用されます。増分のデータセットを使用する場合、増分の各セットが ML モデルのトレーニングデータとして使用されるため、結果はそれほど正確ではない可能性があります。

ジョブ図で FillMissingValues 変換ノードを使用するには

1. (オプション) リソースパネルを開いて、[FillMissingValues] を選択し、必要に応じてジョブ図に新しい変換を追加します。
2. [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。ノードの親がまだ選択されていない場合、[Node parents] (ノードの親) リストから、変換の入力ソースとして使用するノードを選択します。
3. [Transform] (変換) タブを選択します。
4. [Data field] (データフィールド) で、欠落値を分析するソースデータから、列名またはフィールド名を選択します。
5. (オプション) [New field name] (新しいフィールド名) フィールドに、フィールドの名前を入力します。このフィールドは、分析したフィールドの推定置換値を保持する各レコードに追加されます。分析したフィールドに欠落値がない場合、その分析したフィールドの値が新しいフィールドにコピーされます。

フィールドの名前を指定しない場合、デフォルトの名前は、分析した列に `_filled` を追加した名前になります。例えば、[Data field] (データフィールド) に「Age」と入力し、[New field name] (新しいフィールド名) に値を指定しない場合、**Age_filled** という名前の新しいフィールドが各レコードに追加されます。

6. (オプション) 変換ノードのプロパティを設定した後、ノードの詳細パネルの [Output schema] (出力スキーマ) タブを選択して、データ用に変更されたスキーマを表示できます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。

7. (オプション) ノードおよび変換のプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、変更されたデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

データセット内のキーのフィルタリング

フィルター変換を使用して、正規表現に基づいて入力データセットのレコードをフィルタリングすることで、新しいデータセットを作成できます。フィルター条件を満たさない行は、出力から削除されます。

- 文字列のデータ型の場合、キー値と指定された文字列が一致する行をフィルタリングできます。
- 数値のデータ型の場合、比較演算子 $<$ 、 $>$ 、 $=$ 、 \neq 、 \leq 、 \geq を使用してキー値と指定された値を比較することで、行をフィルタリングできます。

複数のフィルター条件を指定した場合、デフォルトでは AND 演算子を使用して結果が結合されますが、代わりに OR を使用することもできます。

フィルター変換では、大文字と小文字が区別されます。大文字と小文字を区別しないプロパティキー名が必要な場合は、ApplyMapping 変換を親ノードとして追加します。

フィルター変換ノードをジョブ図に追加するには

1. (オプション) リソースパネルを開いて、[フィルター] を選択し、必要に応じてジョブ図に新しい変換を追加します。
2. [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入力ソースとして使用するノードを選択します。
3. [Transform] (変換) タブを選択します。
4. [Global AND] または [Global OR] を選択します。これにより、複数のフィルター条件の結合方法が決定します。すべての条件は、AND または OR オペレーションで結合されます。フィルター条件が 1 つしかない場合は、どちらかを選択できます。
5. [Filter condition] (フィルター条件) セクションで [Add condition] (条件の追加) ボタンをクリックして、フィルター条件を追加します。

[Key] (キー) フィールドで、データセットからプロパティキー名を選択します。[Operation] (オペレーション) フィールドで、比較演算子を選択します。[Value] (値) フィールドに、比較値を入力します。フィルター条件のいくつかの例を次に示します。

- `year >= 2018`
- `State matches 'CA*'`

文字列の値をフィルタリングする場合は、比較値に、ジョブのプロパティ (Python または Scala) で選択したスクリプト言語と一致する正規表現の形式が使用されていることを確認します。

6. 必要に応じて、他のフィルター条件を追加します。
7. (オプション) 変換ノードのプロパティを設定した後、ノードの詳細パネルの [Output schema] (出力スキーマ) タブを選択して、データ用に変更されたスキーマを表示できます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。
8. (オプション) ノードおよび変換のプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、変更されたデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

NULL 値を持つフィールドを削除する DropNullFields の使用

フィールド内のすべての値が 'null' の場合、[DropNullFields]transform を使用し、データセットからフィールドを削除します。デフォルトでは、AWS Glue Studio は null オブジェクトを認識しますが、空の文字列、“null”の列、-1 の整数、または 0 などのプレースホルダなどでは、自動的に「null」として認識されません。

DropNullFields を使用するには

1. ジョブダイアグラムに DropNullFields ノードを追加します。
2. [Node properties]タブで、NULL 値を表す追加の値を選択します。値を選択しないか、すべての値を選択するかを選択できます。

Node properties | **Transform** | Output schema | Data preview 

DropNullFields [Info](#)

Remove fields or columns where all the values are the null objects.

Choose additional values that represent a null value below.

- Empty String (" " or "")
- "null" String
- 1 Integer

Add custom null values

Specify custom null values by entering the value and choosing the datatype.



- 空の文字列(" " or "") - 空の文字列を含むフィールドは削除されます。
 - "null string"-単語「null」の文字列を含むフィールドは削除されます。
 - a-1 (負の 1) integerを含む-1 整数 - フィールドは削除されます。
3. 必要に応じて、カスタム NULL 値を指定することもできます。これらは、データセットに一意である可能性がある NULL 値です。カスタム NULL 値を追加するには、[Add new value]を選択します。
 4. カスタム Null 値を入力します。たとえば、ゼロを指定したり、データセット内の NULL を表すために使用されている任意の値を指定できます。
 5. ドロップダウンフィールドでデータタイプを選択します。データ型は、文字列または整数のいずれかです。

Note

カスタムの NULL 値とそのデータ型は、フィールドが NULL 値として認識され、フィールドを削除するために、正確に一致する必要があります。カスタムの NULL 値のみが一致しているがデータ型が一致しない部分的な一致については、フィールドが削除されません。

SQL クエリを使用してデータを変換する

[SQL] 変換を使用して、SQL クエリの形式で独自の変換を記述できます。

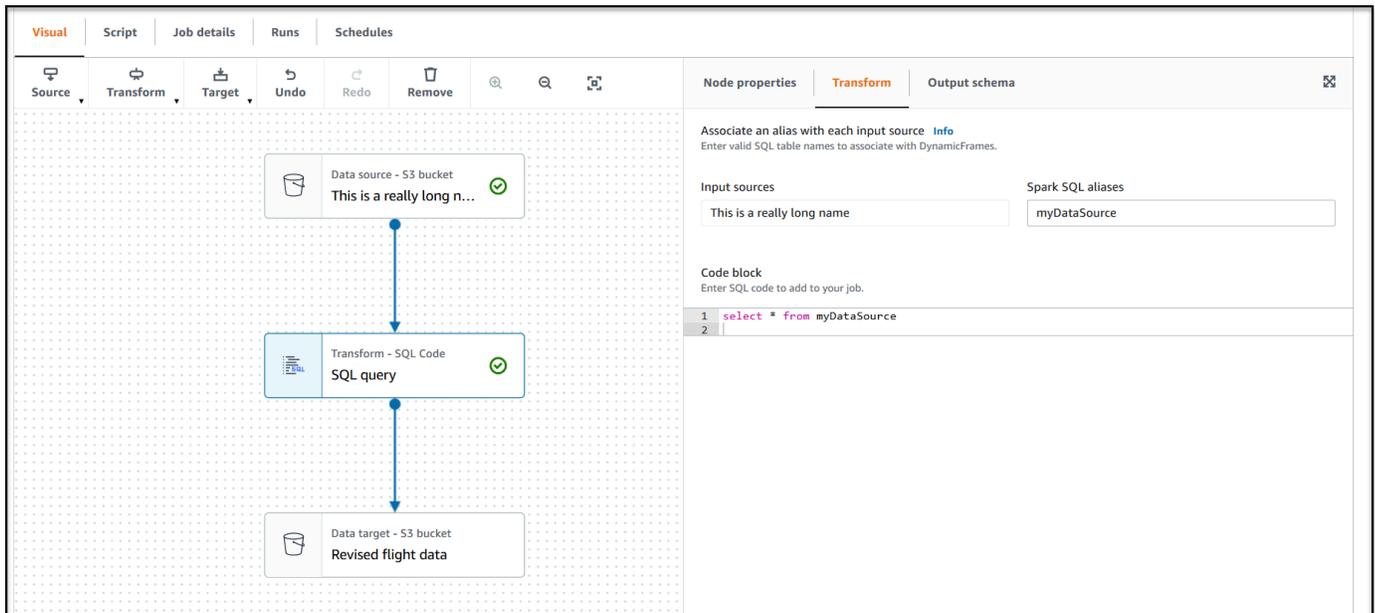
SQL 変換ノードでは、入力として複数のデータセットを持つことができますが、出力として生成されるデータセットは 1 つだけです。これには、Apache SparkSQL クエリを入力するテキストフィールドが含まれています。入力として使用する各データセットにエイリアスを割り当てることで、SQL クエリを簡単に実行できます。SQL 構文の詳細については、[Spark SQL ドキュメント](#)を参照してください。

Note

VPC にあるデータソースで Spark SQL 変換を使用する場合は、AWS Glue VPC エンドポイントを、データソースを含む VPC に追加します。開発エンドポイントの設定の詳細については、AWS Glue デベロッパーガイドの「[Adding a Development Endpoint](#)」、「[Setting Up Your Environment for Development Endpoints](#)」、「[Accessing Your Development Endpoint](#)」を参照してください。

ジョブ図で SQL 変換ノードを使用するには

1. (オプション) 必要な場合、ジョブ図に変換ノードを追加します。ノードタイプを [Spark SQL] に設定します。
2. [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。ノードの親がまだ選択されていない場合、または SQL 変換に複数の入力が必要な場合は、[Node parents] (ノードの親) リストから、変換の入力ソースとして使用するノードを選択します。必要な場合、他の親ノードを追加します。
3. ノードの詳細パネルで [Transform] (変換) タブを選択します。
4. SQL クエリのソースデータセットは、各ノードの [Name] (名前) フィールドで指定する名前により識別されます。これらの名前を使用しない場合、または名前が SQL クエリに適していない場合は、各データセットに名前を関連付けることができます。コンソールでは、MyDataSource などのデフォルトのエイリアスが提供されています。



例えば、SQL 変換ノードの親ノードの名前が `Rename Org PK field` の場合、このデータセットに `org_table` という名前を関連付けることができます。その後、このエイリアスをノード名の代わりに SQL クエリで使用できます。

5. 見出し [Code block] (コードブロック) の下にある テキスト入力フィールドで、SQL クエリを貼り付けるか、入力します。テキストフィールドには、SQL 構文のハイライトとキーワードの候補が表示されます。
6. SQL 変換ノードを選択した状態で、[Output schema] (出カスキーマ) タブを選択してから、[Edit] (編集) を選択します。SQL クエリの出カフィールドを示す、列とデータ型を指定します。

そのページの [Output schema] (出カスキーマ) セクションにある次のアクションを使用して、スキーマを指定します。

- 列の名前を変更するには、その列の [Key] (キー) テキストボックス ([field] (フィールド) または [property key] (プロパティキー) とともに示されます) にカーソルを置き、新しい名前を入力します。
- 列のデータ型を変更するには、ドロップダウンリストから列の新しいデータ型を選択します。
- トップレベルの新しい列をスキーマに追加するには、[Overflow] (オーバーフロー) (...) ボタンを選択して、[Add root key] (ルートキーの追加) をクリックします。新しい列がスキーマの先頭に追加されます。

- スキーマから列を削除するには、キー名の右端にある削除アイコン (□) をクリックします。
7. 出力スキーマの指定が完了したら、[Apply] (適用) を選択して変更を保存し、スキーマエディタを終了します。変更を保存しない場合は、[Cancel] (キャンセル) を選択して、スキーマエディタを編集します。
 8. (オプション) ノードおよび変換のプロパティを設定した後、ノードの詳細パネルの [Data preview] (データのプレビュー) タブを選択して、変更されたデータセットをプレビューできます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。この機能を使用するには費用がかかり、IAM ロールを指定するとすぐに請求が開始します。

[集計] を使用して選択したフィールドに対して集約計算を実行する

Aggregate transform を使用するには

1. Aggregate ノードをジョブダイアグラムに追加します。
2. Node properties タブ上で、ドロップダウンフィールド (オプション) を選択して、一緒にグループ化するフィールドを選択します。一度に複数のフィールドを選択するか、検索バーに入力してフィールド名を検索できます。

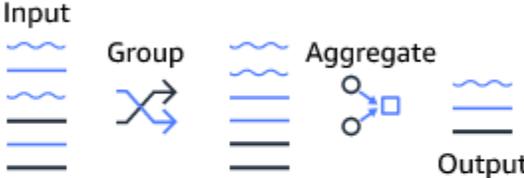
フィールドを選択すると、名前とデータ型が表示されます。フィールドを削除するには、フィールドで [X] を選択します。

Node properties | **Transform 1** | Output schema | 

Data preview

▼ **Aggregate** [Info](#)

This transform first groups your rows by fields you choose, and then computes the aggregated value for fields you choose by specific function (e.g., sum, average, max).



Fields to group by - *optional*

Select the fields you would like to group your rows by, so the aggregation would be done for each unique group.

Choose one or more fields ▼

Aggregate another column

 Add an aggregation.

- [Aggregate another column]を選択します。少なくとも1つのフィールドを選択する必要があります。

Field to aggregate

Choose a field ▼

Aggregation function [Info](#)

Choose a function ▼



Aggregate another column

- [Field to aggregate]ドロップダウン中のフィールドを選択します。

5. 選択したフィールドに適用する集計関数を選択します：

- avg-平均を計算します。
- CountDistinct-一意の非ヌル値の数を計算します。
- count-非 null 値の数を計算します。
- first-「グループ別」の基準を満たす最初の値を返します。
- last-「グループ別」の基準を満たす最後の値を返します。
- kurtosis-周波数デイス Tribi ューション曲線のピークのシャープネスを計算します。
- max-「グループ別」の基準を満たす最高値を返します。
- min-「グループ別」の基準を満たす最小値を返します。
- skewness-正規デイス Tribi ューションの確率デイス Tribi ューションの非対称性の尺度
- stddev_pop-母標準偏差を計算し、母分散の平方根を返します。
- sum-グループ内のすべての値の合計
- sumDistinct-グループ内の個別の値の合計
- var_samp-グループのサンプル分散 (ヌルは無視されます)
- var_pop-グループの母分散 (ヌルは無視されます)

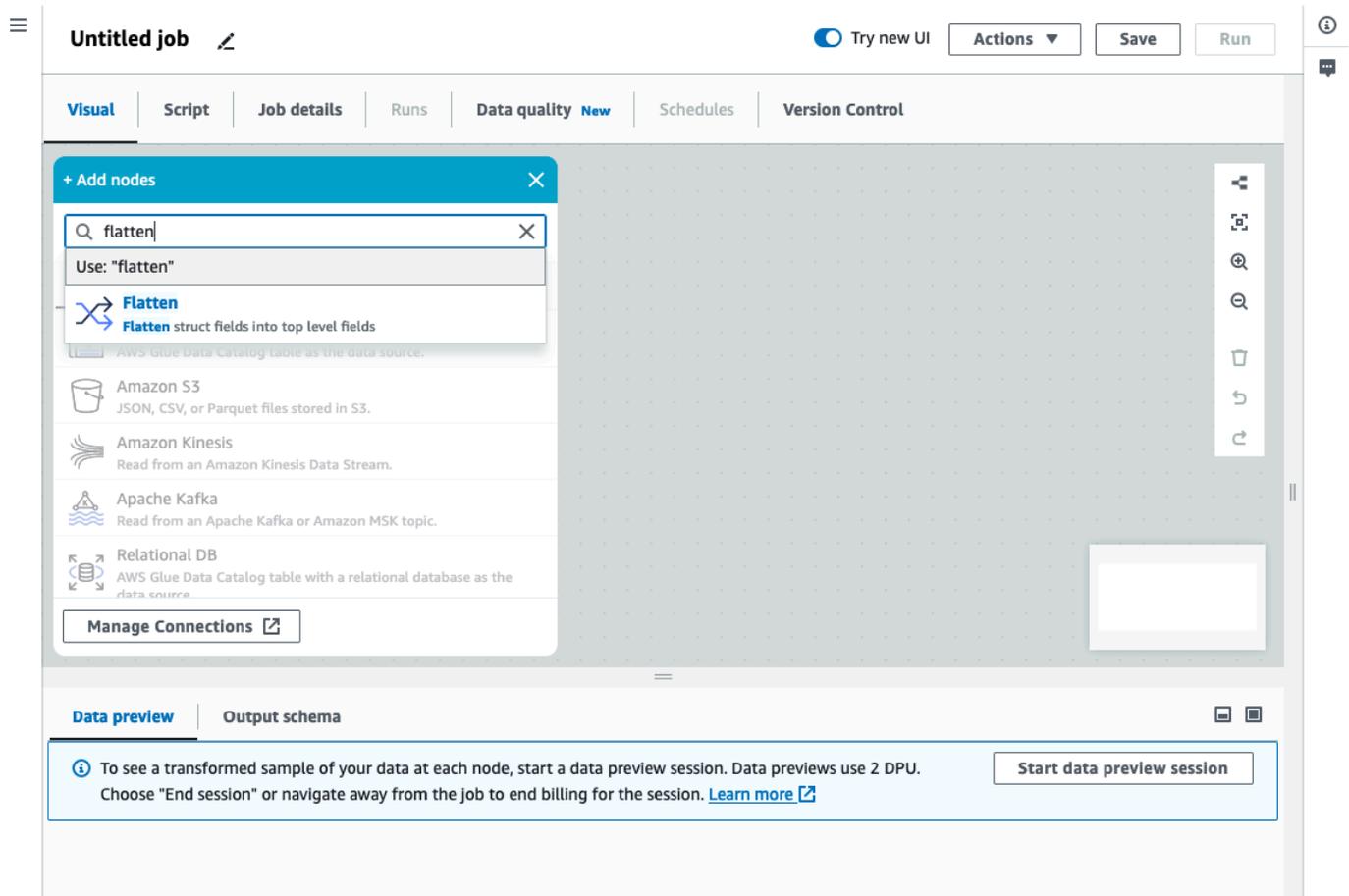
ネストされた Struct のフラット化

データ内のネストされた Struct のフィールドをフラット化し、最上位のフィールドにします。新しいフィールドの名前には、そこに届く Struct フィールドの名前が接頭辞として付いたフィールド名が、ドットで区切られて使用されます。

例えば、データに「phone_numbers」という Struct 型のフィールドがあり、他のフィールドでの「home_phone」という名前の Struct 型の 1 つとともに、「country_code」および「number」という 2 つのフィールドを伴っているとします。フラット化されると、この 2 つのフィールドはそれぞれ「phone_numbers.home_phone.country_code」および「phone_numbers.home_phone.number」という名前の最上位のフィールドになります。

フラット化された変換ノードをジョブ図に追加するには

1. リソースパネルを開いて、[変換] タブの [フラット化] 選択し、ジョブ図に新しい変換を追加します。検索バーを使用して、「Flatten」と入力し、フラット化ノードをクリックすることもできます。ノードを追加する際に選択したノードが、その親になります。



2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. (オプション) [Transform] (変換) タブでは、フラット化するネストレベルの上限を設定できます。例えばその値を 1 に設定すると、最上位の Struct のみをフラット化できます。上限を 2 に設定すると、最上位とその次の Struct をフラット化できます。

UUID カラムを追加する

UUID (ユニバーサル一意識別子) 列を追加すると、各行に 36 字の一意的文字列が割り当てられます。

UUID 変換ノードをジョブ図に追加するには

1. リソースパネルを開いて、[UUID] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。

2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. (オプション) 新しい列の名前は [Transform] (変換) タブでカスタマイズできます。デフォルトでは「uuid」という名前が付けられます。

識別子列を追加する

データセット内の各行に数値識別子を割り当てます。

ジョブ図に識別子変換ノードを追加するには

1. リソースパネルを開いて、[Identifier] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. (オプション) 新しい列の名前は [Transform] (変換) タブでカスタマイズできます。デフォルトでは、「id」という名前が付けられます。
4. (オプション) ジョブがデータを段階的に処理して保存する場合は、ジョブの実行時に同じ ID が再利用されないよう注意する必要があります。

[Transform] (変換) タブで、チェックボックスのオプション [unique] (一意) をチェックします。これで識別子にジョブのタイムスタンプが追加されるため、複数の実行間で一意になります。数字を増やせるようにするには、long の代わりに列を 10 進数にします。

列をタイムスタンプタイプに変換する

To timestamp 変換を使用すると、数値または文字列の列のデータ型をタイムスタンプに変更できます。これで、そのデータ型で保存したり、タイムスタンプを必要とする他の変換に適用したりできるようになります。

To timestamp 変換ノードをジョブ図に追加するには

1. リソースパネルを開いて、[To timestamp] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。

2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. [Transform] (変換) タブで、変換する列の名前を入力します。
4. [Transform] (変換) タブで、タイプを選択し、選択した列の解析方法を定義します。

値が数値の場合は、秒 (Unix/Python タイムスタンプ)、ミリ秒、マイクロ秒のいずれかで表すことができます。その際は対応するオプションを選択します。

値がフォーマットされた文字列の場合は、「iso」タイプを選択します。文字列は ISO 形式のいずれかに準拠している必要があります (例: 「2022-11-02T14:40:59.915Z」)。

この時点でタイプが不明な場合や、行によってタイプが異なる場合は、[autodetect] (自動検出) を選択すると、システムが最小限の実行コストでベストな推測を行います。

5. (オプション) [Transform] タブでは、選択した列を変換する代わりに、新しい列の名前を入力すれば、新しい列を作成し、元の列を残すことができます。

タイムスタンプ列をフォーマットされた文字列に変換する

タイムスタンプ列をパターンに基づいた文字列にフォーマットします。Format timestamp を使用すると、日付と時刻を希望する形式の文字列で取得できます。形式の定義には、[Spark の日付構文](#)とほとんどの [Python 日付コード](#)を使用できます。

例えば、日付文字列を「2023-01-01 00:00」のような形式にしたいときは、Spark 構文を「yyyy-MM-dd HH:mm」として使用するか、Python の同様の日付コード「%Y-%m-%d %H:%M」を使用することで、形式を定義できます。

Format timestamp 変換ノードをジョブ図に追加するには

1. リソースパネルを開いて、[Format timestamp] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. [Transform] (変換) タブで、変換する列の名前を入力します。
4. [Transform] タブに、使用するタイムスタンプ形式のパターンを、[Spark 日付構文](#)または [Python 日付コード](#)を使って入力します。

5. (オプション) [Transform] タブでは、選択した列を変換する代わりに、新しい列の名前を入力すれば、新しい列を作成し、元の列を残すことができます。

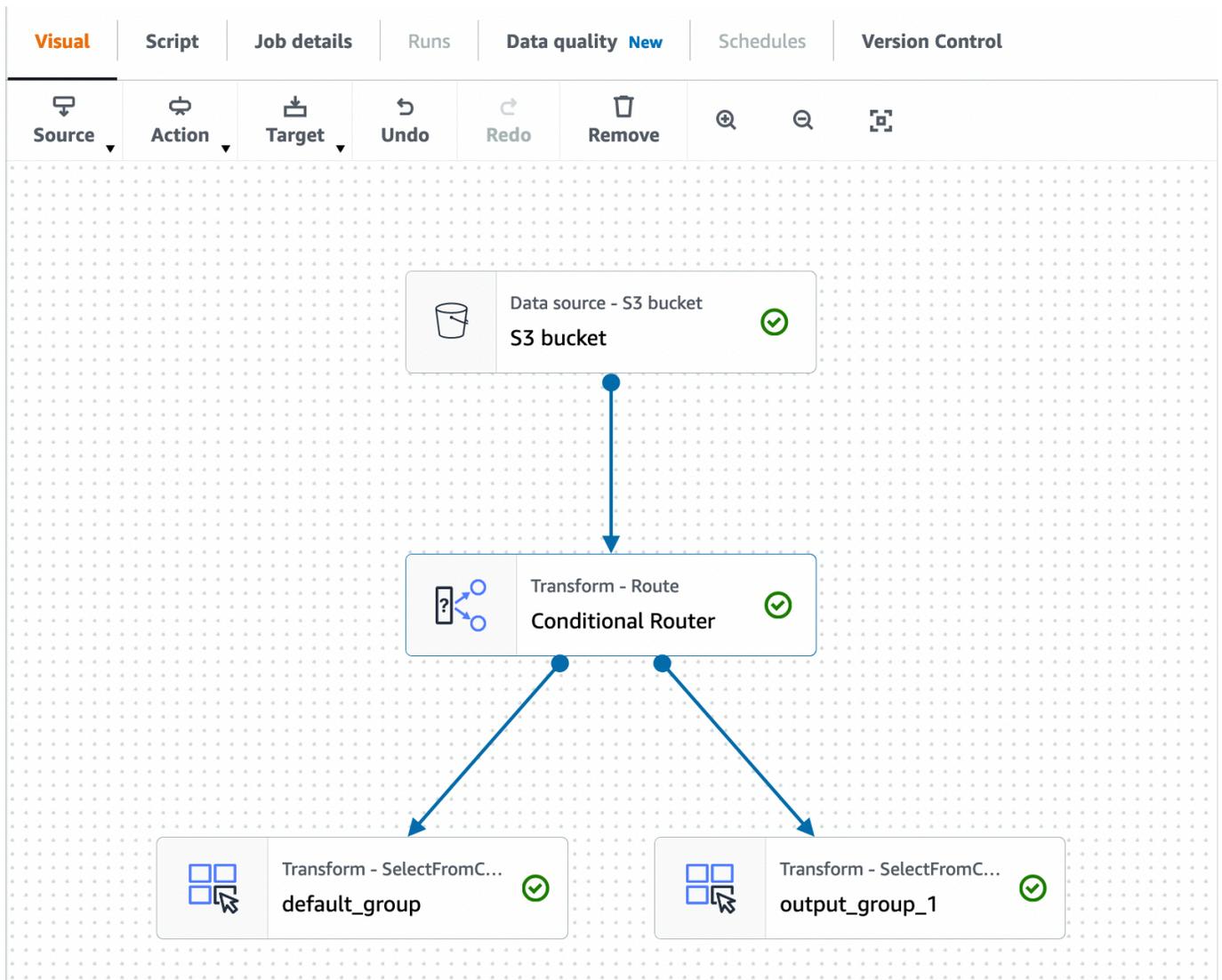
条件付きルーター変換の作成

条件付きルーター変換を使うと、着信データに複数の条件を適用することができます。着信データの行はそれぞれグループフィルター条件により評価され、対応するグループに取り入れて処理されます。行が複数のグループフィルター条件を満たしている場合、変換により、その行が複数のグループに渡されます。行がどの条件も満たしていない場合は、その行を削除するか、デフォルトの出力グループにルーティングします。

この変換はフィルター変換に似ていますが、同じ入力データを複数の条件でテストしたいユーザーにとって便利です。

条件付きルーター変換を追加するには

1. 条件付きルーター変換を実行するノードを選択します。これは、ソースノードでも別の変換でもかまいません。
2. [Action] (アクション) を選択し、検索バーで [Conditional Router] (条件付きルーター) を検索し、選択します。条件付きルーター変換が、2つの出力ノードと共に追加されます。一方の出力ノード [Default group] (デフォルトグループ) には、もう一方の出力ノードで定義された条件のいずれをも満たさないレコードが含まれます。デフォルトのグループは編集できません。



[Add group] (グループを追加) を選択すると、出力グループを追加できます。各出力グループで、グループに名前を付けてフィルタ条件と論理演算子を追加できます。

Node properties | **Transform** | Output schema | Data preview ✕

Add group

output_group_1

Remove group

Define a set of conditions a record has to meet in order to be routed to the output group.

Group name
The name of this output group, as it would appear in your job. Letters, numbers, _ and - are allowed.

Logical operator

AND
Trigger only when ALL conditions are met.

OR
Trigger when at least one of the conditions is met.

Filter condition [Info](#)
Specify your filter condition by choosing the key, operator, and entering a value.

Start by adding a filter condition.

Add condition

Default group

Records which do not meet any of the conditions defined above will be routed here.

- 出力グループの名前を変更するときは、グループに新しい名前を入力します。AWS Glue Studio は、ユーザーに代わって自動的にグループに名前を付けます (例: 「output_group_1」)。
- 論理演算子 (AND、OR) を選択し、[Key] (キー)、[Operation] (演算)、[Value] (値) を指定して [Filter condition] (フィルター条件) を追加します。論理演算子を使用すると、複数のフィルター条件を実装し、指定した各フィルター条件に対して論理演算子を実行できます。

キーを指定するときは、スキーマに表示されているキーの中から選択します。その後、選択したキーの種類に応じて使用可能な演算を選択します。例えば、キーの種類が「string」の場合、選択できる演算は「match」です。

Filter condition **Info**

Specify your filter condition by choosing the key, operator, and entering a value.

Key	Operation	Value	
year ▼	= ▼	2023	
Add condition			

- [Value] (値) フィールドに値を入力します。その他の条件を追加するときは、[Add condition] (条件を追加) を選択します。フィルター条件を削除するときは、ごみ箱のアイコンを選択します。

[列の連結] 変換を使用して列を追加する

[連結] 変換では、オプションのスペーサーを使用して、他の列の値を使った新しい文字列の列を作成できます。「year」、「month」、「day」を順にスペーサー「-」で連結して、連結した列「date」を定義する例を以下に示します。

day	month	year	date
01	01	2020	2020 年 1 月 1 日
02	01	2020	2020 年 1 月 2 日
03	01	2020	2020 年 1 月 3 日
04	01	2020	2020 年 1 月 4 日

[連結] 変換を追加するには

- リソースパネルを開きます。次に [列の連結] を選択して、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
- (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
- [変換] タブで、連結する列および連結した文字列を保持する列の名前を入力します。ドロップダウンで列にチェックを入れた順序が、使用される順序になります。

Node properties	Transform	Output schema	Data preview	
-----------------	------------------	---------------	--------------	---

Name of the concatenated column
Name of the string column that will be generated

List of column named separated by comma or spaces
The fields listed will be concatenated on that order

Array new column Name - optional
String to place between the concatenated fields, by default there is no spacer.

Null value - optional
The string to use when a column value is null, for example: 'NULL' or 'NA', by default an empty string will be used

4. [スペーサー - オプション] — 連結したフィールドの間に配置する文字列を入力します。デフォルトでは、スペーサーはありません。
5. [NULL 値 - オプション] — 列の値が「NULL」の場合に使用する文字列を入力します。デフォルトでは、列の値が「NULL」または「NA」の場合、空の文字列が使用されます。

[文字列の分割] 変換を使用して文字列の列を分割する

[文字列の分割] 変換を使用して、正規表現の文字列をトークンの配列に分割し、分割方法を定義します。その後、列を配列型のままにするか、この後に [配列から列へ] 変換を適用して、配列の値を抽出し最上位のフィールドに追加できます。ただし、各トークンの意味が事前にわかっていることが前提です。また、カテゴリのセットなどトークンの順序が関係ない場合、[分解] 変換を使用して値ごとに個別の行を生成できます。

例えば、カンマを分割のパターンとして使用して「categories」列を分割し、「categories_arr」列を追加できます。

product_id	カテゴリ	categories_arr
1	sports,winter	[sports, winter]
2	garden,tools	[garden, tools]
3	videogames	[videogames]
4	game,boardgame,social	[game, boardgame, social]

[文字列の分割] 変換を追加するには：

1. リソースパネルを開いて、[Split String] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. [変換] タブで、分割する列を選択して文字列の分割に使用するパターンを入力します。ほとんどの場合、正規表現として特別な意味がありエスケープする必要がある場合を除いて、文字を入力するだけで済みます。エスケープする必要がある文字は `\.[]{}()<>*+ -= !?^$|` で、文字の前にバックスラッシュを追加することでエスケープできます。例えば、ドット (「.」) で区切る場合は、「\。」と入力する必要があります。ただし、カンマには特別な意味はなく、「,」のようにそのまま指定できます。

Node properties

Transform

Output schema

Data preview

✕

Column to split
Column whose string will be split into an string array

Splitting regular expression
Regex defining the separator token, examples: ';', '\|' (pipe needs to be escaped) or '\s+' (whitespace split)

Array column Name - optional
Name to use for the column with the extracted array resulting of the split. If not specified, instead of a new column the existing one is replaced

4. (オプション) 元の文字列の列を保持したい場合は、新しい配列の列の名前を入力できます。これにより、元の文字列の列と新しいトークン化された配列の列の両方が保持できます。

[配列から列へ] 変換を使用して、配列の要素を抽出して最上位の列に追加する

[配列から列へ] 変換では、配列型の列の一部またはすべての要素を抽出して新しい列に追加できます。配列に抽出に十分な数の値がある場合、変換によって新しい列が可能な限り埋め込まれ、オプションで指定された位置の要素を取得できます。

例えば、ip v4 サブネットに「文字列の分割」変換を適用した結果である配列の列「subnet」がある場合、最初と 4 番目の要素を抽出して、新しい列「first_octect」と「forth_octect」に追加できます。この例の変換の出力は次のようになります (最後の 2 行の配列が期待する長さよりも短いことに注意してください)。

サブネット	first_octect	fourth_octect
[54, 240, 197, 238]	54	238
[192, 168, 0, 1]	192	1
[192, 168]	192	

サブネット	first_octect	fourth_octect
[]		

[配列から列へ] 変換を追加するには:

1. リソースパネルを開いて、[Array To Columns] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. [変換] タブで、抽出する配列の列を選択し、抽出されたトークンの新しい列のリストを入力します。

Node properties
Transform
Output schema
Data preview
⌵

Array type column
Column of type array from which the new columns are extracted

Output columns
The names (separated by commas) of the columns to create out of the array fields. The data type will be the same as the array. For each row, the transform will try to fill them as much as possible using the array elements, the rest will be NULL

Array indexes to use - optional
List of array positions (starting from 1 and separated by commas), indicating which columns to take to fill the columns. Only need to set this if you want to skip some positions of the array

4. (オプション) 列を割り当てるために配列のトークンを取得しない場合は、指定した同じ順序で列のリストに割り当てるインデックスを指定できます。例えば、出力列が「column1, column2, column3」でインデックスが「4, 1, 3」の場合、配列の 4 番目の要素は column1 に、最初の要素は column2 に、3 番目の要素は column3 に指定できます (配列がインデックスの番号より短い場合は NULL 値が設定されます)。

[現在のタイムスタンプを追加] 変換の使用

[現在のタイムスタンプを追加] 変換では、データが処理された時刻で行をマークできます。これは、監査目的やデータパイプラインでのレイテンシーの追跡に役立ちます。この新しい列は、タイムスタンプのデータ型またはフォーマットされた文字列として追加できます。

[現在のタイムスタンプを追加] 変換を追加するには:

1. リソースパネルを開いて、[Add Current Timestamp] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。

Node properties	Transform	Output schema	Data preview	⌵
<p>Timestamp column - optional Name to use for the new column, by default: timestamp. With type "string" if a dataFormat is specified, otherwise "timestamp"</p> <input type="text"/>				
<p>Timestamp format - optional Optional pattern to format as a string, accepts most Python date format codes, such as '%Y-%m-%d %H:%M:%S'; as well as Spark patterns such as 'yyyy-MM-dd'T'HH:mm:ss.SSSZ'</p> <input type="text"/>				

3. (オプション) 列をフォーマット済の日付の文字列にしたい場合は、[変換] タブで、新しい列のカスタム名とフォーマットを入力します。

[行から列へのピボット] 変換の使用

[行から列へのピボット] 変換では、選択した列の固有の値を回転させて数値列を集約し、新しい列にすることができます (複数の列を選択した場合は、値が連結され新しい列に名前が付けられます)。このように、それぞれの固有の値ごとに部分的な集計を含む列を増やししながら、行を統合します。例えば、月別および国別の売上に関するデータセットがあるとします (わかりやすいように並べ替えられています)。

年	month	country	amount
2020	Jan	uk	32
2020	Jan	de	42
2020	Jan	us	64
2020	Feb	uk	67
2020	Feb	de	4
2020	Feb	de	7
2020	Feb	us	6
2020	Feb	us	12
2020	Jan	us	90

金額と国を集計列としてピボットすると、元の国列から新しい列が作成されます。次の表では、国列の代わりに de、uk、us の新しい列が作成されています。

年	month	de	uk	us
2020	Jan	42	32	64
2020	Jan	11	67	18
2021	Jan			90

代わりに、月と国の両方をピボットする場合は、それらの列の値の組み合わせごとに列が表示されず。

year	Jan_de	Jan_uk	Jan_us	Feb_de	Feb_uk	Feb_us
2020	42	32	64	11	67	18

year	Jan_de	Jan_uk	Jan_us	Feb_de	Feb_uk	Feb_us
2021			90			

[行から列へのピボット] 変換を追加するには:

1. リソースパネルを開いて、[Pivot Rows To Columns] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. [変換] タブで、新しい列の値を生成するために集約される数値列、適用する集計関数、および固有の値を新しい列に変換する列を選択します。

Node properties
Transform
Output schema
Data preview
⌵

Aggregation column

Numeric column on which the aggregation function is applied

Aggregation

The Spark function to apply to the aggregation column.

Columns to convert

List of columns whose values will become new columns. If multiple columns are specified, the values are concatenated using underscore.

Choose options
⌵

[列から行へのピボット解除] 変換の使用

[ピボット解除] 変換では、列を新しい列の値に変換して、固有の値ごとに行を生成できます。この変換はピボットの逆ですが、集約された同じ値の行を分解したり、組み合わせを元の列に分割できないため、同等ではないことに注意してください (これらは、後で [分割] 変換を使用して行うことができます)。例えば、次のようなテーブルがあるとします。

年	month	de	uk	us
2020	Jan	42	32	64
2020	Feb	11	67	18
2021	Jan			90

値「amount」を使用して「de」、「uk」、「us」の列を「country」列にピボット解除すると、次のようになります (ここでは分かりやすくするために並べ替えられています)。

年	month	country	amount
2020	Jan	uk	32
2020	Jan	de	42
2020	Jan	us	64
2020	Feb	uk	67
2020	Feb	de	11
2020	Feb	us	18
2021	Jan	us	90

NULL 値の列 (Jan 2021 の「de」と「uk」) はデフォルトでは生成されません。このオプションを有効にすると、次のようになります。

年	month	country	amount
2020	Jan	uk	32
2020	Jan	de	42
2020	Jan	us	64

年	month	country	amount
2020	Feb	uk	67
2020	Feb	de	11
2020	Feb	us	18
2021	Jan	us	90
2021	Jan	de	
2021	Jan	uk	

[列から行へのピボット解除] 変換を追加するには:

1. リソースパネルを開いて、[Unpivot Columns to Rows] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. [変換] タブで、ピボット解除する列の名前と値を保持するために作成する新しい列を入力します。

Node properties	Transform	Output schema	Data preview	⌵
-----------------	------------------	---------------	--------------	---

Unpivot names column
Column to create out of the source columns names

Unpivot values column
Column to create out of values of the old columns

Columns to unpivot into the new value column
List of columns whose name will become values of the new column

[オートバランス処理] 変換を使用してランタイムを最適化する

[オートバランス処理] 変換は、パフォーマンスを向上させるためにデータをワーカー間で再配分します。これは、データのバランスが取れていない場合や、ソースからのデータでは十分に並行処理ができない場合に役立ちます。これは、ソースが gzip 圧縮されている場合や JDBC である場合に一般的です。データの再配分には中程度のパフォーマンスコストがかかるため、データのバランスが既に取れている場合は、最適化により必ずしもコストが補われるとは限りません。変換で Apache Spark の再パーティション化を使用して、クラスターのキャパシティに最適な複数のパーティション間でデータをランダムに再度割り当てできます。上級ユーザーの場合は、複数のパーティションを手動で入力できます。さらに、指定した列に基づいてデータを再編成することで、分割テーブルの書き込みを最適化することもできます。これにより、出力ファイルがより統合されます。

1. リソースパネルを開いて、[Autobalance Processing] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. (オプション) [変換] タブで、複数のパーティションを入力できます。一般的に、この値はシステムに決定させることをお勧めしますが、値を制御する必要がある場合は乗数を調整したり特定の値を入力できます。列で分割されたデータを保存する場合は、再パーティションの列と同じ列

を選択できます。これにより、各パーティションのファイルの数を最小限に抑え、パーティションごとに多数のファイルが作成されるのを防ぐことができます。ファイルが多数作成されると、データをクエリするツールのパフォーマンスが低下します。

Node properties
Transform
Output schema
Data preview

Number of partitions - optional
 Number of partitions on which to randomly distribute the data. If the number ends with the x letter then it means it's a multiple of the number of cores in the cluster. By default: 2x

Repartition columns - optional
 Instead of randomly reassign the data to partitions, assign data with the same values of the columns specified to the same partition.

Choose options
▼

[派生列] 変換を使用して他の列を結合する

[派生列] 変換では、定数やリテラル、データ内の他の列も使用できる数式または SQL 式に基づいて新しい列を定義できます。例えば「success」および「count」列から「percentage」列を派生させるには、「`success * 100 / count || '%'`」という SQL 式を入力します。

結果の例:

success	count	割合(%)
14	100	14%
6	20	3%
3	40	7.5%

[派生列] 変換を追加するには

- リソースパネルを開いて、[Derived Column] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。

- (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
- [変換] タブで、列の名前とその内容の式を入力します。

Node properties	Transform	Output schema	Data preview	✕
-----------------	------------------	---------------	--------------	---

Name of the derived column
Name to use for the new column or replace an existing one

SQL Expression
A SQL expression that defines the column, which can be derived from other existing columns and use operators to modify or combine them. For instance, to derive a percentage from the columns "success" and "count", you can enter: "success * 100 / count"

[ルックアップ] 変換を使用してカタログテーブルから一致するデータを追加する

[ルックアップ] 変換では、キーがデータ内の定義済みのルックアップ列と一致する場合、定義済みのカタログテーブルから列を追加できます。これは、データとルックアップテーブル間を左外部結合し、条件に一致する列として使用するのと同じです。

[ルックアップ] 変換を追加するには

- リソースパネルを開いて、[Lookup] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
- (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
- [変換] タブで、ルックアップの実行に使用するカタログテーブルの完全修飾名を入力します。例えば、データベースが「mydb」で、テーブルが「mytable」の場合は、「mydb.mytable」と入力します。次に、ルックアップキーが作成されている場合は、ルックアップテーブルで一致するものを見つけるための基準を入力します。キー列のリストは、カンマで区切って入力します。1つ以上のキー列に同じ名前がない場合は、マッチマッピングを定義する必要があります。

例えば、データ列が「user_id」と「region」で、ユーザーテーブルで対応する列の名前が「id」と「region」の場合、[一致する列] フィールドに「user_id=id, region」と入力します。region=region と入力することもできますが、これらは同じなので必要ありません。

- 最後に、ルックアップテーブルで一致した行から取得する列を入力して、それらをデータに組み込みます。一致するものが見つからなかった場合、それらの列は NULL に設定されます。

Note

[ルックアップ] 変換では、効率を上げるため左結合を使用しています。ルックアップテーブルに複合キーがある場合は、一致する列がすべてのキーの列と一致するように設定されていることを確認し、一致が1つだけになるようにしてください。そうしないと、複数のルックアップ行が一致し、一致するごとに余分な行が追加されます。

Node properties

Transform

Output schema

Data preview



AWS Glue Data Catalog table

Qualified name of the catalog table to use for the lookup, specifying the database and table name separated by a dot

Lookup key columns to match

Columns in the lookup table to match separated by commas; if the column names don't match, you can specify the mapping between the data and the lookup table separating the names with an equals sign =

Lookup columns to take

Columns in the lookup table to add to the data when a match is found in the lookup table

[配列またはマップを行に分解] 変換の使用

[分解] 変換では、ネストされたデータ構造から値を抽出し、操作しやすい個々の行を追加できます。配列の場合、変換により配列の各値に対して行が生成され、その行に他の列の値が複製されます。マップの場合、変換によりキーと値を列として持つ各エントリに対して行が生成されます。各エントリには列としてのキーと値があり、その行に他の列も複製されます。

例えば、このデータセットに複数の値がある「category」配列の列があるとします。

product_id	category
1	[sports, winter]
2	[garden, tools]
3	[videogames]
4	[game, boardgame, social]
5	[]

「category」列が同じ名前の列に分解されると、その列が上書きされます。NULL を含めるように選択すると、次のようになります (分かりやすいように順番になっています)。

product_id	category
1	sports
1	winter
2	garden
2	tool
3	videogames
4	ゲーム
4	boardgame
4	social
5	

[配列またはマップを行に分解] 変換を追加するには:

1. リソースパネルを開いて、[Explode Array Or Map Into Rows] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. (オプション) [Node properties] (ノードのプロパティ) タブで、ジョブ図にノードの名前を入力できます。ノードの親がまだ選択されていない場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. [変換] タブで、分解する列を選択します (タイプが配列またはマップである必要があります)。次に、配列の項目を表す列の名前を入力するか、マップを分解する場合はキーと値を表す列の名前を入力します。
4. (オプション) [変換] タブでは、デフォルトで分解する列が NULL、またはデータ構造が空である場合、分解されたデータセットではその列が省略されます。(新しい列を NULL として) 行を残しておきたい場合は、[NULL を含む] にチェックを入れます。

Node properties	Transform	Output schema	Data preview	
-----------------	------------------	---------------	--------------	---

Column to explode
A column of type array or map

New column name
The name of the column to put the array values or the dictionary keys

Values column - optional
If exploding a dictionary, you can specify a name for a column to contain the values. Default name: "value"

Include NULLs - optional
If selected, NULL values will also generate a new rows, otherwise the row with a NULL value is omitted

レコードマッチング変換を使用して既存のデータ分類変換を呼び出す

この変換は、既存のレコードマッチング機械学習データ分類変換を呼び出します。

この変換は、ラベルに基づいて、トレーニング済みモデルに対して現在のデータを評価します。列「match_id」が追加され、アルゴリズムトレーニングに基づいて、同等とみなされる項目のグループに各行が割り当てられます。詳細については、「[AWS Lake Formation FindMatches によるレコードのマッチング](#)」を参照してください。

Note

ビジュアルジョブで使用する AWS Glue のバージョンは、レコードマッチング変換の作成時に AWS Glue で使用するバージョンと一致する必要があります。

Transform	Output schema	Data preview				
Data preview (20) Info		Previewing 6 of 7 fields				
<input type="text" value="Filter sample dataset"/>						
id	title	venue	year	source	match_id	
journals_sigmod_Liu02	Editor's Notes	SIGMOD Record	2002	DBLP	25769803776	
journals_sigmod_Hammer02	Report on the ACM Fourth International Workshop on Data Warehousing and OLAP (DOLAP 2001)	null	2002	DBLP	25769803777	
journals_sigmod_Konig-RiesMMPPRSVW02	Report on the NSF Workshop on Building an Infrastructure for Mobile and Wireless Systems	null	2002	DBLP	68719476736	

レコードマッチング変換ノードをジョブ図に追加するには

- リソースパネルを開いて、[Record Matching] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
- ノードのプロパティパネルで、ジョブ図にノード名を入力します。ノードの親がまだ選択されていない場合、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
- [変換] タブに、[機械学習の変換] ページから取得した ID を入力します。

AWS Glue > ML transforms

Machine learning transforms (1) [Info](#)
Clean all your data using machine learning transforms.

Q Filter transforms

	Transform name ▲	ID	Status ▼	Label count ▼
<input type="radio"/>	Test	tfm-3d291b652cec092a79aeda5062f2c96e7c528474	✔ Ready for use	352

- (オプション) [変換] タブで、信頼スコアを追加するオプションを確認できます。計算量は増えますが、モデルは各マッチングの信頼スコアを追加の列として推定します。

null 行の削除

この変換は、すべての列が null である行をデータセットから削除します。さらに、この条件を拡張して空のフィールドを含め、1 つ以上の列が空でない行を保持することもできます。

null 行削除変換ノードをジョブ図に追加するには

- リソースパネルを開いて、[Remove Null Rows] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
- ノードのプロパティパネルで、ジョブ図にノード名を入力します。ノードの親がまだ選択されていない場合、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
- (オプション) 行が null でもなく空でもないようにする必要がある場合は、[変換] タブで [Extended] オプションをオンにします。この変換により、空の文字列、配列、マップは null とみなされます。

JSON データを含む文字列の列の解析

この変換は、JSON データを含む文字列の列を解析し、JSON がオブジェクトか配列かに応じて、それぞれを構造体または配列の列に変換します。(オプション) 解析された列と元の列の両方を保持できます。

JSON スキーマは、オプションのサンプリングを使用して提供または推測できます (JSON オブジェクトの場合)。

JSON 列解析変換ノードをジョブ図に追加するには

1. リソースパネルを開いて、[Parse JSON Column] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. ノードのプロパティパネルで、ジョブ図にノード名を入力します。ノードの親がまだ選択されていない場合、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. [変換] タブで、JSON 文字列を含む列を選択します。
4. (オプション) [変換] タブで、JSON データが従うスキーマを、SQL 構文を使用して入力します (例: オブジェクトの場合は「field1 STRING, field2 INT」、配列の場合は「ARRAY<STRING>」)。

配列の場合はスキーマが必要ですが、オブジェクトの場合はデータを使用して推測されます (スキーマが指定されていない場合)。スキーマの推論による影響を軽減するため (特に大規模なデータセットの場合)、[Ratio of samples to use to infer schema] を入力することで、データ全体を 2 回読み取る必要がなくなります。値が 1 より小さい場合は、対応するランダムサンプルの比率を使用してスキーマが推測されます。データに信頼性があり、行間でオブジェクトが一貫している場合は、0.1 などの小さな比率を使用してパフォーマンスを高めることができます。

5. (オプション) 元の文字列の列と解析された列の両方を保持する場合は、[変換] タブで新しい列名を入力します。

JSON パスの抽出

この変換は、JSON 文字列の列から新しい列を抽出します。この変換は、必要なデータ要素が少なく、JSON コンテンツ全体をテーブルスキーマにインポートしない場合に便利です。

JSON パス抽出変換ノードをジョブ図に追加するには

1. リソースパネルを開いて、[Extract JSON Path] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. ノードのプロパティパネルで、ジョブ図にノード名を入力します。ノードの親がまだ選択されていない場合、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。
3. [変換] タブで、JSON 文字列を含む列を選択します。1 つ以上の JSON パス式をカンマで区切って入力します。各式は、JSON 配列またはオブジェクトから値を抽出する方法を示しています。例えば、JSON 列に「prop_1」と「prop2」というプロパティを持つオブジェクトが含まれている場合、「prop_1, prop_2」と名前を指定して両方を抽出できます。

JSON フィールドに特殊文字が含まれていて、この JSON {"a. a": 1} からプロパティを抽出する場合などには、パス \$['a. a'] を使用できます。例外は、パスを区切るためのカンマです。続いて、各パスに対応する列名をカンマで区切って入力します。

4. (オプション) [変換] タブでは、抽出された JSON 列を削除するようにチェックできます。これは、必要な JSON データを抽出した後、残りのデータが不要な場合に便利です。

正規表現による文字列の断片の抽出

この変換は、正規表現を使用して文字列の断片を抽出し、そこから新しい列を作成します。正規表現グループを使用する場合は複数の列を作成します。

正規表現抽出変換ノードをジョブ図に追加するには

1. リソースパネルを開いて、[Regex Extractor] を選択し、ジョブ図に新しい変換を追加します。ノードを追加する際に選択したノードが、その親になります。
2. ノードのプロパティパネルで、ジョブ図にノード名を入力します。ノードの親がまだ選択されていない場合、[Node parents] (ノードの親) リストから、変換の入力ソースとして使用するノードを選択します。
3. [変換] タブに、正規表現とそれを適用する必要がある列を入力します。続いて、一致する文字列を格納する新しい列の名前を入力します。新しい列は、ソース列が null の場合にのみ null になり、正規表現が一致しない場合は空になります。

正規表現がグループを使用する場合、対応する列名はカンマで区切られますが、列名を空のままにすることでグループをスキップできます。

例えば、長い ISO 日付形式と短い ISO 日付形式の両方を使用する文字列を含む列

「purchase_date」があり、使用可能な場合は、年、月、日、時間を抽出します。時間グループはオプションであり、使用できない場合、(正規表現が一致しないため) 抽出されたすべてのグループは空の文字列になります。グループで、時間をオプションではなく含めたい場合は、名前を空のままにすることで抽出されなくなります (グループには T 文字が含まれます)。

Transform | Output schema | Data preview 

Name

Node parents

Choose which nodes will provide inputs for this one.

S3 bucket 

S3 - DataSource

Column to extract from

String column on which to apply the regex.

string

Regular expression

Regex to apply on the column, if multiple columns need to be extracted then the expression needs an equal number of groups.

Extracted column

The name of the column where to extract the matched regex. Multiple column names can be specified separated by commas, if the name is empty it means that group is skipped. If the source column is null, the new column will be null as well, otherwise an empty string means there was no match.

データプレビューの結果:

Data preview (5) [Info](#) Previewing 5 of 5 fields



<code>purchase_date</code>	<code>year</code>	<code>month</code>	<code>day</code>	<code>hour</code>
2023-03-04T12:23:31	2023	03	04	12
2021-06-09T02:21:01	2021	06	09	02
2022-02-04	2022	02	04	
2020-09-05T23:07:02	2020	09	05	23
2020-09-08	2020	09	08	

カスタム変換を作成する

データに対してより複雑な変換を実行する必要がある場合、またはデータプロパティキーをデータセットに追加する場合は、ジョブ図にカスタムコード変換を追加します。カスタムコードノードを使用すると、変換を実行するスクリプトを入力できます。

カスタムコードを使用する場合は、スキーマエディタを使用して、カスタムコードを通して出力に加えられる変更を指定する必要があります。スキーマの編集時に、次のアクションを実行できます。

- データプロパティキーの追加または削除
- データプロパティキーのデータ型の変更
- データプロパティキーの名前の変更
- ネストされたプロパティキーの再構築

出力をターゲットの場所へ送信するには、`SelectFromCollection` 変換を使用して、カスタム変換ノードの結果から単一の `DynamicFrame` を選択する必要があります。

次のタスクを実行して、カスタム変換ノードをジョブ図に追加できます。

カスタムコード変換ノードをジョブ図に追加する

カスタム変換ノードをジョブ図に追加するには

1. (オプション) リソースパネルを開いて、[Custom transform] を選択し、ジョブ図にカスタム変換を追加します。
2. [Node properties] (ノードのプロパティ) タブで、ジョブ図のノードの名前を入力します。ノードの親がまだ選択されていない場合、またはカスタム変換に複数の入力が必要な場合は、[Node parents] (ノードの親) リストから、変換の入カソースとして使用するノードを選択します。

カスタム変換ノードのコードを入力する

入力フィールドにコードを入力またはコピーできます。ジョブでは、データ変換の実行にこのコードを使用します。Python または Scala のいずれかでコードのスニペットを指定できます。コードは、入力として 1 つ以上の DynamicFrames を取り、DynamicFrames のコレクションを返す必要があります。

カスタム変換ノードのスクリプトを入力するには

1. ジョブ図でカスタム変換ノードを選択した状態で、[Transform] (変換) タブを選択します。
2. 見出し [Code block] (コードブロック) の下のテキスト入力フィールドで、変換用のコードを貼り付けるか、入力します。使用するコードは、[Job details] (ジョブの詳細) タブのジョブで指定されている言語と一致する必要があります。

コード内の入力ノードを参照する場合、AWS Glue Studio は、作成された順序に基づきジョブの図表ノードによって連続的に返される DynamicFrames に名前を付けます。コード内で次のいずれかのネーミングメソッドを使用します。

- Classic code generation — 関数名を使用して、ジョブダイアグラム内のノードを参照します。
 - データソースノード: DataSource0、DataSource1、DataSource2 などです。
 - 変換ノード: Transform0、Transform1、Transform2 などです。
- New code generation — '_node1','_node2'などが添付されたノードの[Node properties]タブ上で指定された名前を使用します。例えば、S3bucket_node1、ApplyMapping_node2、S3bucket_node2、MyCustomNodeName_node

各エラーコード generator のさらなる詳細については、「[スクリプトコードの生成](#)」を参照してください。

次の例は、コードボックスに入力するコードの形式を示しています。

Python

次の例では、まず DynamicFrame を受信し、それを DataFrame に変換してネイティブフィルターメソッド (1000 票を超えるレコードのみを保持) を適用しています。その後、返す前に DynamicFrame に再び変換しています。

```
def FilterHighVoteCounts (glueContext, dfc) -> DynamicFrameCollection:
    df = dfc.select(list(dfc.keys())[0]).toDF()
    df_filtered = df.filter(df["vote_count"] > 1000)
    dyf_filtered = DynamicFrame.fromDF(df_filtered, glueContext, "filter_votes")
    return(DynamicFrameCollection({"CustomTransform0": dyf_filtered}, glueContext))
```

Scala

次の例では、まず DynamicFrame を受信し、それを DataFrame に変換してネイティブフィルターメソッド (1000 票を超えるレコードのみを保持) を適用しています。その後、返す前に DynamicFrame に再び変換しています。

```
object FilterHighVoteCounts {
    def execute(glueContext : GlueContext, input : Seq[DynamicFrame]) :
    Seq[DynamicFrame] = {
        val frame = input(0).toDF()
        val filtered = DynamicFrame(frame.filter(frame("vote_count") > 1000),
        glueContext)
        Seq(filtered)
    }
}
```

カスタム変換ノードでスキーマを編集する

カスタム transform ノードを使用する場合、AWS Glue Studio は、transform によって作成された出力スキーマを自動で推測することはできません。スキーマエディタを使用して、カスタム変換コードによって実装されるスキーマの変更を記述できます。

カスタムコードノードでは、カスタムコードの入力として DynamicFrame を提供する任意の数の親ノードを持つことができます。カスタムコードノードでは、DynamicFrames のコレクションが返されます。入力として使用される各 DynamicFrame には、スキーマが関連付けられています。カスタムコードノードによって返される各 DynamicFrame を記述するスキーマを追加する必要があります。

Note

カスタム変換で独自のスキーマを設定した場合、AWS Glue Studio は、前のノードからのスキーマ継承を行いません。スキーマを更新するには、カスタム変換のノードを選択した上で、[Data preview] (データプレビュー) タブを開きます。プレビューが生成されたら、[Use Preview Schema] (プレビュースキーマを使用) をクリックします。その後、このスキーマがプレビューデータを使用するスキーマに置き換えられます。

カスタム変換ノードの入カスキーマを編集するには

1. ジョブ図でカスタム変換ノードを選択した状態で、ノードの詳細パネルの [Output schema] (出カスキーマ) タブを選択します。
2. [Edit] (編集) を選択して、スキーマを変更します。

配列やオブジェクトなどのネストされたデータのプロパティキーがある場合、各スキーマパネルの右上部にある [Expand-Rows] (行の展開) アイコン

()

を選択して、子データのプロパティキーのリストを展開します。このアイコンは、選択すると [Collapse-Rows] (列の折りたたみ) アイコン

()

に変わり、子のプロパティキーのリストを折りたたむことができます。

3. ページ右側のセクションにある次のアクションを使用して、スキーマを変更します。
 - プロパティキーの名前を変更するには、プロパティキーの [Key] (キー) テキストボックスにカーソルを置き、新しい名前を入力します。
 - プロパティキーのデータ型を変更するには、リストを使用して、プロパティキーの新しいデータ型を選択します。
 - トップレベルの新しいプロパティキーをスキーマに追加するには、[Cancel] (キャンセル) ボタンの左側にある [Overflow] (オーバーフロー)

- (...)
アイコンをクリックして、[Add root key] (ルートキーの追加) を選択します。
 - スキーマに子のプロパティキーを追加するには、親キーに関連付けられている [Add-Key] (キーの追加)
(☞)
アイコンをクリックします。子のキーの名前を入力し、データ型を選択します。
 - スキーマからプロパティキーを削除するには、キー名の右端にある [Remove] (削除) アイコン
(☐)
をクリックします。
4. カスタム変換コードで DynamicFrames を複数使用していない場合、他の出力スキーマを追加できます。
- 新しい空のスキーマを追加するには、[Overflow] (オーバーフロー)
(...)
アイコンをクリックし、[Add output schema] (出力スキーマの追加) を選択します。
 - 既存のスキーマを新しい出力スキーマにコピーするには、コピーするスキーマがスキーマセレクタに表示されていることを確認します。[Overflow] (オーバーフロー)
(...)
アイコンをクリックして、[Duplicate] (複製) を選択します。
- 出力スキーマを削除する場合は、コピーするスキーマがスキーマセレクタに表示されていることを確認します。[Overflow] (オーバーフロー)
(...)
アイコンをクリックして、[Delete] (削除) を選択します。
5. 新しいスキーマに新しいルートキーを追加するか、複製したキーを編集します。
6. 出力スキーマを変更する場合は、[Apply] (適用) ボタンをクリックして変更を保存し、スキーマエディタを終了します。

変更を保存しない場合は、[Cancel] (キャンセル) ボタンをクリックします。

カスタム変換の出力を設定する

カスタムコード変換では、1つの結果セットに DynamicFrame が1つだけであっても、DynamicFrames のコレクションが返されます。

カスタム変換ノードからの出力を処理するには

1. カスタム変換ノードを親ノードとして持つ `SelectFromCollection` 変換ノードを追加します。この変換を更新して、使用するデータセットを指定します。詳細については、「[SelectFromCollection を使用して保持するデータセットを選択する](#)」を参照してください。
2. カスタム変換ノードによって生成される他の `DynamicFrames` を使用する場合は、`SelectFromCollection` 変換をジョブ図に追加します。

カスタム変換ノードを追加してフライトのデータセットを複数のデータセットに分割し、フライトの日付や番号などの各出カスキーマで識別プロパティキーを複製するシナリオを考えてみましょう。各出カスキーマに、カスタム変換ノードを親として持つ `SelectFromCollection` 変換ノードを追加します。

3. (オプション) その後、各 `SelectFromCollection` 変換ノードをジョブ内の他のノードの入力として、またはデータターゲットノードの親として使用できます。

AWS Glue カスタムビジュアル変換

カスタムビジュアル変換を使用すると、変換を作成して AWS Glue Studio ジョブで使用できるようになります。ETL 開発者は、コーディングに慣れていない場合がありますが、カスタムビジュアル変換により、AWS Glue Studio インターフェイスを使用して増え続ける変換のライブラリを検索して使用することができます。

カスタムのビジュアル変換を作成し、Amazon S3 にアップロードして、AWS Glue Studio のビジュアルエディタでこれらのジョブを処理できるようにすることができます。

トピック

- [カスタムビジュアル変換を開始する](#)
- [Step 1. JSON 設定ファイルを作成する](#)
- [Step 2. 変換ロジックを実装する](#)
- [ステップ 3. AWS Glue Studio でのカスタムビジュアル変換の検証とトラブルシューティング](#)
- [Step 4. 必要に応じてカスタムビジュアル変換を更新する](#)
- [Step 5. AWS Glue Studio でカスタムビジュアル変換を使用する](#)
- [使用例](#)
- [カスタムビジュアルスクリプトの例](#)
- [動画](#)

カスタムビジュアル変換を開始する

カスタムビジュアル変換を作成するには、次の手順を実行します。

- Step 1. JSON 設定ファイルを作成する
- Step 2. 変換ロジックを実装する
- ステップ 3。カスタムビジュアル変換を検証する
- Step 4. 必要に応じてカスタムビジュアル変換を更新する
- Step 5. AWS Glue Studio でカスタムビジュアル変換を使用する

Amazon S3 バケットをセットアップして開始し、ステップ 1。JSON 設定ファイルを作成するに進みます。

前提条件

お客様提供の変換はお客様の AWS アカウント内にあります。そのアカウントは変換を所有しているため、変換を表示 (検索、使用)、編集、削除するすべてのアクセス許可を持っています。

AWS Glue Studio でカスタム変換を使用するには、その AWS アカウントの Amazon S3 アセットバケットに次の 2 つのファイルを作成してアップロードする必要があります。

- Python ファイル – 変換関数が含まれています
- JSON ファイル – 変換を記述します。これは、設定ファイルとも呼ばれ、変換の定義に必要です。

ファイルをペアリングするには、両方に同じ名前を使用します。例:

- myTransform.json
- myTransform.py

必要に応じて、アイコンを含む SVG ファイルを指定することで、カスタムビジュアル変換にカスタムアイコンを付加できます。ファイルをペアリングするには、次のように、アイコンに同じ名前を使用します。

- myTransform.svg

AWS Glue Studio では、それぞれのファイル名を使用して自動的に照合されます。ファイル名を既存のモジュールと同じにすることはできません。

推奨される変換ファイル名の命名規則

AWS Glue Studio では、ファイルはモジュール (例えば、`import myTransform`) としてジョブスクリプトにインポートされます。そのため、ファイル名は Python の変数名 (識別子) に設定されているのと同じ命名規則に従う必要があります。具体的には、文字またはアンダースコアで始まり、文字、数字、および/またはアンダースコアのみで構成されている必要があります。

Note

予期しない実行時の問題を避けるため、変換ファイル名がロードされている既存の Python モジュール (`sys`, `array`, `copy` など) と競合していないことを確認してください。

Amazon S3 バケットのセットアップ

作成した変換は Amazon S3 に保存され、AWS アカウントによって所有されます。すべてのジョブスクリプトが現在保存されている Amazon S3 アセットフォルダ (例えば、`s3://aws-glue-assets-
<accountid>-<region>/transforms`) にファイル (`json` と `py`) をアップロードするだけで、新しいカスタムビジュアル変換を作成できます。カスタムアイコンを使用している場合は、それもアップロードします。デフォルトでは、AWS Glue Studio は同じ S3 バケット内の `/transforms` フォルダからすべての `.json` ファイルを読み取ります。

Step 1. JSON 設定ファイルを作成する

カスタムビジュアル変換を定義して記述するには、JSON 設定ファイルが必要です。設定ファイルのスキーマは次の通りです。

JSON ファイル構造

フィールド

- `name`: `string` – (必須) 変換の識別に使用される変換システム名。Python の変数名 (識別子) に設定されているのと同じ命名規則に従ってください。具体的には、文字またはアンダースコアで始まり、文字、数字、および/またはアンダースコアのみで構成されている必要があります。
- `displayName`: `string` – (オプション) AWS Glue Studio ビジュアルジョブエディタに表示される変換の名前。`displayName` を指定しない場合は、AWS Glue Studio では `name` が変換の名前として使用されます。
- `description`: `string` – (オプション) 変換の説明が AWS Glue Studio に表示され、検索できます。

- `functionName`: `string` – (必須) Python 関数名が Python スクリプトで呼び出す関数を識別するために使用されます。
- `path`: `string` – (オプション) Python ソースファイルへの Amazon S3 のフルパス。指定しない場合は、AWS Glue では、ファイル名マッチングを使用して `.json` ファイルと `.py` ファイルがペアリングされます。例えば、JSON ファイルの名前が `myTransform.json` の場合、同じ Amazon S3 ロケーションにある `myTransform.py` という Python ファイルとペアになります。
- `parameters`: `Array of TransformParameter object` – (オプション) AWS Glue Studio ビジュアルエディタで設定するときに表示されるパラメータのリスト。

TransformParameter フィールド

- `name`: `string` – (必須) ジョブスクリプト内で名前付き引数として Python 関数に渡されるパラメータ名。Python の変数名 (識別子) に設定されているのと同じ命名規則に従ってください。具体的には、文字またはアンダースコアで始まり、文字、数字、および/またはアンダースコアのみで構成されている必要があります。
- `displayName`: `string` – (オプション) AWS Glue Studio ビジュアルジョブエディタに表示される変換の名前。`displayName` を指定しない場合は、AWS Glue Studio では `name` が変換の名前として使用されます。
- `type`: `string` – (必須) 一般的な Python データ型を受け入れるパラメータ型。有効な値: `'str'` | `'int'` | `'float'` | `'list'` | `'bool'`。
- `isOptional`: `boolean` – (オプション) パラメータがオプションかどうかを決定します。デフォルトでは、すべてのパラメータが必須です。
- `description`: `string` – (オプション) ユーザーによる変換パラメータの設定に役立つ説明が AWS Glue Studio に表示されます。
- `validationType`: `string` – (オプション) このパラメータの検証方法を定義します。現在、正規表現のみをサポートしています。デフォルトで検証タイプは `RegularExpression` に設定されています。
- `validationRule`: `string` – (オプション) `validationType` が `RegularExpression` に設定されている場合、送信前にフォーム入力を検証するために使用される正規表現。正規表現の構文は [RegExp EcmaScript の仕様](#) に合致している必要があります。
- `validationMessage`: `string` – (オプション) 検証が失敗したときに表示するメッセージ。
- `listOptions`: `An array of TransformParameterListOption object` または `string`、もしくは「column」という文字列の値 — (オプション) 選択、または複数選択の UI コントロールに表示するためのオプション。カンマ区切り値 (CSV) のリスト、または強く型付けさ

れた `TransformParameterListOption` 型の JSON オブジェクトを受け入れます。文字列の値「column」を指定すると、親ノードのスキーマから、列のリストを動的に入力することもできます。

- `listType: string` – (オプション) `type = 'list'` のオプションタイプを定義します。有効な値: `'str'` | `'int'` | `'float'` | `'list'` | `'bool'`。一般的な Python データ型を受け入れるパラメータ型。

TransformParameterListOption フィールド

- `value: string | int | float | bool` – (必須) オプション値。
- `label: string` – (オプション) 選択ドロップダウンに表示されるオプションラベル。

AWS Glue Studio での変換パラメータ

.json ファイルで `isOptional` と指定しない限り、デフォルトでパラメータは必須です。AWS Glue Studio では、パラメータが [Transform] (変換) タブに表示されます。この例は、E メールアドレス、電話番号、年齢、性別、出身国などのユーザー定義パラメータを示しています。

The screenshot displays the AWS Glue Studio interface. On the left, a visual graph shows a 'Data source - S3 bucket Amazon S3' connected to a 'Transform - Dynamic Trans... My Transform' node. On the right, the 'Transform' tab is active, showing a configuration panel with the following fields:

- Email Address**: Enter your work email address below (text input)
- Phone Number**: Enter your mobile phone number below (text input)
- Your age**: (text input)
- Your gender**: (dropdown menu)
- Your origin country? - optional**: What country were you born in? (dropdown menu with 'Choose options' selected)
- Do you want to receive promotional newsletter from us? - optional**: (checkbox)

`validationRule` パラメータを指定し、`validationMessage` で検証メッセージを指定することで、JSON ファイルで正規表現を使用する際に AWS Glue Studio で一部の検証を必須にできます。

```
"validationRule": "^\\((?\\d{3})\\)?[- ]?(\\d{3})[- ]?(\\d{4})$",
"validationMessage": "Please enter a valid US number"
```

Note

検証はブラウザで行われるため、正規表現の構文は [RegExp Ecmascript の仕様](#) に合致している必要があります。Python 構文は、これらの正規表現ではサポートされていません。

検証を追加することで、ユーザーが誤ったユーザー入力でジョブを保存するのを防ぐことができます。次の例のように AWS Glue Studio に検証メッセージが表示されます。



The screenshot shows the 'Transform' tab in AWS Glue Studio. Under the 'Node properties' section, there is a field for 'Email Address' with the placeholder text 'Enter your work email address below'. The input field contains the text 'wrongEmail.com'. Below the input field, a red error message is displayed: 'Please enter a valid email address'.

パラメータは、パラメータ設定に基づいて AWS Glue Studio に表示されます。

- type が str、int、または float の場合、テキスト入力フィールドが表示されます。例えば、スクリーンショットは「E メールアドレス」と「年齢」パラメータの入力フィールドを示しています。

Email Address

Enter your work email address below

Your age

- type が bool の場合、チェックボックスが表示されます。

Do you want to receive promotional newsletter from us?

- type が str であり、listOptions が指定されていると、単一選択リストが表示されます。

Your gender

Male	▲
Male	✓
Female	
Other	

- type が list であり、listOptions と listType が指定されていると、複数選択リストが表示されます。

Country recently visited - optional

What countries did you visit in the past 2 years?

Choose options	▲
🔍	
<input type="checkbox"/> Iceland	
<input checked="" type="checkbox"/> India	
<input type="checkbox"/> Indor India	
<input type="checkbox"/> Iran	
<input type="checkbox"/> Iraq	
<input type="checkbox"/> Ireland	
<input checked="" type="checkbox"/> Israel	
<input type="checkbox"/> Italy	
<input type="checkbox"/> Jamaica	
<input type="checkbox"/> Japan	

列のセレクトをパラメータとして表示

スキーマから列をユーザーが選択するように設定されている場合は、列セレクトを表示することで、ユーザーは列の名前を入力する必要がなくなります。listOptions フィールドを「column」に設定すると、親ノードの出カスキーマに基づいて、列セレクトを AWS Glue Studio が動的に表示します。AWS Glue Studio は単一の列、または複数列のセレクトを表示できます。

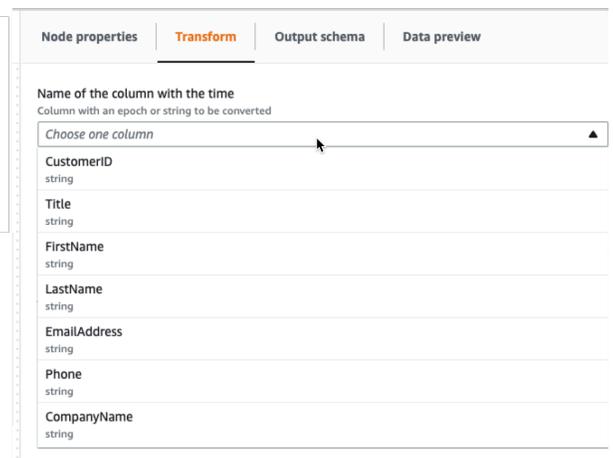
次の例でスキーマを使用しています。

Node properties	Data source properties - S3	Output schema	Data preview
Schema			Edit
Key	Data type	Partition	
CustomerID	string	-	
Title	string	-	
FirstName	string	-	
LastName	string	-	
EmailAddress	string	-	
Phone	string	-	
CompanyName	string	-	

単一の列を表示するように、Custom Visual Transform パラメータを定義するには

1. JSON ファイルで、parameters オブジェクトの listOptions の値を、「column」に設定します。これにより、ユーザーが AWS Glue Studio の選択リストから列を選択できるようになります。

```
{
  "name": "mb_to_timestamp",
  "displayName": "MB Convert column to timestamp",
  "description": "Convert a timestamp string or a system epoch column into a new timestamp type column.",
  "functionName": "mb_to_timestamp",
  "parameters": [
    {
      "name": "colName",
      "displayName": "Name of the column with the time",
      "type": "str",
      "listOptions": "column",
      "description": "Column with an epoch or string to be converted"
    }
  ],
}
```



2. 次のようにパラメータを定義すると、複数の列を選択することもできます。

- listOptions: "column"
- type: "list"

```

{
  "name": "mb_to_timestamp",
  "displayName": "MB Convert column to timestamp",
  "description": "Convert a timestamp string or a system epoch column into a new timestamp type column.",
  "functionName": "mb_to_timestamp",
  "parameters": [
    {
      "name": "colNames",
      "displayName": "Name of the column with the time",
      "type": "list",
      "listOptions": "column",
      "listType": "str",
      "description": "Column with an epoch or string to be converted"
    }
  ]
}

```

Step 2. 変換ロジックを実装する

Note

カスタムビジュアル変換は Python スクリプトのみをサポートしています。Scala はサポートされていません。

.json 設定ファイルで定義された関数を実装するコードを追加するには、Python ファイルを .json ファイルと同じ場所に同じ名前でも「.py」拡張子を付けて配置することをお勧めします。AWS Glue Studio では、.json ファイルと .py ファイルが自動的にペアリングされるため、設定ファイルで Python ファイルのパスを指定する必要はありません。

Python ファイルに、指定されたパラメータを設定した宣言済み関数を追加し、その関数を DynamicFrame で使用するために登録します。Python ファイルの例を次に示します。

```

from awsglue import DynamicFrame

# self refers to the DynamicFrame to transform,
# the parameter names must match the ones defined in the config
# if it's optional, need to provide a default value
def myTransform(self, email, phone, age=None, gender="",
                country="", promotion=False):
    resulting_dynf = # do some transformation on self
    return resulting_dynf

DynamicFrame.myTransform = myTransform

```

Python コードの開発とテストを最も早く行うために、AWS Glue ノートブックの使用をお勧めします。「[AWS Glue Studio 中でのノートブックの使用開始](#)」を参照してください。

変換ロジックの実装方法を説明するために、以下の例でのカスタムビジュアル変換は、入力データをフィルタリングして特定の米国の州に関連するデータのみを保持する変換になっています。json ファイルには `functionName` のパラメータが `custom_filter_state` として含まれ、2 つの引数 (「str」型の「state」と「colName」) が含まれています。

.json 設定ファイルの例は次のとおりです。

```
{
  "name": "custom_filter_state",
  "displayName": "Filter State",
  "description": "A simple example to filter the data to keep only the state indicated.",
  "functionName": "custom_filter_state",
  "parameters": [
    {
      "name": "colName",
      "displayName": "Column name",
      "type": "str",
      "description": "Name of the column in the data that holds the state postal code"
    },
    {
      "name": "state",
      "displayName": "State postal code",
      "type": "str",
      "description": "The postal code of the state whole rows to keep"
    }
  ]
}
```

Python でコンパニオンスクリプトを実装するには

1. AWS Glue ノートブックを起動し、セッションを開始するために提供された最初のセルを実行します。最初のセルを実行すると、必要な基本コンポーネントが作成されます。
2. 例で説明されているようにフィルタリングを実行する関数を作成し、`DynamicFrame` に登録します。以下のコードをコピーして、AWS Glue ノートブックのセルに貼り付けます。

```
from awsglue import DynamicFrame

def custom_filter_state(self, colName, state):
    return self.filter(lambda row: row[colName] == state)

DynamicFrame.custom_filter_state = custom_filter_state
```

3. サンプルデータを作成またはロードして、同じセルまたは新しいセルでコードをテストします。サンプルデータを新しいセルに追加する場合は、セルを実行することを忘れないでください。例:

```
# A few of rows of sample data to test
data_sample = [
    {"state": "CA", "count": 4},
    {"state": "NY", "count": 2},
    {"state": "WA", "count": 3}
]
df1 = glueContext.sparkSession.sparkContext.parallelize(data_sample).toDF()
dynf1 = DynamicFrame.fromDF(df1, glueContext, None)
```

4. さまざまな引数で次のようにテストして、「custom_filter_state」を検証します。

```
[14]: dynf1.custom_filter_state("state", "NY").show()
      {"count": 2, "state": "NY"}
```

5. いくつかのテストを実行した後、.py 拡張子を付けてコードを保存し、.py ファイルに .json ファイル名と同じ名前を付けます。.py ファイルと .json ファイルは同じ変換フォルダにある必要があります。

次のコードをコピーしてファイルに貼り付け、.py ファイル拡張子に変更します。

```
from awsglue import DynamicFrame

def custom_filter_state(self, colName, state):
    return self.filter(lambda row: row[colName] == state)

DynamicFrame.custom_filter_state = custom_filter_state
```

6. AWS Glue Studio でビジュアルジョブを開き、使用可能な [Transforms] (変換) のリストから変換を選択して変換をジョブに追加します。

この変換を Python スクリプトコードで再利用するには、Amazon S3 パスを [Referenced files path] (参照されるファイルパス) の下のジョブの .py ファイルに追加し、スクリプト内の python ファイルの名前 (拡張子なし) をファイルの先頭に追加してインポートします。例えば、import <ファイル名 (拡張子なし)>

ステップ 3。AWS Glue Studio でのカスタムビジュアル変換の検証とトラブルシューティング

カスタムビジュアル変換が AWS Glue Studio にロードされる前に AWS Glue Studio で JSON 設定ファイルが検証されます。検証には次の項目が含まれます。

- 必須フィールドの有無
- JSON 形式の検証
- パラメータの誤りまたは無効
- 同じ Amazon S3 パスに .py ファイルと .json ファイルの両方が存在すること
- .py ファイル名と .json のファイル名の一致

検証が成功すると、変換がビジュアルエディタの利用可能な [Actions] (アクション) のリストに表示されます。カスタムアイコンが指定されている場合は、[アクション] の横に表示されているはずで

検証に失敗した場合、AWS Glue Studio では、カスタムビジュアル変換はロードされません。

Step 4. 必要に応じてカスタムビジュアル変換を更新する

一度作成して使用すると、次のように変換が対応する JSON 定義に従っている限り、変換スクリプトを更新できます。

- DynamicFrame に割り当てるときに使用される名前は、JSON の functionName と一致している必要があります。
- 関数の引数は、[Step 1. JSON 設定ファイルを作成する](#) で説明されているように JSON ファイルで定義する必要があります。

- Python ファイルの Amazon S3 パスは、ジョブが直接それに依存しているため、変更できません。

Note

更新が必要な場合は、スクリプトと .json ファイルの更新を整合させ、ビジュアルジョブを新しい変換で正しく再保存します。更新後にビジュアルジョブを保存しないと、更新の適用と検証が行われません。Python スクリプトファイルの名前が変更されたり、.json ファイルの横に配置されなかったりする場合は、.json ファイルにフルパスを指定する必要があります。

カスタムアイコン

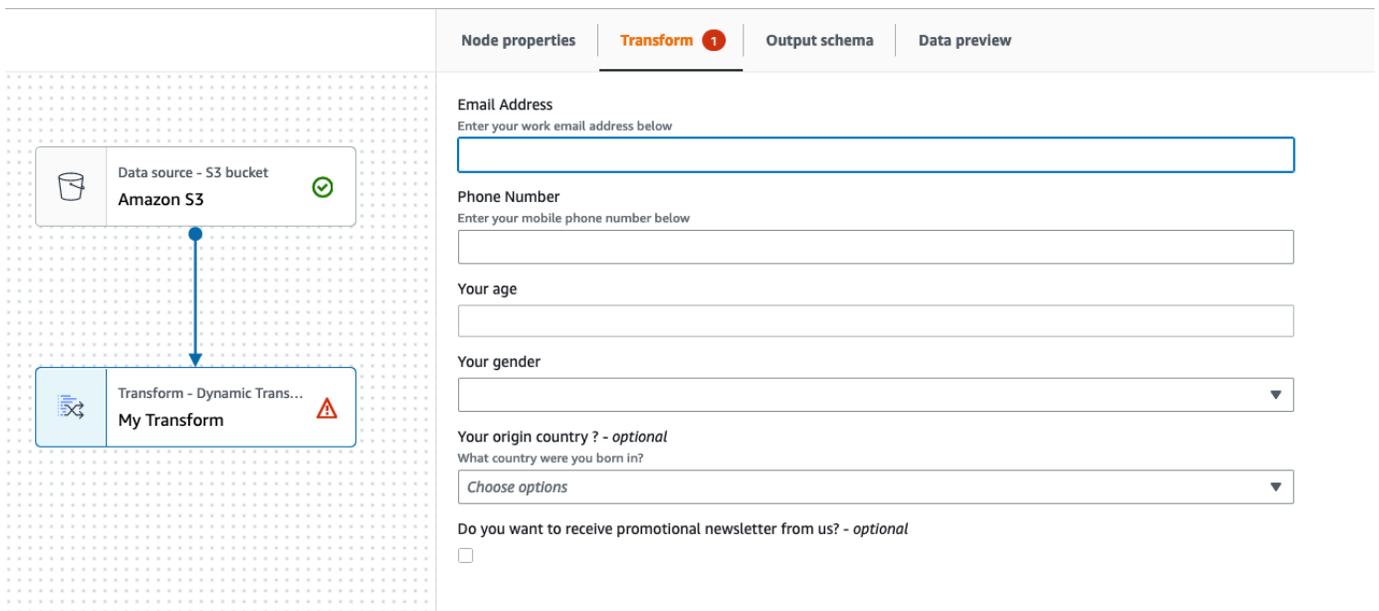
[アクション] のデフォルトアイコンではワークフローの一部として視覚的に区別できないと判断した場合は、「[the section called “カスタムビジュアル変換を開始する”](#)」で説明されているようにカスタムアイコンを指定できます。Amazon S3 でホストされている対応する SVG を更新することで、アイコンを更新できます。

最良の結果を得るには、Cloudscape Design System のガイドラインに従って、32 x 32 ピクセルで表示されるように画像をデザインしてください。Cloudscape ガイドラインの詳細については、[Cloudscape のドキュメント](#)を参照してください。

Step 5. AWS Glue Studio でカスタムビジュアル変換を使用する

AWS Glue Studio でカスタムビジュアル変換を使用するには、設定ファイルとソースファイルをアップロードし、[Action] (アクション) メニューから変換を選択します。値や入力が必要なパラメータはどれも [Transform] (変換) タブで利用できます。

1. 2つのファイル (Python ソースファイルと JSON 設定ファイル) を、ジョブスクリプトが保存されている Amazon S3 アセットフォルダにアップロードします。デフォルトでは、AWS Glue は同じ Amazon S3 バケット内の /transforms フォルダからすべての JSON ファイルを取得します。
2. [Action] (アクション) メニューから、カスタムビジュアル変換を選択します。この名前は、変換の `displayName` または .json 設定ファイルで指定した名前になります。
3. 設定ファイルに設定されたパラメータの値を入力します。



The screenshot shows the AWS Glue console interface. On the left, a workflow diagram displays a data source node labeled 'Data source - S3 bucket Amazon S3' with a green checkmark, connected by a blue arrow to a transform node labeled 'Transform - Dynamic Trans... My Transform' with a red warning triangle. On the right, the 'Transform' configuration page is open, showing fields for 'Email Address', 'Phone Number', 'Your age', 'Your gender', 'Your origin country? - optional', and a checkbox for 'Do you want to receive promotional newsletter from us? - optional'.

使用例

.json 設定ファイルに含めることが可能なすべてのパラメータの例を次に示します。

```
{
  "name": "MyTransform",
  "displayName": "My Transform",
  "description": "This transform description will be displayed in UI",
  "functionName": "myTransform",
  "parameters": [
    {
      "name": "email",
      "displayName": "Email Address",
      "type": "str",
      "description": "Enter your work email address below",
      "validationType": "RegularExpression",
      "validationRule": "^\\w+([\\.-]?\\w+)*@\\w+([\\.-]?\\w+)*(\\.\\w{2,3})+$",
      "validationMessage": "Please enter a valid email address"
    },
    {
      "name": "phone",
      "displayName": "Phone Number",
      "type": "str",
      "description": "Enter your mobile phone number below",
      "validationRule": "^\\((?\\d{3})\\)?[- ]?(\\d{3})[- ]?(\\d{4})$",
      "validationMessage": "Please enter a valid US number"
    }
  ]
}
```

```

    },
    {
      "name": "age",
      "displayName": "Your age",
      "type": "int",
      "isOptional": true
    },
    {
      "name": "gender",
      "displayName": "Your gender",
      "type": "str",
      "listOptions": [
        {"label": "Male", "value": "male"},
        {"label": "Female", "value": "female"},
        {"label": "Other", "value": "other"}
      ],
      "isOptional": true
    },
    {
      "name": "country",
      "displayName": "Your origin country ?",
      "type": "list",
      "listOptions": "Afghanistan,Albania,Algeria,American
Samoa,Andorra,Angola,Anguilla,Antarctica,Antigua and
Barbuda,Argentina,Armenia,Aruba,Australia,Austria,Azerbaijan,Bahamas,Bahrain,Bangladesh,Barbad
and Herzegovina,Botswana,Bouvet Island,Brazil,British Indian Ocean Territory,Brunei
Darussalam,Bulgaria,Burkina Faso,Burundi,Cambodia,Cameroon,Canada,Cape
Verde,Cayman Islands,Central African Republic,Chad,Chile,China,Christmas
Island,Cocos (Keeling Islands),Colombia,Comoros,Congo,Cook Islands,Costa
Rica,Cote D'Ivoire (Ivory Coast),Croatia (Hrvatska,Cuba,Cyprus,Czech
Republic,Denmark,Djibouti,Dominica,Dominican Republic,East Timor,Ecuador,Egypt,El
Salvador,Equatorial Guinea,Eritrea,Estonia,Ethiopia,Falkland Islands (Malvinas),Faroe
Islands,Fiji,Finland,France,France,Metropolitan,French Guiana,French Polynesia,French
Southern
Territories,Gabon,Gambia,Georgia,Germany,Ghana,Gibraltar,Greece,Greenland,Grenada,Guadeloupe,G
Bissau,Guyana,Haiti,Heard and McDonald Islands,Honduras,Hong
Kong,Hungary,Iceland,India,Indonesia,Iran,Iraq,Ireland,Israel,Italy,Jamaica,Japan,Jordan,Kazak
(North),Korea
(South),Kuwait,Kyrgyzstan,Laos,Latvia,Lebanon,Lesotho,Liberia,Libya,Liechtenstein,Lithuania,Lu
Islands,Martinique,Mauritania,Mauritius,Mayotte,Mexico,Micronesia,Moldova,Monaco,Mongolia,Mont
Antilles,New Caledonia,New Zealand,Nicaragua,Niger,Nigeria,Niue,Norfolk
Island,Northern Mariana Islands,Norway,Oman,Pakistan,Palau,Panama,Papua
New Guinea,Paraguay,Peru,Philippines,Pitcairn,Poland,Portugal,Puerto
Rico,Qatar,Reunion,Romania,Russian Federation,Rwanda,Saint Kitts and Nevis,Saint

```

```

Lucia,Saint Vincent and The Grenadines,Samoa,San Marino,Sao Tome and Principe,Saudi
Arabia,Senegal,Seychelles,Sierra Leone,Singapore,Slovak Republic,Slovenia,Solomon
Islands,Somalia,South Africa,S. Georgia and S. Sandwich Isls.,Spain,Sri
Lanka,St. Helena,St. Pierre and Miquelon,Sudan,Suriname,Svalbard and Jan Mayen
Islands,Swaziland,Sweden,Switzerland,Syria,Tajikistan,Tanzania,Thailand,Togo,Tokelau,Tonga,Tri
and Tobago,Tunisia,Turkey,Turkmenistan,Turks and Caicos
Islands,Tuvalu,Uganda,Ukraine,United Arab Emirates,United Kingdom
(Britain / UK),United States of America (USA),US Minor Outlying
Islands,Uruguay,Uzbekistan,Vanuatu,Vatican City State (Holy See),Venezuela,Viet
Nam,Virgin Islands (British),Virgin Islands (US),Wallis and Futuna Islands,Western
Sahara,Yemen,Yugoslavia,Zaire,Zambia,Zimbabwe",
    "description": "What country were you born in?",
    "listType": "str",
    "isOptional": true
},
{
    "name": "promotion",
    "displayName": "Do you want to receive promotional newsletter from us?",
    "type": "bool",
    "isOptional": true
}
]
}

```

カスタムビジュアルスクリプトの例

以下の例で行われる変換はどれも等価です。ただし、2 番目の例 (SparkSQL) が最もクリーンで効率的です。その次が Pandas UDF で、最後は最初の例の低レベルのマッピングです。次の例は、2 つの列を合計する単純な変換の完全な例です。

```

from awsglue import DynamicFrame

# You can have other auxiliary variables, functions or classes on this file, it won't
# affect the runtime
def record_sum(rec, col1, col2, resultCol):
    rec[resultCol] = rec[col1] + rec[col2]
    return rec

# The number and name of arguments must match the definition on json config file
# (expect self which is the current DynamicFrame to transform
# If an argument is optional, you need to define a default value here

```

```
# (resultCol in this example is an optional argument)
def custom_add_columns(self, col1, col2, resultCol="result"):
    # The mapping will alter the columns order, which could be important
    fields = [field.name for field in self.schema()]
    if resultCol not in fields:
        # If it's a new column put it at the end
        fields.append(resultCol)
    return self.map(lambda record: record_sum(record, col1, col2,
resultCol)).select_fields(paths=fields)

# The name we assign on DynamicFrame must match the configured "functionName"
DynamicFrame.custom_add_columns = custom_add_columns
```

次の例は、SparkSQL API を利用した同等の変換です。

```
from awsglue import DynamicFrame

# The number and name of arguments must match the definition on json config file
# (expect self which is the current DynamicFrame to transform
# If an argument is optional, you need to define a default value here
# (resultCol in this example is an optional argument)
def custom_add_columns(self, col1, col2, resultCol="result"):
    df = self.toDF()
    return DynamicFrame.fromDF(
        df.withColumn(resultCol, df[col1] + df[col2]) # This is the conversion logic
        , self.glue_ctx, self.name)

# The name we assign on DynamicFrame must match the configured "functionName"
DynamicFrame.custom_add_columns = custom_add_columns
```

次の例では同じ変換を使用していますが、Pandas UDF を使用しているため、プレーン UDF を使用するよりも効率的です。Pandas UDF を書く方法の詳細については、[Apache Spark SQL のドキュメント](#)を参照してください。

```
from awsglue import DynamicFrame
import pandas as pd
from pyspark.sql.functions import pandas_udf

# The number and name of arguments must match the definition on json config file
```

```
# (expect self which is the current DynamicFrame to transform
# If an argument is optional, you need to define a default value here
# (resultCol in this example is an optional argument)
def custom_add_columns(self, col1, col2, resultCol="result"):
    @pandas_udf("integer") # We need to declare the type of the result column
    def add_columns(value1: pd.Series, value2: pd.Series) # pd.Series:
        return value1 + value2

    df = self.toDF()
    return DynamicFrame.fromDF(
        df.withColumn(resultCol, add_columns(col1, col2)) # This is the conversion
        logic
        , self.glue_ctx, self.name)

# The name we assign on DynamicFrame must match the configured "functionName"
DynamicFrame.custom_add_columns = custom_add_columns
```

動画

次の動画では、ビジュアルカスタム変換の概要と使用方法を紹介しています。

データレイクフレームワークを AWS Glue Studio で使用する

概要

オープンソースのデータレイクフレームワークは、Amazon S3 上に構築されたデータレイクに保存されたファイルのインクリメンタルデータ処理を簡素化します。AWS Glue 3.0 以降では、次のオープンソースのデータレイクストレージフレームワークをサポートしています。

- Apache Hudi
- Linux Foundation Delta Lake
- Apache Iceberg

AWS Glue 4.0 では、AWS Glue はこれらのフレームワークをネイティブにサポートしているため、トランザクションが一貫した方法で、Amazon S3 に保存したデータを読み書きできます。AWS Glue ジョブでこれらのフレームワークを使用する場合でも、別のコネクタをインストールしたり、設定手順を追加で実行したりする必要はありません。

データレイクフレームワークは、Spark Script Editor ジョブから AWS Glue Studio 内のソースまたはターゲットとして使用できます。Apache Hudi、Apache Iceberg、Delta Lake の使用方法の詳細に

については、「[AWS Glue ETL ジョブでデータレイクフレームワークを使用する](#)」を参照してください。

AWS Glue ストリーミングソースからのオープンテーブルフォーマットの作成

AWS Glue ストリーミング ETL ジョブは、ストリーミングソースからのデータを継続的に消費し、転送中のデータをクリーンアップおよび変換して、数秒で分析できるようにします。

AWS は、お客様のニーズをサポートする幅広いサービスを提供します。AWS Database Migration Service などのデータベースレプリケーションサービスは、ソースシステムから Amazon S3 にデータをレプリケートできます。Amazon S3 は通常、データレイクのストレージレイヤーをホストします。オンラインソースアプリケーションをサポートするリレーショナルデータベース管理システム (RDBMS) に更新を適用するのは簡単ですが、この CDC プロセスをデータレイクに適用するのは困難です。オープンソースのデータ管理フレームワークは、増分データ処理とデータパイプライン開発を簡素化し、この問題を解決するための優れたオプションです。

詳細については、以下を参照してください。

- [AWS Glue ストリーミングを使用して、Apache Hudi ベースのほぼリアルタイムのトランザクションデータレイクを作成する](#)
- [Build a real-time GDPR-aligned Apache Iceberg data lake](#)

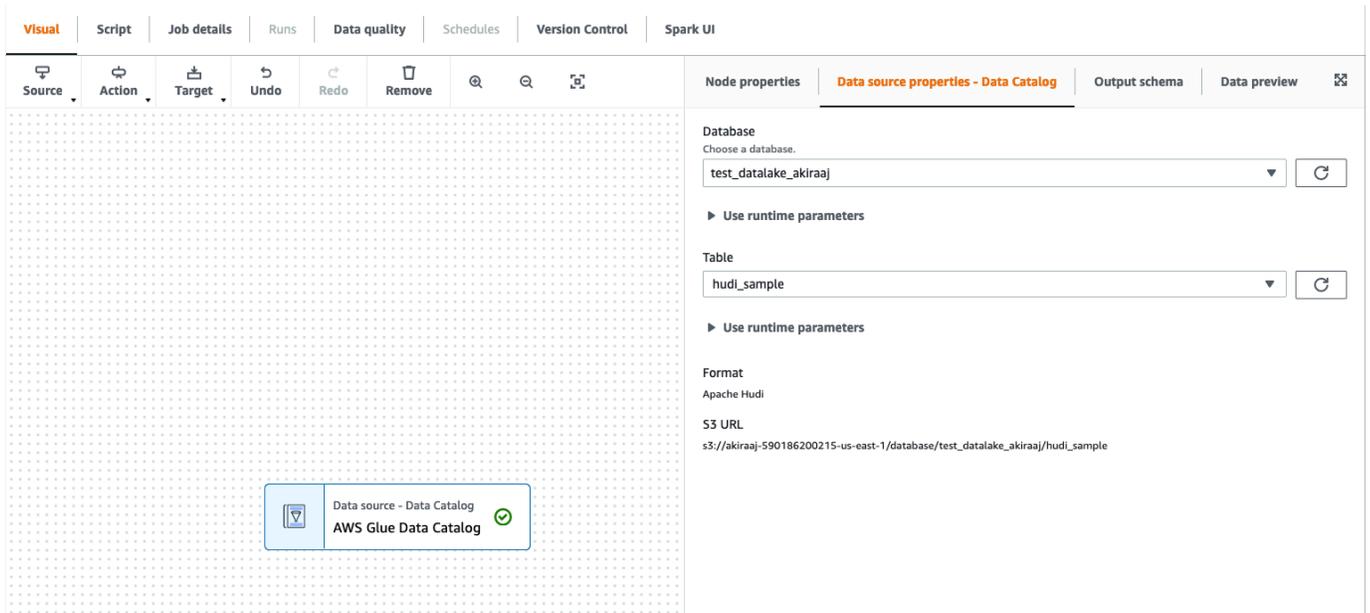
AWS Glue Studio Glue で Hudi フレームワークを使用する

ジョブを作成または編集すると、AWS Glue Studio では、使用している AWS Glue のバージョンに応じて対応する Hudi ライブラリが自動的に追加されます。詳細は、「[AWS Glue での Hudi フレームワークの使用](#)」を参照してください。

データカタログのデータソースで Apache Hudi フレームワークを使用する

Hudi データソースフォーマットをジョブに追加するには:

1. [ソース] メニューで、[AWS Glue Studio データカタログ] を選択します。
2. [Data source properties] (データソースのプロパティ) タブで、データベースとテーブルを選択します。
3. AWS Glue Studio に、フォーマットタイプとして Apache Hudi と Amazon S3 URL が表示されます。



Amazon S3 データソースで Hudi フレームワークを使用する

1. [Source] (ソース) メニューで、[Amazon S3] を選択します。
2. Amazon S3 ソースタイプとしてデータカタログテーブルを選択する場合は、データベースとテーブルを選択します。
3. AWS Glue Studio に、Apache Hudi と Amazon S3 URL がフォーマットとして表示されます。
4. Amazon S3 ソースタイプとして Amazon S3 の場所を選択する場合は、[Browse Amazon S3] (Amazon S3 を参照) をクリックして Amazon S3 URL を選択します。
5. [Data format] (データフォーマット) で、[Apache Hudi] を選択します。

Note

AWS Glue Studio が、選択した Amazon S3 フォルダまたはファイルからスキーマを推測できない場合は、[Additional options] (追加オプション) を選択し、新しいフォルダまたはファイルを選択します。

[Additional options] (追加オプション) の [Schema inference] (スキーマ推論) から、次のオプションを選択します。

- [AWS Glue Studio にサンプルファイルを自動選択させる] — スキーマを推測できるよう、AWS Glue Studio が Amazon S3 の場所にあるサンプルファイルを選択します。自動選択されたファイルは [Auto-sampled file] (自動サンプル化ファイル) フィールドで確認できます。

- [Choose a sample file from Amazon S3] (Amazon S3 からサンプルファイルを選択) - [Browse Amazon S3] (Amazon S3 を参照) をクリックし、使用する使用する Amazon S3 ファイルを選択します。

6. [Infer schema] (スキーマを推測) をクリックします。続いて、[Output schema] (出力スキーマ) タブをクリックすると、出力スキーマを確認できます。
7. [Additional options] (追加オプション) をクリックし、キー値ペアを入力します。

Additional options **Info**

Enter additional key-value pairs for your data source connection.

Key

Value



Add new option

データターゲットで Apache Hudi フレームワークを使用する

データカタログのデータターゲットで Apache Hudi フレームワークを使用する

1. [ターゲット] メニューで [AWS Glue Studio データカタログ] を選択します。
2. [Data source properties] (データソースのプロパティ) タブで、データベースとテーブルを選択します。
3. AWS Glue Studio に、フォーマットタイプとして Apache Hudi と Amazon S3 URL が表示されます。

Amazon S3 データターゲットで Apache Hudi フレームワークを使用する

値を入力するか、使用可能なオプションの中から選択し、Apache Hudi 形式を設定します。Apache Hudi の詳細については、「[Apache Hudi のドキュメント](#)」を参照してください。

Node properties | **Data target properties - S3 2** | Output schema | Data preview 

Format

Apache Hudi 

Hudi Table Name

Hudi Storage Type

Copy on write
Recommended for optimizing read performance

Merge on read
Recommended for minimizing write latency

Hudi Write Operation

Upsert 

Hudi Record Key Fields
Set your primary key(s)

Select a source location to view schema 

Hudi Deduplicate Records by Field Value
Set your field to choose the largest value when inserting records with duplicate key(s)

Select a source location to view schema 

Compression Type

GZIP 

S3 Target Location
Choose an S3 location in the format s3://bucket/prefix/object/ with a trailing slash (/).

"/>

Data Catalog update options [Info](#)

Choose how you want to update the Data Catalog table's schema and partitions. These options will only apply if the Data Catalog table is an S3 backed source.

Do not update the Data Catalog

Create a table in the Data Catalog and on subsequent runs, update the schema and add new partitions

Create a table in the Data Catalog and on subsequent runs, keep existing schema and add new partitions

Partition keys - optional
Add partition keys.

- [Hudi Table Name] (Hudi テーブル名) — お使いの Hudi Table の名前です。
- [Hudi Storage Type] (Hudi ストレージタイプ) — 次の 2 つのうちいずれかを選択します。
 - [Copy on write] (書き込み時にコピー) — 読み取りパフォーマンスの最適化のため、推奨されています。こちらはデフォルトの Hudi ストレージタイプです。書き込み中、更新するたびに新しいバージョンのファイルが作成されます。
 - [Merge on read] (読み取り時に結合) — 書き込みレイテンシーの最小化のため、推奨されています。更新は、行形式の差分ファイルに記録され、必要に応じて圧縮されて、新しいバージョンの列形式のファイルが作成されます。
- [Hudi Write Operation] (Hudi 書き込みオペレーション) - 次の中から選択します。
 - [Upsert] (アップサート) — デフォルトのオペレーションで、入力レコードはインデックスを検索して挿入または更新として最初にタグ付けされます。既存データを更新する場合に推奨されています。
 - [Insert] (インサート) — レコードが挿入されます。ただし、既存のレコードはチェックされないため重複する可能性があります。
 - [Bulk Insert] (一括インサート) — レコードを挿入します。データの量が多い場合に推奨されています。
- [Hudi Record Key Fields] (Hudi レコードキーフィールド) — 検索バーを使ってプライマリレコードキーを検索し選択します。Hudi のレコードは、レコードキーとレコードが属しているパーティションパスのペアから成るプライマリキーで識別されます。
- [Hudi Precombine Field] (Hudi 事前結合フィールド) — 実際には書き込む前の事前結合で使用されるフィールドです。2 つのレコードのキー値が同じである場合、AWS Glue Studio は、事前結合フィールドの値が大きい方のレコードを選択します。増分値 (例: updated_at) が属するフィールドを設定します。
- [Compression Type] (圧縮タイプ) — 圧縮タイプ (非圧縮、GZIP、LZO、Snappy) の中から 1 つを選択します。
- [Amazon S3 Target Location] (Amazon S3 ターゲットの場所) — [Browse S3] (S3 を参照) をクリックして Amazon S3 ターゲットの場所を選択します。
- [Data Catalog update options] (データカタログ更新オプション) — 次の中から選択します。
 - Do not update the Data Catalog(データカタログを更新しない): (デフォルト) スキーマが変更されたり、新しいパーティションが追加された場合、ジョブでデータカタログを更新したくない場合は、このオプションを選択します。
 - [Create a table in the Data Catalog and on subsequent runs, update the schema and add new partitions] (データカタログにテーブルを作成し、それ以降の実行時にスキーマを更新して新しいパーティションを追加する): このオプションを選択すると、最初のジョブの実行時に、ジョブに

よりデータカタログにテーブルが作成されます。それ以降のジョブの実行時にスキーマが変更されたり、新しいパーティションが追加されたりすると、ジョブによりデータカタログテーブルが更新されます。

また、データカタログからデータベースを選択し、テーブル名を入力する必要があります。

- [Create a table in the Data Catalog and on subsequent runs, keep existing schema and add new partitions] (データカタログでテーブル作成し、それ以降の実行時に既存のスキーマの保持して新しいパーティションを追加する): このオプションを選択すると、最初のジョブの実行時に、ジョブによりデータカタログにテーブルが作成されます。それ以降のジョブの実行時に、新しいパーティションを追加するため、ジョブでデータカタログテーブルが更新されます。

また、データカタログからデータベースを選択し、テーブル名を入力する必要があります。

- Partition keys (パーティションキー): 出力でパーティションキーとして使用する列を選択します。さらにパーティションキーを追加するには、[Add a partition key] (パーティションキーの追加) を選択します。
- [Additional options] (追加オプション) — 必要に応じてキーと値のペアを入力します。

AWS Glue Studio によるコード生成

ジョブを保存すると、Hudi のソースまたはターゲットが検出された場合に、次のジョブパラメータがジョブに追加されます。

- `--datalake-formats` — ビジュアルジョブで、([Format] (形式) を選択することで、直接的またはデータレイクでバックアップされたカタログテーブルを選択することで間接的に) 検出されたデータレイク形式の個別リスト。
- `--conf` — `--datalake-formats` の値に基づいて生成される。例えば、`--datalake-formats` の値が「hudi」である場合、AWS Glue はこのパラメータの `spark.serializer=org.apache.spark.serializer.KryoSerializer` `-conf spark.sql.hive.convertMetastoreParquet=false` の値を生成します。

AWS Glue が提供したライブラリのオーバーライド

AWS Glue でサポートされていない Hudi のバージョンを使用するには、独自の Hudi ライブラリ JAR ファイルを指定します。独自の JAR ファイルを使用するには:

- `--extra-jars` ジョブパラメータを使用します。例えば、`'--extra-jars': 's3pathtojarfile.jar'` と指定します。詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。
- `--datalake-formats` ジョブパラメータの値として、`hudi` を含めることはできません。値として空白の文字列を入力すると、AWS Glue はデータレイクライブラリを自動的に提供することができません。詳細は、「[AWS Glue での Hudi フレームワークの使用](#)」を参照してください。

AWS Glue Studio でデータレイクフレームワークを使用する

データソースでデータレイクフレームワークを使用する

Amazon S3 データソースでデータレイクフレームワークを使用する

1. [Source] (ソース) メニューで、[Amazon S3] を選択します。
2. Amazon S3 ソースタイプとしてデータカタログテーブルを選択する場合は、データベースとテーブルを選択します。
3. AWS Glue Studio に、Delta Lake と Amazon S3 URL がフォーマットとして表示されます。
4. [Additional options] (追加オプション) をクリックし、キー値ペアを入力します。例えば、キーと値のペアは次のようになります。キー: `timestampAsOf`、値: `2023-02-24 14:16:18`。

Additional options [Info](#)

Enter additional key-value pairs for your data source connection.

Key

Value



Add new option

5. Amazon S3 ソースタイプとして Amazon S3 の場所を選択する場合は、[Browse Amazon S3] (Amazon S3 を参照) をクリックして Amazon S3 URL を選択します。
6. [Data format] (データ形式) で [Delta Lake] を選択します。

Note

AWS Glue Studio が、選択した Amazon S3 フォルダまたはファイルからスキーマを推測できない場合は、[Additional options] (追加オプション) を選択し、新しいフォルダまたはファイルを選択します。

[Additional options] (追加オプション) の [Schema inference] (スキーマ推論) から、次のオプションを選択します。

- [AWS Glue Studio にサンプルファイルを自動選択させる] — スキーマを推測できるよう、AWS Glue Studio が Amazon S3 の場所にあるサンプルファイルを選択します。自動選択されたファイルは [Auto-sampled file] (自動サンプル化ファイル) フィールドで確認できます。
- [Choose a sample file from Amazon S3] (Amazon S3 からサンプルファイルを選択) - [Browse Amazon S3] (Amazon S3 を参照) をクリックし、使用する使用する Amazon S3 ファイルを選択します。

7. [Infer schema] (スキーマを推測) をクリックします。続いて、[Output schema] (出カスキーマ) タブをクリックすると、出カスキーマを確認できます。

データカタログのデータソースで Delta Lake フレームワークを使用する

1. [ソース] メニューで、[AWS Glue Studio データカタログ] を選択します。
2. [Data source properties] (データソースのプロパティ) タブで、データベースとテーブルを選択します。
3. AWS Glue Studio に、フォーマットタイプとして Delta Lake と Amazon S3 URL が表示されます。

Note

Delta Lake ソースが AWS Glue データカタログテーブルとしてまだ登録されていない場合は、次の 2 つの方法があります。

1. Delta Lake データストアの AWS Glue クローラーを作成します。詳細は、「[Delta Lake データストアの設定オプションを指定する方法](#)」を参照してください。

2. Amazon S3 データソースを使用して Delta Lake のデータソースを選択します。「[Amazon S3 データソースでデータレイクフレームワークを使用する](#)」を参照してください。

データターゲットで Delta Lake フォーマットを使用する

データカタログのデータターゲットで Delta Lake フォーマットを使用する

1. [ターゲット] メニューで [AWS Glue Studio データカタログ] を選択します。
2. [Data source properties] (データソースのプロパティ) タブで、データベースとテーブルを選択します。
3. AWS Glue Studio に、フォーマットタイプとして Delta Lake と Amazon S3 URL が表示されます。

Amazon S3 データソースで Delta Lake フォーマットを使用する

値を入力するか使用可能なオプションの中から選択し、Delta Lake 形式を設定します。

- [Compression Type] (圧縮タイプ) — 圧縮タイプ (非圧縮、Snappy) の中から 1 つ選択します。
- [Amazon S3 Target Location] (Amazon S3 ターゲットの場所) — [Browse S3] (S3 を参照) をクリックして Amazon S3 ターゲットの場所を選択します。
- [Data Catalog update options] (データカタログ更新オプション) — Glue Studio ビジュアルエディターでは、この形式のデータカタログの更新はサポートされていません。
 - Do not update the Data Catalog(データカタログを更新しない): (デフォルト) スキーマが変更されたり、新しいパーティションが追加された場合、ジョブでデータカタログを更新したくない場合は、このオプションを選択します。
 - AWS Glue ジョブの実行後にデータカタログを更新するには、AWS Glue クローラーを実行またはスケジューリングします。詳細は、「[Delta Lake データストアの設定オプションを指定する方法](#)」を参照してください。
- [Partition keys] (パーティションキー): 出力でパーティションキーとして使用する列を選択します。さらにパーティションキーを追加するには、[Add a partition key] (パーティションキーの追加) を選択します。
- [Additional options] (追加オプション) をクリックして、キー値ペアを入力します。例えば、キーと値のペアは次のようになります。キー: timestampAsOf、値: 2023-02-24 14:16:18。

AWS Glue Studio での Apache Iceberg フレームワークの使用

データターゲットで Apache Iceberg フレームワークを使用する

データカタログのデータターゲットで Apache Iceberg フレームワークを使用する

1. [ターゲット] メニューで [AWS Glue Studio データカタログ] を選択します。
2. [Data source properties] (データソースのプロパティ) タブで、データベースとテーブルを選択します。
3. AWS Glue Studio に、フォーマットタイプとして Apache Iceberg と Amazon S3 URL が表示されます。

Amazon S3 データターゲットで Apache Iceberg フレームワークを使用する

値を入力するか、使用可能なオプションの中から選択し、Apache Iceberg 形式を設定します。

- 形式 — ドロップダウンメニューから [Apache Iceberg] を選択します。
- [Amazon S3 ターゲットの場所] — [S3 を参照] をクリックして Amazon S3 ターゲットの場所を選択します。
- [データカタログの更新オプション] — 続けるには、[データカタログでテーブル作成し、それ以降の実行時に既存のスキーマを保持して新しいパーティションを追加する] を選択する必要があります。AWS Glue を使用して新しい Iceberg テーブルを作成するには、Iceberg テーブルのカタログとして Data Catalog を設定する必要があります。Data Catalog に登録されている既存の Iceberg テーブルを更新するには、ターゲットとして Data Catalog を選択します。
- [データベース] — Data Catalog からデータベースを選択します。
- [テーブル名] — テーブルの名前を値として入力します。Apache Iceberg のテーブル名は、すべて小文字でなければなりません。スペースは使用できないため、必要な場合はアンダースコアを使用してください。例えば、「data_lake_format_tables」のようにします。

Node properties	Data target properties - S3	Output schema	Data preview	✕
-----------------	------------------------------------	---------------	--------------	---

Format

Apache Iceberg ▼

Compression Type

GZIP ▼

S3 Target Location

Choose an S3 location in the format `s3://bucket/prefix/object/` with a trailing slash (/).

🔍 ✕

Data Catalog update options

Choose how you want to update the Data Catalog table's schema and partitions. These options will only apply if the Data Catalog table is an S3 backed source.

- Do not update the Data Catalog
- Create a table in the Data Catalog and on subsequent runs, update the schema and add new partitions
- Create a table in the Data Catalog and on subsequent runs, keep existing schema and add new partitions

Database

Choose the database from the AWS Glue Data Catalog.

▼

▶ **Use runtime parameters**

Table name

Enter a table name for the AWS Glue Data Catalog.

Amazon S3 データソースで Apache Iceberg フレームワークを使用する

データカタログのデータソースで Apache Iceberg フレームワークを使用する

1. [ソース] メニューで、[AWS Glue Studio データカタログ] を選択します。
2. [Data source properties] (データソースのプロパティ) タブで、データベースとテーブルを選択します。
3. AWS Glue Studio に、フォーマットタイプとして Apache Iceberg と Amazon S3 URL が表示されます。

Node properties	Data source properties - S3	Output schema	Data preview
S3 source type			
<input type="radio"/> S3 location Choose a file or folder in an S3 bucket.			
<input checked="" type="radio"/> Data Catalog table			
Database Choose a database.			
<input type="text" value="data_lake_format_tables"/> <input type="button" value="↻"/>			
▶ Use runtime parameters			
Table			
<input type="text" value="source_iceberg"/> <input type="button" value="↻"/>			
▶ Use runtime parameters			
Format Apache Iceberg			
S3 URL s3://data-lake-format-data/iceberg/ <input type="button" value="↗"/>			
Partition predicate - optional Enter a boolean expression supported by Spark SQL, using only partition columns.			
<input type="text"/>			
Partition predicate syntax for Spark SQL is <code>year == year(date_sub(current_date, 7)) AND month == month(date_sub(current_date, 7)) AND day == day(date_sub(current_date, 7))</code> .			

Amazon S3 データソースで Apache Iceberg フレームワークを使用する

Apache Iceberg は、AWS Glue Studio の Amazon S3 ソースノードのデータオプションとしては使用できません。

データターゲットノードの設定

データターゲットは、変換されたデータがジョブにより書き込まれる場所です。

データターゲットでのオプションの概要

データターゲット (データシンクとも呼ばれます) は、次のようになります。

- S3 – ジョブにより、選択した Amazon S3 の場所にある、指定した形式のファイルにデータが書き込まれます。

データターゲットにパーティション列を設定すると、ジョブによりパーティションキーに基づいて Amazon S3 へのデータセットがディレクトリに書き込まれます。

- AWS Glue Data Catalog – データカタログ内のテーブルに関連付けられた情報を使用して、ジョブによりターゲットの場所に出カデータが書き込まれます。

テーブルは、手動で、またはクローラーを使用して作成できます。また、AWS CloudFormation テンプレートを使用して、データカタログにテーブルを作成できます。

- コネクタ – コネクタは、データストアと AWS Glue 間の通信を容易にするコードの一部です。コネクタや関連付けられた接続を使用して、ジョブにより出カデータがターゲットの場所へ書き込まれます。提供されるコネクタをサブスクライブするには、AWS Marketplace または、独自のカスタムコネクタを作成することもできます。詳細については、「[AWS Glue Studio にコネクタを追加する](#)」を参照してください。

ジョブにより Amazon S3 のデータターゲットに書き込みが行われる際、データカタログを更新するように選択できます。このオプションを使用すると、スキーマまたはパーティションの変更時に、クローラーをリクエストしなくてもデータカタログが更新されます。これにより、テーブルを最新の状態に保つことが容易になります。また、このオプションではデータカタログに新しいテーブルを追加したり、テーブルパーティションを更新したり、ジョブから直接テーブルのスキーマを更新できます。これにより、データを分析に使用できるようにするプロセスが簡略化されます。

データターゲットノードの編集

データターゲットは、変換されたデータがジョブにより書き込まれる場所です。

ジョブ図でデータターゲットノードを追加または設定するには

1. (オプション) ターゲットノードを追加する必要がある場合、ビジュアルエディタの上部のツールバーで [Target] (ターゲット) を選択し、その後 [S3] または [Glue Data Catalog] (Glue データカタログ) を選択します。
 - ターゲットに [S3] を選択した場合、ジョブにより指定した Amazon S3 の場所にある 1 つ以上のファイルにデータセットが書き込まれます。
 - ターゲットに [AWS Glue Data Catalog] を選択した場合、ジョブによりデータカタログから選択したテーブルによって記述された場所へ書き込みが行われます。

2. ジョブ図でデータターゲットノードを選択します。ノードを選択すると、ページの右側にノードの詳細パネルが表示されます。
3. [Node properties] (ノードのプロパティ) タブを選択して、次の情報を入力します。
 - Name (名前): ジョブ図でノードに関連付ける名前を入力します。
 - Node type (ノードタイプ): 値は既に選択されていますが、必要に応じて変更できます。
 - Node parents (ノードの親): 親ノードは、ターゲットの場所に書き込む出力データを提供するジョブ図内のノードです。事前に設定されたジョブ図では、ターゲットノードで既に親ノードが選択されている必要があります。親ノードが表示されていない場合は、リストから親ノードを選択します。

ターゲットノードは、単一の親ノードを持ちます。
4. データターゲットのプロパティ情報を設定する 詳細については、次のセクションを参照してください。
 - [データターゲットに Amazon S3 を使用する](#)
 - [データターゲットにデータカタログテーブルを使用する](#)
 - [データターゲットにコネクタを使用する](#)
5. (オプション) データターゲットノードのプロパティの設定後、ノードの詳細パネルの [Output schema] (出力スキーマ) タブを選択して、データの出力スキーマを表示できます。ジョブ内の任意のノードに対してこのタブを初めて選択すると、データにアクセスする IAM ロールを指定するよう求められます。[Job details] (ジョブの詳細) タブで IAM ロールをまだ指定していない場合、ここで IAM ロールを入力するよう求められます。

データターゲットに Amazon S3 を使用する

Amazon S3 とコネクタを除くすべてのデータソースでは、選択するソースタイプの AWS Glue Data Catalog にテーブルが存在する必要があります。AWS Glue Studio はデータカタログテーブルを生成しません。

Amazon S3 に書き込むデータターゲットノードを設定するには

1. 新規または保存済みのジョブのビジュアルエディタに移動します。
2. ジョブ図でデータソースノードを選択します。
3. [Data source properties] (データソースのプロパティ) タブを選択して、次の情報を入力します。

- **Format (形式):** リストから形式を選択します。データ結果に使用できる形式のタイプは次のとおりです。
 - JSON: JavaScript Object Notation。
 - CSV: カンマで区切られた値。
 - Avro: Apache Avro JSON バイナリ。
 - Parquet: Apache Parquet 列指向ストレージ。
 - Glue Parquet: データ形式として DynamicFrames 用に最適化されたカスタム Parquet ライターのタイプ。データに対して事前に計算されたスキーマをリクエストしなくても、スキーマを動的に計算して変更します。
 - ORC: Apache Optimized Row Columnar (ORC) 形式。

これらの形式のオプションについての詳細は、AWS Glue デベロッパーガイドの「[Format Options for ETL Inputs and Outputs in AWS Glue](#)」を参照してください。

- **Compression Type (圧縮タイプ):** gzip または bzip2 の形式を使用して、任意でデータを圧縮するかどうかを選択できます。デフォルトでは圧縮なし、または [None] (なし) です。
- **S3 Target Location (S3 ターゲットの場所):** Amazon S3 バケットおよびデータを出力する場所。[Browse S3] (S3 をブラウズ) ボタンをクリックして、アクセスできる Amazon S3 バケットを表示し、バケットの 1 つをターゲットの宛先として選択できます。
- **データカタログの更新オプション**
 - **Do not update the Data Catalog(データカタログを更新しない):** (デフォルト) スキーマが変更されたり、新しいパーティションが追加された場合、ジョブでデータカタログを更新したくない場合は、このオプションを選択します。
 - **Create a table in the Data Catalog and on subsequent runs, update the schema and add new partitions (データカタログにテーブルを作成し、それ以降の実行時にスキーマを更新して新しいパーティションを追加する):** このオプションを選択すると、最初のジョブの実行時に、ジョブによりデータカタログにテーブルが作成されます。それ以降のジョブの実行時にスキーマが変更されたり、新しいパーティションが追加されたりすると、ジョブによりデータカタログテーブルが更新されます。

また、データカタログからデータベースを選択し、テーブル名を入力する必要があります。

- **[Create a table in the Data Catalog and on subsequent runs, keep existing schema and add new partitions] (データカタログでテーブル作成し、それ以降の実行時に既存のスキーマの保持して新しいパーティションを追加する):** このオプションを選択すると、最初のジョブの実行時に、ジョブによりデータカタログにテーブルが作成されます。それ以降のジョブの実

行時に、新しいパーティションを追加するため、ジョブでデータカタログテーブルが更新されます。

また、データカタログからデータベースを選択し、テーブル名を入力する必要があります。

- Partition keys (パーティションキー): 出力でパーティションキーとして使用する列を選択します。さらにパーティションキーを追加するには、[Add a partition key] (パーティションキーの追加) を選択します。

データターゲットにデータカタログテーブルを使用する

Amazon S3 とコネクタを除くすべてのデータソースでは、選択するターゲットタイプの AWS Glue Data Catalog にテーブルが存在する必要があります。AWS Glue Studio はデータカタログテーブルを生成しません。

データカタログテーブルを使用するターゲットのデータプロパティを設定するには

1. 新規または保存済みのジョブのビジュアルエディタに移動します。
2. ジョブ図でデータターゲットノードを選択します。
3. [Data target properties] (データターゲットのプロパティ) タブを選択して、次の情報を入力します。
 - Database (データベース): ターゲットとして使用するテーブルを含むデータベースをリストから選択します。このデータベースは、データカタログに既に存在している必要があります。
 - Table (テーブル): 出力データのスキーマを定義するテーブルをリストから選択します。このテーブルは、データカタログに既に存在している必要があります。

データカタログのテーブルは、列名、データ型の定義、パーティション情報、およびターゲットのデータセットに関するその他のメタデータで構成されています。ジョブにより、データカタログのこのテーブルで示されている場所に行き込みが行われます。

データカタログでのテーブルの作成についての詳細は、AWS Glue デベロッパーガイドの「[Defining Tables in the Data Catalog](#)」を参照してください。

- データカタログの更新オプション
 - Do not change table definition (テーブルの定義を変更しない): (デフォルト) スキーマが変更されたり、新しいパーティションが追加された場合、ジョブでデータカタログを更新したくない場合は、このオプションを選択します。

- [Update schema and add new partitions] (スキーマを更新して新しいパーティションを追加する): このオプションを選択すると、スキーマが変更されたり新しいパーティションが追加される際に、ジョブによりデータカタログテーブルが更新されます。
- [Keep existing schema and add new partitions] (既存のスキーマを保持して新しいパーティションを追加する): このオプションを選択すると、ジョブにより データカタログテーブルが更新され、新しいパーティションが追加されます。
- Partition keys (パーティションキー): 出力でパーティションキーとして使用する列を選択します。さらにパーティションキーを追加するには、[Add a partition key] (パーティションキーの追加) を選択します。

データターゲットにコネクタを使用する

ノードタイプにコネクタを選択する場合、[カスタムコネクタを使用したジョブのオーサリング](#) の説明に従って、データターゲットのプロパティ設定を完了します。

ジョブスクリプトの編集またはアップロード

AWS Glue Studio のビジュアルエディタを使用して、ジョブスクリプトを編集したり、独自のスクリプトをアップロードします。

ジョブが AWS Glue Studio で作成された場合にのみ、ビジュアルエディタを使用してジョブノードを編集できます。AWS Glue コンソールをAPI コマンドを通して、またはコマンドラインインターフェイス (CLI) を使用してジョブが作成された場合、AWS Glue Studio 中のスクリプトエディタを使用して、ジョブスクリプト、パラメータ、スケジュールを編集できます。ジョブをスクリプト専用モードに変換して、AWS Glue Studio 内に作成されたジョブのスクリプトを編集することもできます。

ジョブスクリプトを編集または独自のスクリプトをアップロードするには

1. 新しいジョブを作成する場合、[Jobs] (ジョブ) ページで、[Spark script editor] (Spark スクリプトエディタ) オプションを選択して Spark ジョブを作成するか、[Python Shell script editor] (Python シェルスクリプトエディタ) を選択して Python シェルジョブを作成します。新しいスクリプトを作成するか、既存のスクリプトをアップロードできます。[Spark script editor] (Spark スクリプトエディタ) を選択した場合、Scala スクリプトまたは Python スクリプトのいずれかを作成またはアップロードできます。[Python Shell script editor] (Python シェルスクリプトエディタ) を選択した場合、Python スクリプトの作成またはアップロードのみが可能です。

新しいジョブを作成するオプションを選択した後、表示される [Options] (オプション) セクションで、スタータースクリプト (定型コードを含む新しいスクリプトを作成) から開始するか、ジョブスクリプトとして使用するローカルファイルをアップロードするかを選択できます。

[Spark script editor] (Spark スクリプトエディタ) を選択した場合、Python または Scala スクリプトファイルのどちらかをアップロードできます。Scala スクリプトには、ファイル拡張子 `.scala` が必要です。Python スクリプトは、Python タイプのファイルとして認識される必要があります。[Python Shell script editor] (Python シェルスクリプトエディタ) を選択した場合、Python スクリプトファイルのみをアップロードできます。

選択が完了したら、[Create] (作成) をクリックしてジョブを作成し、ビジュアルエディタを開きます。

2. 新規または保存済みのジョブのビジュアルジョブエディタに移動して、[Script] (スクリプト) タブを選択します。
3. スクリプトエディタのオプションを使用して新しいジョブを作成しておらず、既存のジョブのスクリプトを編集したことがない場合、[Script] (スクリプト) タブには、[Script (Locked)] (スクリプト (ロックされています)) という見出しが表示されます。これは、スクリプトエディタが読み取り専用モードであることを意味します。[Edit script] (スクリプトの編集) を選択して、編集用のスクリプトのロックを解除します。

スクリプトを編集可能にするには、AWS Glue Studio がジョブをビジュアルジョブからスクリプト専用ジョブに変換します。編集用にスクリプトをロック解除すると、保存後にこのジョブでビジュアルエディタを使用できなくなります。

確認ウィンドウで、[Confirm] (確認) を選択して続行するか、[Cancel] (キャンセル) を選択して引き続きビジュアルの編集にジョブを使用します。

[Confirm] (確認) を選択すると、[Visual] (ビジュアル) タブはエディタに表示されなくなります。AWS Glue Studio を使って、スクリプトエディタを使用したスクリプトを変更したり、ジョブの詳細やスケジュールを変更したり、実行中のジョブを表示したりできます。

Note

ジョブを保存するまで、スクリプト専用ジョブへの変換は永続的ではありません。コンソールのウェブページを更新したり、ビジュアルエディタでジョブを保存する前に閉じて再度開いても、ビジュアルエディタで個々のノードを編集できます。

4. 必要に応じて、スクリプトを編集します。

スクリプトの編集が終了したら、[Save] (保存) をクリックしてジョブを保存し、ジョブをビジュアルからスクリプト専用に変換します。

- (オプション)[Script] タブ 上の[Download]ボタンを選択することにより、AWS Glue Studio コンソールからのスクリプトをダウンロードできます。このボタンをクリックすると、新しいブラウザウィンドウが開き、Amazon S3 内にあるその場所からのスクリプトが表示されます。[Job details] (ジョブの詳細) タブの [Script filename] (スクリプトファイル名) および [Script path] (スクリプトパス) パラメータにより、Amazon S3 のスクリプトファイルの名前と場所が決定します。

Join test job2

The screenshot shows the 'Job details' tab in AWS Glue Studio. At the top, there are navigation tabs: 'Visual', 'Script', 'Job details' (selected), 'Runs', and 'Schedules'. Below the tabs, there is a section titled 'Advanced properties'. Under this section, there are two main fields: 'Script filename' and 'Script path'. The 'Script filename' field contains the text 'Join test job.py'. The 'Script path' field contains the text 's3://aws-glue-assets-111122223333-t'. Below the 'Script path' field, there are three buttons: a search button with a magnifying glass icon, a 'View' button with an external link icon, and a 'Browse S3' button. Below these fields, there are three checkboxes with labels and 'Info' links: 'Job metrics' (unchecked), 'Continuous logging' (checked), and 'Spark UI' (checked). Each checkbox has a brief description of its function.

ジョブを保存すると、AWS Glue によりこれらのフィールドで指定された場所にジョブスクリプトが保存されます。Amazon S3 内のこの場所にあるスクリプトファイルを変更すると、AWS Glue Studio は次回ジョブを編集するときに、変更されたスクリプトをロードします。

AWS Glue Studio 中の Scala スクリプトの作成および編集

ジョブの作成のためにスクリプトエディタを選択する場合、デフォルトでは、ジョブのプログラミング言語は Python 3 に設定されています。スクリプトをアップロードするのではなく、新しいスク

リプトの記述を選択した場合、AWS Glue Studio は、Python 内に記述された定型テキストを用いて新しいスクリプトを開始します。Scala スクリプトを記述する場合は、Scala を使用するよう最初にスクリプトエディタを設定する必要があります。

Note

ジョブのプログラミング言語として Scala を選択し、ビジュアルエディタを使用してジョブを設計する場合、生成されるジョブスクリプトは Scala で記述されます。それ以上のアクションは必要ありません。

AWS Glue Studio 中の新しい Scala スクリプトの作成

1. [Spark script editor] (Spark スクリプトエディタ) オプションを選択して、新しいジョブを作成します。
2. [Options] (オプション) の、[Create a new script with boilerplate code] (定型コードで新しいスクリプトを作成する) を選択します。
3. [Job details] (ジョブの詳細) タブを選択し、[Language] (言語) を [Scala] (Python 3 ではなく) に設定します。

Note

[Spark script editor] (Spark スクリプトエディタ) オプションを選択してジョブを作成すると、ジョブの [Type] (タイプ) プロパティは Spark に自動的に設定されます。

4. [Script] (スクリプト) タブを選択します。
5. Python の定型テキストを削除します。これを、次の Scala の定型テキストに置き換えることができます。

```
import com.amazonaws.services.glue.{DynamicRecord, GlueContext}
import org.apache.spark.SparkContext
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job

object MyScript {
  def main(args: Array[String]): Unit = {
    val sc: SparkContext = new SparkContext()
    val glueContext: GlueContext = new GlueContext(sc)
```

```
}  
}
```

6. エディタで Scala ジョブスクリプトを記述します。必要に応じて、他の import ステートメントを追加します。

AWS Glue Studio 中の Python シェルジョブの作成および編集

ジョブの作成に Python シェルスクリプトエディタを選択すると、既存の Python スクリプトをアップロードするか、新しい Python スクリプトを記述できます。新しいスクリプトを作成する場合、定型コードが新しい Python ジョブスクリプトに追加されます。

新しい Python シェルジョブを作成するには

「[AWS Glue Studio でのジョブの開始](#)」の手順を参照してください。

Python シェルジョブでサポートされているジョブプロパティは、Spark ジョブでサポートされているものとは異なります。次のリストは、[Job details] (ジョブの詳細) タブでの、Python シェルジョブで使用可能なジョブパラメータの変更について示しています。

- ジョブの [Type] (タイプ) プロパティは自動的に [Python Shell] に設定され、変更はできません。
- [Language] (言語) の代わりに、ジョブ用の [Python version] (Python バージョン) プロパティがあります。現在、AWS Glue Studio で作成される Python シェルジョブは Python 3.6 を使用します。
- [Glue version] (Glue バージョン) プロパティは Python シェルジョブには適用されないため、使用できません。
- [Worker type] (ワーカータイプ) および [Number of workers] (ワーカー数) の代わりに、[Data processing units] (データ処理ユニット) プロパティが表示されます。このジョブプロパティにより、ジョブの実行時に Python シェルによって消費されるデータ処理ユニット (DPU) の数が決定します。
- [Job bookmark] (ジョブのブックマーク) プロパティは Python シェルジョブではサポートされていないため、使用できません。
- [Advanced properties] (詳細プロパティ) の次のプロパティは Python シェルジョブでは使用できません。
 - ジョブのメトリクス

- 連続ログ記録
- [Spark UI] および [Spark UI logs path] (Spark UI ログパス)
- 見出し [Libraries] (ライブラリ) の下の [依存 JARS パス]

ジョブ図でノードの親ノードを変更する

ノードの親を変更して、ジョブ図内のノードを移動したり、ノードのデータソースを変更できます。

親ノードを変更するには

1. ジョブ図で変更するノードを選択します。
2. ノードの詳細パネルの [Node properties] (ノードのプロパティ) タブの、見出し [Node parents] (ノードの親) の下で、そのノードの現在の親を削除します。
3. リストから新しい親ノードを選択します。
4. 必要に応じて、新しく選択した親ノードと一致するように、ノードのその他のプロパティを変更します。

誤ってノードを変更した場合は、[Undo] (元に戻す) ボタンをクリックして、アクションを元に戻します。

ジョブ図からノードを削除する

Visual ETL ジョブを使用するとき、削除されたノードに接続されているノードを再追加または再構築せずに、キャンバスからノードを削除できます。

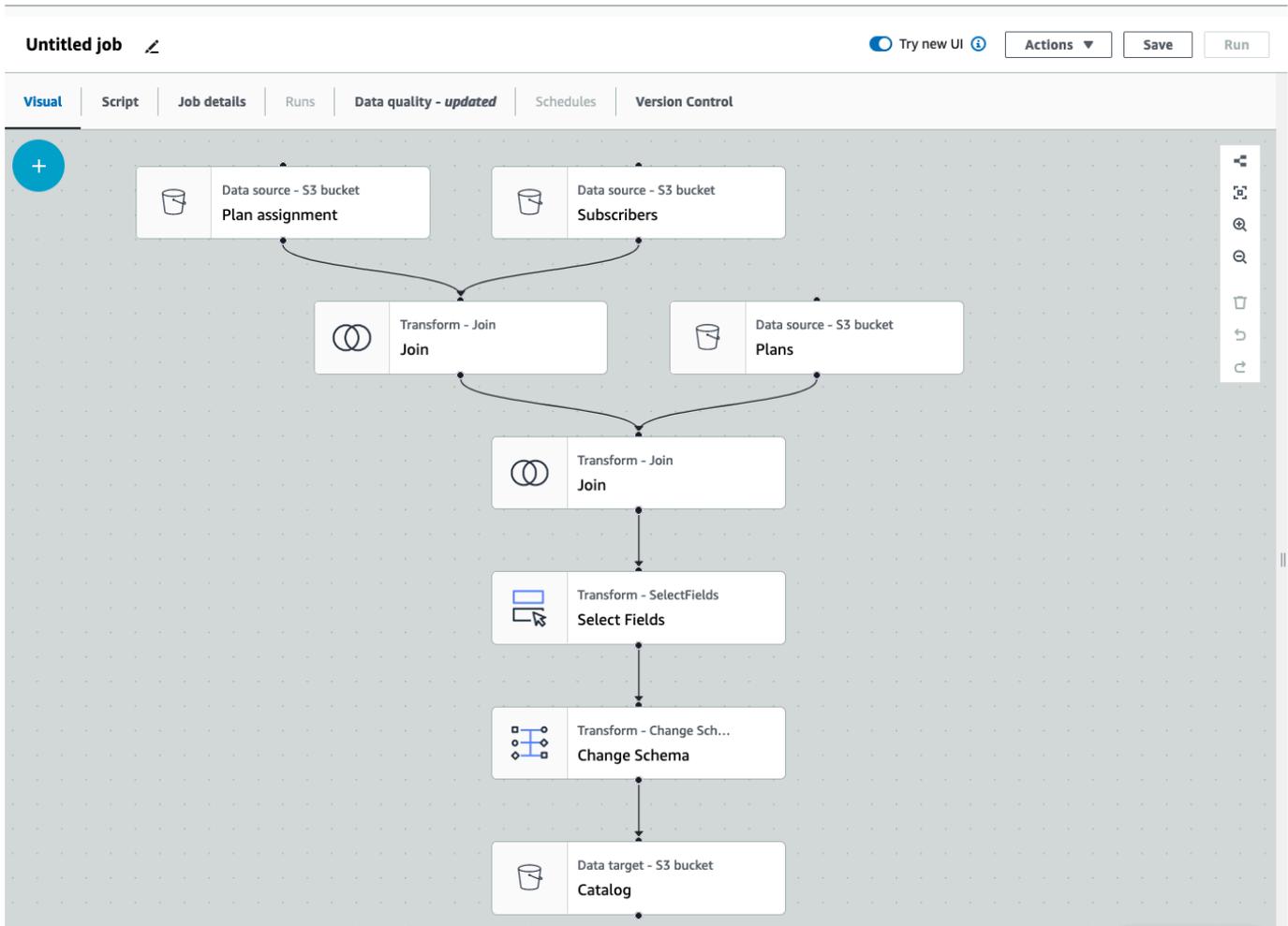
以下の例では、[ETL ジョブ > Visual ETL] を選択し、次に [サンプルジョブ] で [複数のソースを統合する Visual ETL ジョブ] を選択して手順に従います。[サンプルジョブの作成] を選択してジョブを作成し、以下の手順に従います。

The screenshot shows the AWS Glue Studio interface. On the left is a navigation menu with categories like 'Getting started', 'Data Catalog', and 'Data Integration and ETL'. Under 'Data Integration and ETL', 'Visual ETL' is highlighted with a red box. The main content area is titled 'AWS Glue Studio' and includes a 'Create job' section with three options: 'Visual ETL' (highlighted with a red box), 'Notebook', and 'Script editor'. Below this is an 'Example jobs' section with three job templates: 'Visual ETL job to join multiple sources' (highlighted with a red box), 'Ray notebook for parallelizing Python', and 'Spark notebook using Pandas'. At the bottom, there is a 'Your jobs (1)' table showing a single job entry.

Job name	Type	Last modified	AWS Glue version
job_101521	Glue ETL	1/31/2022, 11:44:06 AM	2.0

キャンバスからノードを削除する方法

1. AWS Glue コンソールで、ナビゲーションメニューから [Visual ETL] を選択し、既存のジョブを選択します。ジョブキャンバスは、以下のように示されるサンプルジョブを表示します。



- 削除するノードを選択します。キャンバスがノードにズームインします。キャンバスの右側のツールバーで [ゴミ箱] アイコンを選択します。これによってノードが削除され、そのノードに接続されているすべてのノードが移動し、ワークフローで削除されたノードを置き換えます。この例では、最初の [結合] ノードがキャンバスから削除されました。

ワークフローでノードを削除した場合、ワークフローが無効にならないように AWS Glue がノードを再配置します。それでもノードの設定を修正する必要がある場合があります。

この例では、[サブスクライバー] ノードの下にある [結合] ノードは削除されています。その結果、[プラン] ソースノードは最上位に移動され、子の [結合] ノードにまだ接続されています。[結合] は選択したテーブルを持つ 2 つの親ソースのノードが必要なため、[結合] ノードには追加の設定が必要になりました。キャンバスの右側の [変換] タブは、[結合条件] に不足している要件を表示しています。

The screenshot displays the AWS Glue console interface for an 'Untitled job'. The workflow consists of several nodes: 'Data source - S3 bucket Plan assignment', 'Data source - S3 bucket Subscribers', 'Data source - S3 bucket Plans', 'Transform - Join Join', 'Transform - SelectFields Select Fields', and 'Transform - Change Sch... Change Schema'. The 'Join' node is highlighted with a red box, and the 'Select Fields' node is also highlighted with a red box. The 'Join' node configuration panel on the right shows the 'Join type' set to 'Inner join' and the 'Join conditions' section highlighted with a red box. A warning message at the bottom states 'Node is misconfigured' for the 'Join' node.

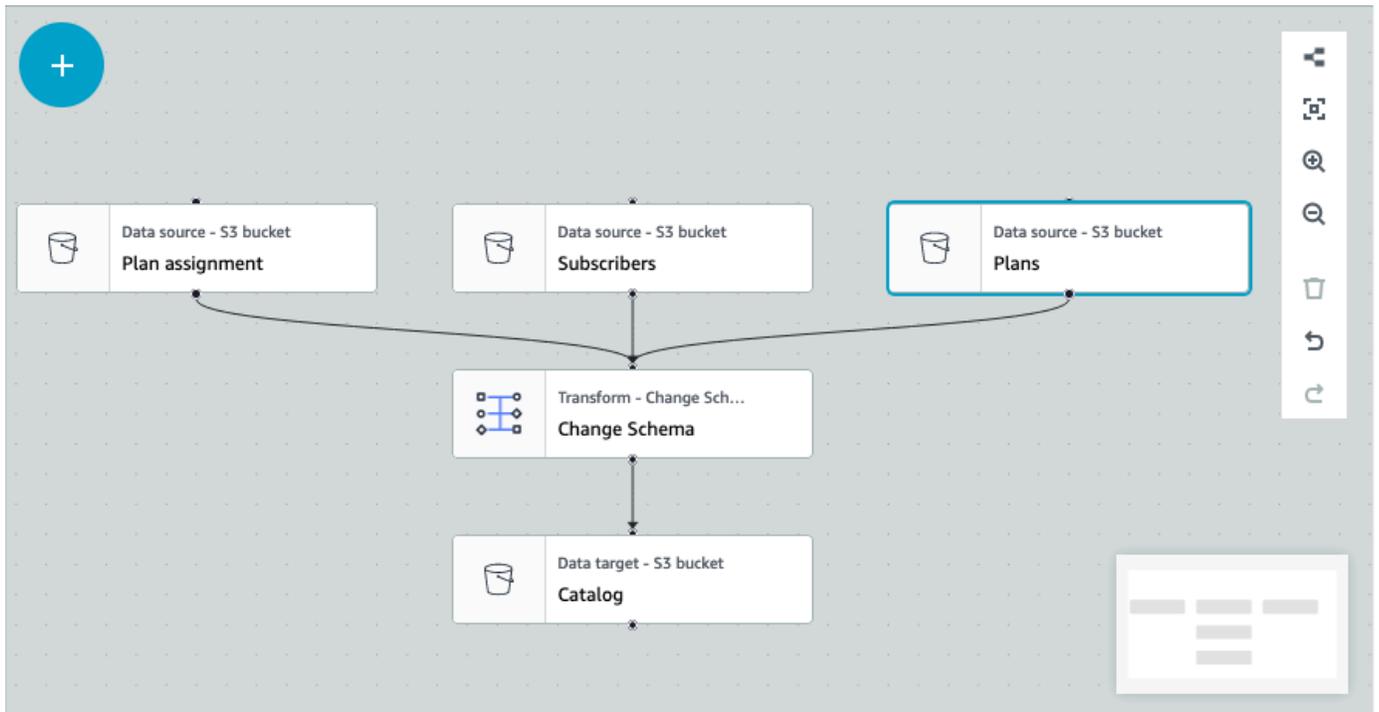
Node is misconfigured
Data preview will be displayed when following node is correctly configured:

- Join

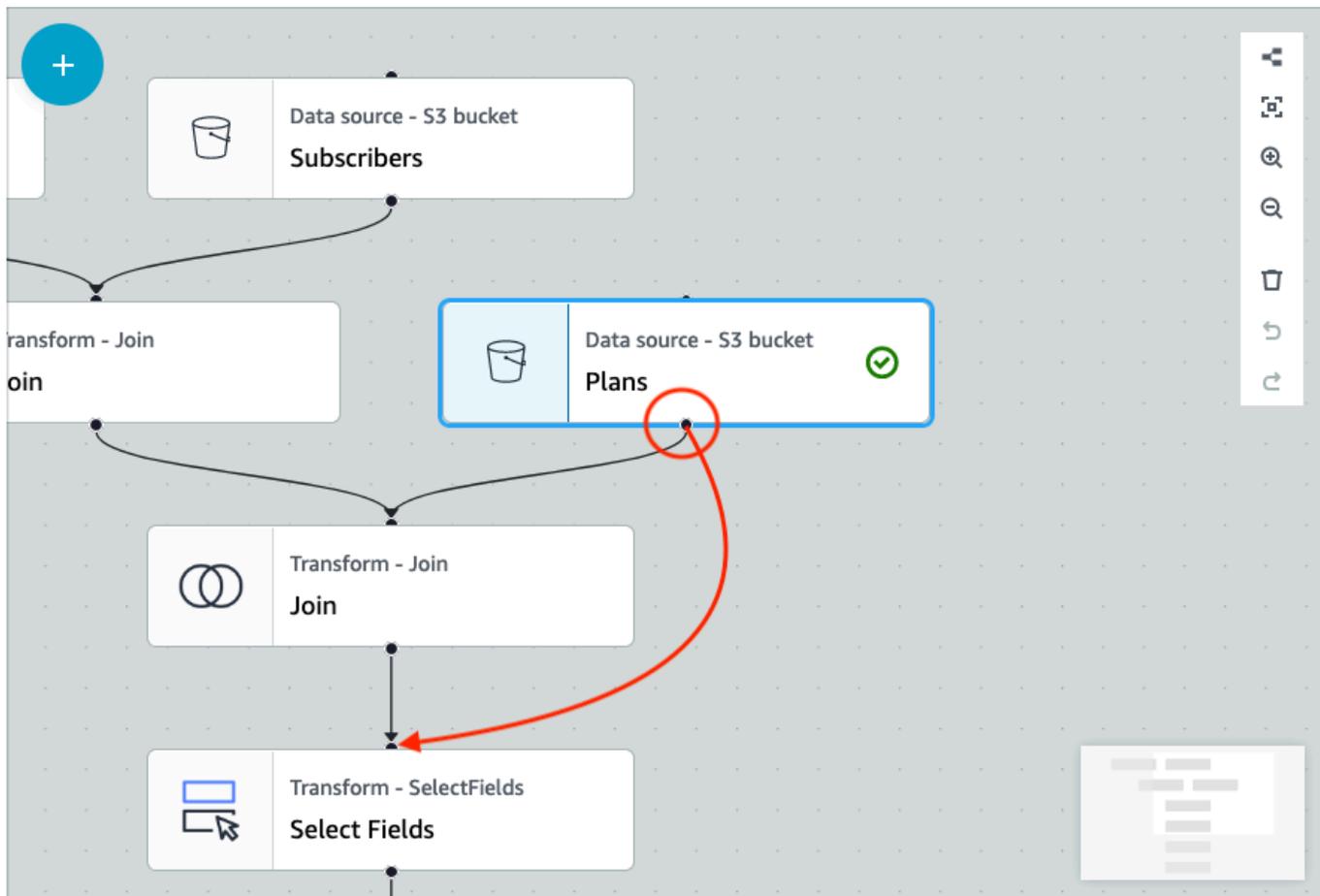
Join conditions
Select a field from each parent node for the join condition.

Insufficient source nodes
The Join transform requires two parent source nodes with selected tables.

- 2つ目の [結合] ノードおよび [フィールド選択] ノードを削除します。ノードが削除されると、ワークフローは以下の例のようになります。



4. ノード接続を変更するには、ノードのハンドルをクリックして接続を新しいノードにドラッグします。ノードを削除したり、論理フローでノードを再配置したりできるようにします。この例では、[プラン] ノードのハンドルをクリックし、赤い矢印で示されているように [結合] ノードに接続をドラッグすることにより、新しい接続を確立しています。



5. 操作を取り消す必要がある場合、キャンバスの右側のツールバーにある [ゴミ箱] アイコンの真下にある [元に戻す] アイコンを選択します。

AWS Glue データカタログノードにソースパラメータとターゲットパラメータを追加する

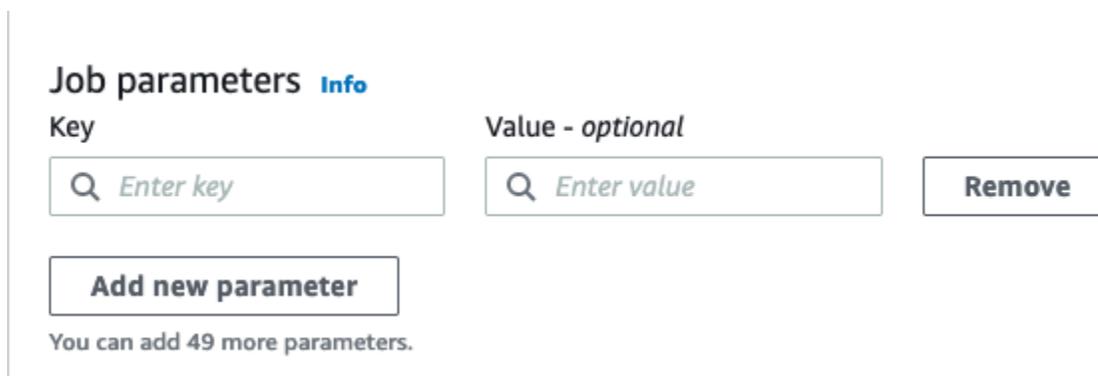
AWS Glue Studio によって、ビジュアルジョブをパラメータ化できます。本番環境と開発環境ではカタログテーブル名が異なる場合があるため、ジョブの実行時に実行されるデータベースとテーブルのランタイムパラメータを定義して選択できます。

ジョブのパラメータ化により、ソースとターゲットをパラメータ化して、AWS Glue データカタログノードを使用する際にそれらのパラメータをジョブに保存できます。ソースとターゲットをパラメータとして指定すると、特に同じジョブを複数の環境で使用する場合に、ジョブの再利用が可能になります。ソースとターゲットを管理する時間と労力を節約して、デプロイメント環境全体でコードをプロモートする際に役立ちます。さらに、指定したカスタムパラメータは、特定の実行のデフォルト引数をオーバーライドします。AWS Glue仕事。

ソースパラメータとターゲットパラメータを追加するには

AWS Glue データカタログノードをソースまたはターゲットとして使用しているかどうかにかかわらず、[Job details] (ジョブの詳細) タブの [Advanced properties] (詳細プロパティ) セクションでランタイムパラメータを定義できます。

1. ソースノードまたはターゲットノードとして AWS Glue データカタログノードを選択します。
2. [Job details] (ジョブの詳細) タブを選択します。
3. [Advanced properties] (詳細プロパティ) を選択します。
4. [Job parameters] (ジョブパラメータ) セクションで、キー値を入力します。例えば、`--db.source` はデータベースソースのパラメータになります。キー名の後に `'` (ダッシュ) が続いている限り、任意の名前をキーに入力できます。



The screenshot shows the 'Job parameters' section in the AWS Glue console. It features a table with two columns: 'Key' and 'Value - optional'. The 'Key' column has an input field with a magnifying glass icon and the placeholder text 'Enter key'. The 'Value - optional' column has a similar input field with the placeholder text 'Enter value'. To the right of the 'Value - optional' input is a 'Remove' button. Below the table is an 'Add new parameter' button. At the bottom of the section, there is a note: 'You can add 49 more parameters.'

5. 値を入力します。例えば、`databasename` はパラメータ化されるデータベースの値になります。
6. さらにパラメータを追加する場合は、[Add new parameter] (新しいパラメータを追加) を選択します。最大 50 個のパラメータを使用できます。キーと値のペアを定義したら、AWS Glue データカタログノードでパラメータを使用できます。

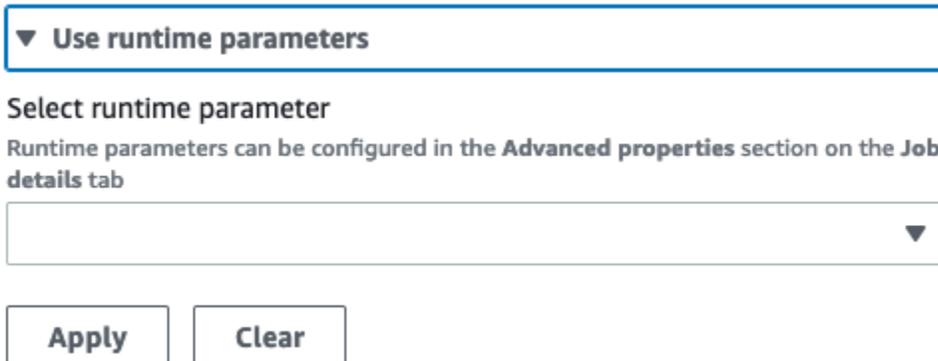
ランタイムパラメータを選択するには

Note

データベースとテーブルのランタイムパラメータを選択するプロセスは、AWS Glue データカタログノードがソースであっても、ターゲットであっても同じです。

1. ソースノードまたはターゲットノードとして AWS Glue データカタログノードを選択します。

2. [Data source properties - Data Catalog] (データソースのプロパティ - データカタログ) タブの [Database] (データベース) で、[Use runtime parameters] (ランタイムパラメータを使用) を選択します。



▼ Use runtime parameters

Select runtime parameter

Runtime parameters can be configured in the **Advanced properties** section on the **Job details** tab

Apply Clear

3. ドロップダウンメニューからパラメータを選択します。例えば、ソースデータベースに定義したパラメータを選択し、[Apply] (適用) を選択すると、データベースのドロップダウンメニューにデータベースが自動的に入力されます。
4. [Table] (テーブル) セクションで、ソーステーブルとして定義済みのパラメータを選択します。[Apply] (適用) を選択した場合は、使用するテーブルとしてテーブルが自動的に入力されます。
5. ジョブを保存して実行すると、AWS Glue Studio はジョブの実行中に選択されたパラメータを参照します。

AWS Glue で Git のバージョン管理システムを使用する

Note

ノートブックは現在、AWS Glue Studio のバージョン管理には対応していません。ただし、AWS Glue ジョブスクリプトとビジュアル ETL ジョブのバージョン管理はサポートされています。

リモートリポジトリがあり、そのリポジトリを使用して AWS Glue ジョブを管理する場合は、AWS Glue Studio または AWS CLI を使用して、リポジトリと AWS Glue のジョブへの変更を同期できます。この方法で変更を同期すると、ジョブは AWS Glue Studio からリポジトリにプッシュされるか、リポジトリから AWS Glue Studio にプルされます。

AWS Glue Studio での Git との統合により、次のことができるようになります。

- AWS CodeCommit、GitHub、GitLab、Bitbucket など、Git のバージョン管理システムとの統合
- ビジュアルジョブやスクリプトジョブを使用するかどうかにかかわらず、AWS Glue Studio で AWS Glue ジョブを編集し、リポジトリに同期する
- ジョブのソースとターゲットをパラメータ化する
- リポジトリからジョブをプルし、AWS Glue Studio で編集する
- AWS Glue Studio のマルチブランチワークフローを利用して、ブランチからプルしたり、ブランチにプッシュしたりして、ジョブをテストする
- リポジトリからファイルをダウンロードし、ジョブを AWS Glue Studio にアップロードしてクロスアカウントジョブを作成する
- 選択した自動化ツール (Jenkins、AWS CodeDeploy など) を使用する

このビデオでは、AWS Glue を Git と統合し、継続的な共有コードパイプラインを構築する方法について紹介します。

IAM アクセス許可

ジョブに次の IAM アクセス許可のいずれかがあることを確認します。IAM アクセス許可の設定方法の詳細については、「[AWS Glue Studio に対する IAM のアクセス許可の設定](#)」を参照してください。

- AWSGlueServiceRole
- AWSGlueConsoleFullAccess

Git との統合では、少なくとも次のアクションが必要です。

- glue:UpdateJobFromSourceControl - バージョン管理システムに存在するジョブで AWS Glue を更新できるようにする
- glue:UpdateSourceControlFromJob - AWS Glue に保存されているジョブでバージョン管理システムを更新できるようにする
- s3:GetObject - バージョン管理システムにプッシュしている間にジョブのスクリプトを取得できるようにする
- s3:PutObject - ソース管理システムからジョブをプルする際にスクリプトを更新できるようにする

前提条件

ジョブをソース管理リポジトリにプッシュするには、以下が必要です。

- 管理者によってすでに作成されているリポジトリ
- リポジトリ内のブランチ
- 個人アクセストークン (Bitbucket の場合、これはリポジトリアクセストークン)
- リポジトリの所有者のユーザー名
- AWS Glue Studio にリポジトリの読み取りと書き込みを許可するためのアクセス許可をリポジトリに設定
 - [GitLab] — トークンスコープを `api`、`read_repository`、`write_repository` に設定
 - [Bitbucket] — アクセス許可を次のように設定します。
 - [Workspace membership] — 読み取り、書き込み
 - [プロジェクト] — 書き込み、管理者による読み取り
 - [リポジトリ] — 読み取り、書き込み、管理、削除

Note

AWS CodeCommit の使用時には、個人アクセストークンとリポジトリの所有者は必要ありません。「[Git および AWS CodeCommit の開始方法](#)」を参照してください。

AWS Glue Studio のソースコントロールリポジトリのジョブを使用する

AWS Glue Studio 以外にあるソース管理リポジトリからジョブを取得し、そのジョブを AWS Glue Studio で使用するための前提条件は、ジョブのタイプによって異なります。

ビジュアルジョブの場合:

- ジョブ名と一致するジョブ定義のフォルダと JSON ファイルが必要です

例えば、以下のジョブ定義を参照してください。リポジトリのブランチには、フォルダと JSON ファイルの両方がジョブ名と一致する `my-visual-job/my-visual-job.json` パスが含まれている必要があります。

```
{  
  "name" : "my-visual-job",
```

```

"description" : "",
"role" : "arn:aws:iam::aws_account_id:role/Rolename",
"command" : {
  "name" : "glueetl",
  "scriptLocation" : "s3://foldername/scripts/my-visual-job.py",
  "pythonVersion" : "3"
},
"codegenConfigurationNodes" : "{\"node-nodeID\":{\"S3CsvSource\":
{\\\"AdditionalOptions\\\":{\\\"EnableSamplePath\\\":false,\\\"SamplePath\\\":\\\"s3://notebook-
test-input/netflix_titles.csv\\\"},\\\"Escaper\\\":\\\"\\\",\\\"Exclusions\\\":[],\\\"Name\\\":\\\"Amazon
S3\\\",\\\"OptimizePerformance\\\":false,\\\"OutputSchemas\\\":[{\\\"Columns\\\":[{\\\"Name\\\":
\\\"show_id\\\",\\\"Type\\\":\\\"string\\\"},{\\\"Name\\\":\\\"type\\\",\\\"Type\\\":\\\"string\\\"},{\\\"Name\\\":
\\\"title\\\",\\\"Type\\\":\\\"choice\\\"},{\\\"Name\\\":\\\"director\\\",\\\"Type\\\":\\\"string\\\"},{\\\"Name\\\":
\\\"cast\\\",\\\"Type\\\":\\\"string\\\"},{\\\"Name\\\":\\\"country\\\",\\\"Type\\\":\\\"string\\\"},{\\\"Name\\\":
\\\"date_added\\\",\\\"Type\\\":\\\"string\\\"},{\\\"Name\\\":\\\"release_year\\\",\\\"Type\\\":\\\"bigint\\\"},
{\\\"Name\\\":\\\"rating\\\",\\\"Type\\\":\\\"string\\\"},{\\\"Name\\\":\\\"duration\\\",\\\"Type\\\":\\\"string
\\\"},{\\\"Name\\\":\\\"listed_in\\\",\\\"Type\\\":\\\"string\\\"},{\\\"Name\\\":\\\"description\\\",\\\"Type
\\\":\\\"string\\\"}]}]}},\\\"Paths\\\":[\\\"s3://dalamgir-notebook-test-input/netflix_titles.csv
\\\"],\\\"QuoteChar\\\":\\\"quote\\\",\\\"Recurse\\\":true,\\\"Separator\\\":\\\"comma\\\",\\\"WithHeader
\\\":true}}}"
}

```

スクリプトジョブの場合:

- フォルダ、ジョブ定義の JSON ファイル、スクリプトが必要です。
- フォルダと JSON ファイルはジョブ名と一致している必要があります。スクリプト名は、ファイル拡張子と共にジョブ定義の `scriptLocation` と一致している必要があります

例えば、以下のジョブ定義では、リポジトリ内のブランチにパス `my-script-job/my-script-job.json` および `my-script-job/my-script-job.py` が含まれている必要があります。スクリプト名は、スクリプトの拡張子を含む `scriptLocation` の名前と一致している必要があります。

```

{
  "name" : "my-script-job",
  "description" : "",
  "role" : "arn:aws:iam::aws_account_id:role/Rolename",
  "command" : {
    "name" : "glueetl",
    "scriptLocation" : "s3://foldername/scripts/my-script-job.py",

```

```
"pythonVersion" : "3"  
}  
}
```

制限事項

- 現在、AWS Glue は「[GitLab-Groups](#)」からのプッシュ・プルをサポートしていません。

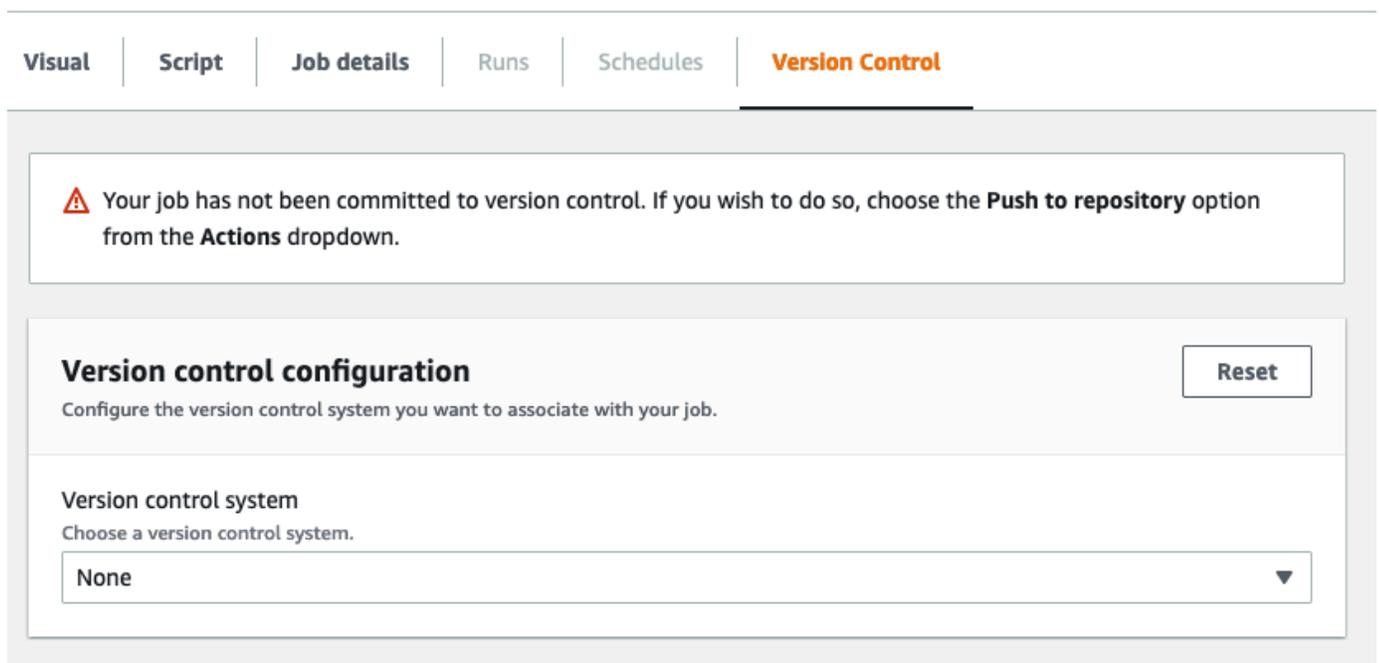
バージョン管理リポジトリと AWS Glue の接続

バージョン管理リポジトリの詳細を入力して、AWS Glue Studio ジョブエディタの [Version Control] (バージョン管理) タブで管理できます。Git のリポジトリと統合するには、AWS Glue Studio にログインするたびにリポジトリに接続する必要があります。

Git バージョン管理システムに接続するには:

1. AWS Glue Studio で、新しいジョブを開始し、[Version Control] (バージョン管理) タブを選択します。

Untitled job 



Visual | Script | Job details | Runs | Schedules | **Version Control**

 Your job has not been committed to version control. If you wish to do so, choose the **Push to repository** option from the **Actions** dropdown.

Version control configuration Reset

Configure the version control system you want to associate with your job.

Version control system
Choose a version control system.

None ▼

2. [Version control system] で、ドロップダウンメニューをクリックし、利用可能なオプションから Git サービスを選択します。

- AWS CodeCommit
 - GitHub
 - GitLab
 - Bitbucket
3. 選択した Git バージョン管理システムによって、入力するフィールドは異なります。

AWS CodeCommit の場合

ジョブのリポジトリとブランチを選択して、リポジトリの設定を完了します。

- [Repository] (リポジトリ) - AWS CodeCommit でリポジトリを設定した場合は、ドロップダウンメニューからリポジトリを選択します。リストに自動的にリポジトリが入力されます
- [Branch] (ブランチ) - ドロップダウンメニューからブランチを選択します。
- [Folder] (フォルダ) - オプション - ジョブを保存するフォルダの名前を入力します。空のままにすると、フォルダが自動的に作成されます。フォルダ名はデフォルトでジョブ名になります。

GitHub の場合:

次のフィールドに入力して GitHub の設定を完了します。

- [Personal access token] (個人用のアクセストークン) - GitHub リポジトリによって提供されるトークンです。個人用のアクセストークンの詳細については、「[GitHub Docs](#)」を参照してください。
- [Repository owner] (リポジトリ所有者) - GitHub リポジトリの所有者です。

GitHub からリポジトリとブランチを選択して、リポジトリの設定を完了します。

- [Repository] (リポジトリ) - GitHub でリポジトリを設定している場合は、ドロップダウンメニューからリポジトリを選択します。リストに自動的にリポジトリが入力されます
- [Branch] (ブランチ) - ドロップダウンメニューからブランチを選択します。
- [Folder] (フォルダ) - オプション - ジョブを保存するフォルダの名前を入力します。空のままにすると、フォルダが自動的に作成されます。フォルダ名はデフォルトでジョブ名になります。

GitLab の場合:

Note

現在、AWS Glue は「[GitLab-Groups](#)」からのプッシュ・プルをサポートしていません。

- [個人用アクセストークン] - これは GitLab リポジトリによって提供されるトークンです。個人用のアクセストークンの詳細については、「[GitLab の個人アクセストークン](#)」を参照してください。
- [リポジトリの所有者] - これは GitLab リポジトリの所有者です。

GitLab からリポジトリとブランチを選択して、リポジトリの設定を完了します。

- [リポジトリ] - GitLab でリポジトリを設定している場合は、ドロップダウンメニューからリポジトリを選択します。リストに自動的にリポジトリが入力されます
- [Branch] (ブランチ) - ドロップダウンメニューからブランチを選択します。
- [Folder] (フォルダ) - オプション - ジョブを保存するフォルダの名前を入力します。空のままにすると、フォルダが自動的に作成されます。フォルダ名はデフォルトでジョブ名になります。

Bitbucket の場合:

- [アプリパスワード] - Bitbucket はリポジトリアクセストークンではなく、アプリパスワードを使用します。アプリパスワードの詳細については、「[アプリパスワード](#)」を参照してください。
- [リポジトリの所有者] - これは Bitbucket リポジトリの所有者です。Bitbucket では、所有者はリポジトリの作成者です。

Bitbucket からワークスペース、リポジトリ、ブランチ、フォルダを選択することで、リポジトリの設定を完了します。

- [ワークスペース] - Bitbucket にワークスペースを設定している場合は、ドロップダウンメニューからワークスペースを選択します。ワークスペースは自動的に入力されます

- [リポジトリ] - Bitbucket でリポジトリを設定した場合は、ドロップダウンメニューからリポジトリを選択します。リポジトリは自動的に入力されます
- [分岐] - ドロップダウンメニューから分岐を選択します。分岐は自動的に入力されます
- [Folder] (フォルダ) - オプション - ジョブを保存するフォルダの名前を入力します。空のままにすると、フォルダがジョブ名で自動的に作成されます。

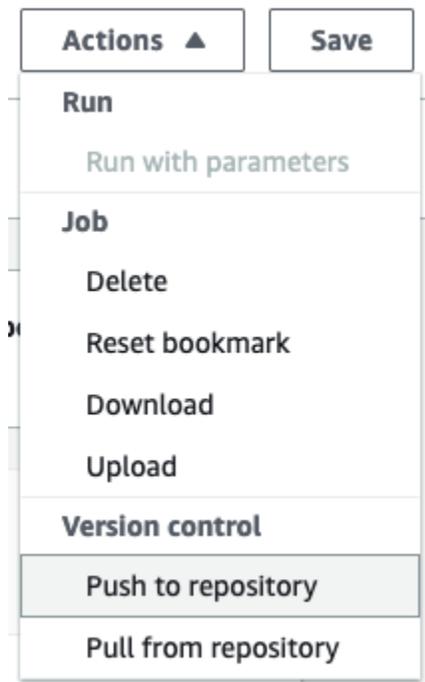
4. AWS Glue Studio ジョブの上部で、[Save] (保存) を選択します。

ソースリポジトリに AWS Glue ジョブをプッシュする

バージョン管理システムの詳細を入力したら、AWS Glue Studio でジョブを編集してソースリポジトリにプッシュできます。プッシュやプルなどの Git の概念に慣れていない場合は、「[Git および AWS CodeCommit の開始方法](#)」で、このチュートリアルを参照してください。

ジョブをリポジトリにプッシュするには、バージョン管理システムの詳細を入力してジョブを保存する必要があります。

1. AWS Glue Studio ジョブで、[Actions] (アクション) を選択します。これにより、追加のメニューオプションが開きます。



2. [Push to repository] (リポジトリにプッシュ) を選択します。

このアクションによってジョブが保存されます。リポジトリにプッシュすると、AWS Glue Studio は最後に保存した変更をプッシュします。リポジトリ内のジョブが自分または別のユー

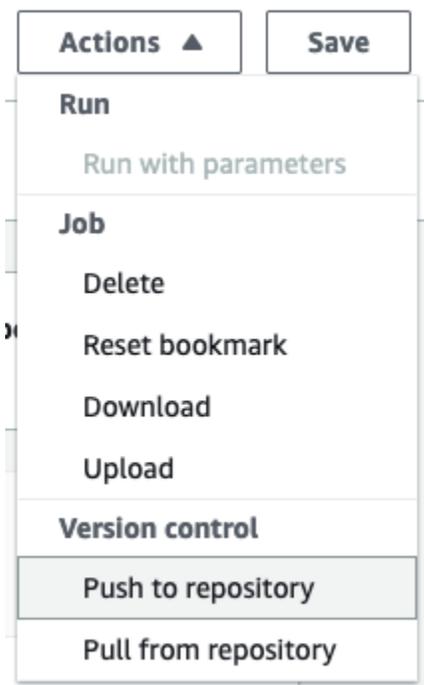
ザーによって変更され、AWS Glue Studio のジョブと同期していない場合は、AWS Glue Studio からジョブをプッシュする際にリポジトリ内のジョブが AWS Glue Studio に保存されたジョブで上書きされます。

3. [Confirm] (確認) を選択してアクションを完了します。これにより、リポジトリに新しいコミットが作成されます。AWS CodeCommit を使用している場合、AWS CodeCommit での最新のコミットへのリンクが確認メッセージに表示されます。

ソースリポジトリからの AWS Glue ジョブのプル

Git リポジトリの詳細を [Version control] (バージョン管理) タブに入力したら、リポジトリからジョブを取得して AWS Glue Studio で編集することもできます。

1. AWS Glue Studio ジョブで、[Actions] (アクション) を選択します。これにより、追加のメニューオプションが開きます。



2. [Pull from repository] (リポジトリからプル) を選択します。
3. [確認] を選択します。これにより、リポジトリから最新のコミットが取得され、AWS Glue Studio のジョブが更新されます。
4. AWS Glue Studio でジョブを編集します。変更を加えた場合は、[Actions] (アクション) ドロップダウンメニューから [Push to repository] (リポジトリにプッシュ) を選択して、ジョブをリポジトリに同期できます。

AWS Glue Studio ノートブックによるコードの作成

データエンジニアは、AWS Glue Studio のインタラクティブノートブックインターフェイス、または AWS Glue のインタラクティブセッションを使用して、従来よりも迅速かつ簡単に AWS Glue ジョブを作成できます。

トピック

- [ノートブックの使用の概要](#)
- [AWS Glue Studio 中のノートブックを使用した ETL ジョブの作成](#)
- [ノートブックエディタコンポーネント](#)
- [ノートブックとジョブスクリプトの保存](#)
- [ノートブックセッションの管理](#)
- [CodeWhisperer との併用 AWS Glue Studio notebooks](#)

ノートブックの使用の概要

AWS Glue Studio は、Jupyter Notebooks に基づいて、ノートブックインターフェイス中でジョブをインタラクティブにオーサリングできるようにします。AWS Glue Studio 中のノートブックを通じて、完全なジョブを実行せずにジョブスクリプトを編集して出力を表示できます。また、完全なジョブを実行せずにデータ統合コードを編集して出力を表示できます。また、マークダウンを追加してノートブックを [ipynb] ファイルとジョブスクリプトとして保存できます。ソフトウェアをローカルにインストールしたり、サーバーを管理したりすることなく、ノートブックを起動できます。コードに満足したら、ボタンをクリックするだけで、AWS Glue Studio はノートブックを Glue ジョブに変換できます。

ノートブックを使用する利点のいくつかを以下に示します：

- プロビジョニングまたは管理するクラスターがない
- 支払うべきアイドルクラスターがない
- 事前設定は不要
- Jupyter Notebook のインストールは不要
- AWS Glue の ETL と同じランタイム/プラットフォーム

AWS Glue Studio を通してスルーノートブックを起動すると、ほんの数秒後にデータを調査し、ジョブスクリプトの開発を開始できるように、すべての設定手順が実行されます。AWS Glue Studio

は、AWS Glue Jupyter カーネルによって Jupyter Notebook を設定します。このノートブックを使用するために VPC、ネットワーク接続、または開発エンドポイントを構成する必要はありません。

ノートブックインターフェイスを使用してジョブを作成するには、次の手順を実行します:

- 必要な IAM 許可を設定します。
- ノートブックセッションを開始してジョブを作成します
- ノートブックのセルにコードを記述します。
- コードを実行してテストし、出力を表示します。
- ジョブを保存します。

ノートブックを保存すると、ノートブックは満杯の AWS Glue のジョブ です。スケジュールジョブの実行、の設定、ジョブパラメータの設定、およびノートブックのすぐ横にあるジョブ実行履歴の表示など、ジョブのすべての側面を管理できます。

AWS Glue Studio の中のノートブックを使用した ETL ジョブの作成

AWS Glue Studio のコンソール中でノートブックの使用を開始するには

1. AWS Identity and Access Management のポリシーを AWS Glue Studio ユーザーにアタッチし、ETL ジョブとノートブックの IAM ロールを作成します。
2. [IAM ロールに対するアクセス許可の付与](#) の説明に従って、ノートブックに追加の IAM セキュリティを設定します。
3. <https://console.aws.amazon.com/gluestudio/> で AWS Glue Studio コンソールを開きます。

Note

ブラウザがサードパーティの Cookie をブロックしていないことを確認してください。ブラウザのデフォルトまたはユーザーの設定でサードパーティの Cookie をブロックしている場合、ノートブックは起動しません。Cookie の管理の詳細については、以下を参照してください。

- [Chrome](#)
- [Firefox](#)
- [Safari](#)

4. ナビゲーションメニューの左側の [ジョブ] リンクを選択します。
5. [Jupyter notebook] を選択し、[Create] をクリックして、新しいノートブックセッションを開始します。
6. [Create job in Jupyter notebook] (Jupyter Notebook でジョブを作成する) ページで、ジョブ名を指定して、使用する IAM ロールを選択します。[ジョブの作成] を選択します。

短時間経過後に、ノートブックエディタが表示されます。

7. コードを追加したら、セルを実行してセッションを開始します。セルを実行するには複数の方法があります。
 - 再生ボタンを押します。
 - キーボードショートカットを使用します。
 - MacOS では、Command + Enter キーでセルを実行します。
 - Windows では、Shift + Enter キーでセルを実行します。

Jupyter ノートブックインターフェースを使用したコードの記述については、「[The Jupyter Notebook User Documentation](#)」を参照してください。

8. スクリプトをテストするには、スクリプト全体、または個々のセルを実行します。コマンド出力は、セルの下の領域に表示されます。
9. スクリプトの開発が完了したら、ジョブを保存して実行できます。スクリプトは、[Script] (スクリプト) タブにあります。ノートブックに追加したマジックはすべて取り除かれ、生成された AWS Glue ジョブのスクリプトの一部として保存されません。AWS Glue Studioは、ノートブックの内容から生成されたスクリプトの最後に `job.commit()` を自動的に追加します。

ジョブの実行のさらなる詳細については、「[ジョブの実行の開始](#)」を参照してください。

ノートブックエディタコンポーネント

ノートブックエディタのインターフェイスには、次の主要なセクションがあります。

- ノートブックインターフェイス (メインパネル) とツールバー
- Job 編集タブ

ノートブックエディタ

AWS Glue Studio のノートブックエディタは Jupyter Notebook Application に基づいています。AWS Glue Studio のノートブックのインターフェースは、Jupyter Notebook が提供するものと似ています。これについては、セクション [[Notebook user interface](#)] で説明されています。インタラクティブセッションにより使用されるノートブックは Jupyter Notebook です。

AWS Glue Studio のノートブックは Jupyter Notebooks に似ていますが、いくつかの重要な点で異なります。

- 現在、AWS Glue Studio のノートブックには拡張をインストールできません。
- 複数のタブを使用することはできません。ジョブとノートブックの間には 1:1 の関係があります
- AWS Glue Studio のノートブックには、Jupyter Notebooks 中に存在する、同じトップファイルメニューがありません。
- 現在、AWS Glue Studio のノートブックでは、AWS Glue のカーネルを用いてのみ動作します。カーネルは、ご自分では更新できないことに注意してください。

AWS Glue Studio のジョブ編集タブ

ETL ジョブを操作するために使用するタブは、ノートブックページの上部にあります。これらは、AWS Glue Studio のビジュアルジョブエディタに表示されるタブに似ていて、同じアクションを実行します。

- Notebook — このタブを使用して、ノートブックインターフェイスを用いてジョブスクリプトを表示します。
- Job details — ジョブ実行の環境とプロパティを設定します。
- Runs — このジョブの以前の実行に関する情報を表示します。
- Schedules — 特定の時間にジョブを実行するためのスケジュールを設定します。

ノートブックとジョブスクリプトの保存

ノートブックと作成中のジョブスクリプトは、いつでも保存できます。右上隅の [Save] ボタンを単に選択し、ビジュアルエディタまたはスクリプトエディタを使用しているかのように同じものを選択してください。

[Save] (保存) を選択した場合、ノートブックファイルはデフォルトの場所に保存されます。

- デフォルトでは、ジョブスクリプトは、[Job Details] (ジョブの詳細) タブの [Advanced properties] (詳細プロパティ) における、[Job details] (ジョブの詳細) プロパティの [Script path] (スクリプトパス) に示される Amazon S3 ロケーションに保存されます。ジョブスクリプトは、Scripts というサブフォルダに保存されます。
- デフォルトでは、ノートブックファイル (.ipynb) は、[Job Details] (ジョブの詳細) タブの [Advanced properties] (詳細プロパティ) における、[Job details] (ジョブの詳細) の [Script path] (スクリプトパス) に示される Amazon S3 ロケーションに保存されます。ノートブックファイルは、Notebooks というサブフォルダに保存されます。

Note

ジョブを保存すると、ジョブスクリプトにはノートブックのコードセルのみが含まれます。Markdown セルとマジックはジョブスクリプトに含まれません。ただし、.ipynb ファイルには Markdown とマジックが含まれます。

ジョブを保存したら、ノートブックで作成したスクリプトを使用してジョブを実行できます。

ノートブックセッションの管理

AWS Glue Studio 中のノートブックは、AWS Glue のインタラクティブセッション機能に基づいています。インタラクティブセッションの使用にはコストがかかります。コストの管理に役立つように、アカウント用に作成されたセッションをモニタリングし、すべてのセッションのデフォルト設定を構成できます。

すべてのノートブックセッションのデフォルトのタイムアウトを変更します。

デフォルトでは、プロビジョニングされた AWS Glue Studio ノートブックの起動の後に実行されたセルがない場合、そのノートブックは 12 時間後にタイムアウトします。これに関連して料金が発生することはありません。また、このタイムアウトは設定できません。

セルを実行することにより、インタラクティブなセッションが開始されます。このセッションのデフォルトのタイムアウトは 48 時間です。このタイムアウトは、セルを実行する前に `%idle_timeout` マジックを渡すことで設定が可能です。

AWS Glue Studio 中のノートブックのデフォルトのセッションタイムアウトを変更するには

1. ノートブックに、セル内に `%idle_timeout` のマジックを入力し、タイムアウト値を分単位で指定します。

- 例: `%idle_timeout 15` はデフォルトのタイムアウトを 15 分に変更します。15 分以内にセッションを使用しない場合、セッションは自動的に停止します。

追加 Python モジュールのインストール

[pip] を使用してセッションに追加のモジュールをインストールする場合は、`%additional_python_modules` を使用してセッションに追加します。

```
%additional_python_modules awscli, s3://mybucket/mymodule.whl
```

[additional_python_modules] へのすべての引数が `pip3 install -m <>` に渡されます。

利用可能な Python モジュールのリストについては、「[AWS Glue での Python ライブラリの使用](#)」を参照してください。

AWS Glue の設定の変更

マジックを使って AWS Glue ジョブの設定値を管理できます。ジョブ設定値を変更する場合は、ノートブックで適切なマジックを使用する必要があります。「[Magics supported by AWS Glue interactive sessions for Jupyter](#)」を参照してください。

Note

実行中のセッションのプロパティを上書きできなくなりました。セッションの設定を変更する場合は、セッションを停止し、新しい構成を設定してから、新しいセッションを開始します。

AWS Glue はさまざまなワーカータイプをサポートしています。ワーカータイプは、`%worker_type` を用いて設定できます。例: `%worker_type G.2X`。デフォルトは G.1X です。

`%number_of_workers` を使用してワーカー数を指定することもできます。例えば、40 人のワーカーを指定するには: `%number_of_workers 40`。

さらなる詳細については、「[Defining Job Properties](#)」を参照してください。

ノートブックセッションを停止します。

ノートブックセッションを停止するには、マジックの `%stop_session` を使用します。

AWS コンソールでノートブックから離れると、セッションの停止を選択できる警告メッセージが表示されます。

CodeWhisperer との併用 AWS Glue Studio notebooks

AWS Glue Studio は、Jupyter Notebooks に基づいて、ノートブックインターフェイス中でジョブをインタラクティブにオーサリングできるようにします。を使用すると、CodeWhisperer ノートブック内でのオーサリングエクスペリエンスが向上します。AWS Glue Studio

Amazon CodeWhisperer エクステンションは、推奨コードを生成したり、コードの問題に関連する改善を提案したりすることで、コードの記述をサポートします。

Amazon とは何ですか CodeWhisperer?

Amazon CodeWhisperer は、開発者の生産性の向上に役立つ機械学習を活用したサービスです。CodeWhisperer これは、自然言語での開発者のコメントと IDE 内のコードに基づいて推奨コードを生成することで実現しています。プレビュー期間中、Amazon CodeWhisperer では Java、Python、JavaScript、C#、TypeScript およびプログラミング言語をご利用いただけます。このサービスは JupyterLab、Amazon SageMaker Studio、Amazon SageMaker ノートブックインスタンス、およびその他の統合開発環境 (IDE) と統合されます。

詳しくは、「[でのセットアップ CodeWhisperer](#)」を参照してください。AWS Glue Studio

コンソール上の AWS Glue ジョブ実行のステータス

AWS Glue 抽出、変換、ロード (ETL) ジョブのステータスは、実行中または停止後に表示できません。AWS Glue コンソールを使用してステータスを表示できます。ジョブ実行のステータスの詳細については、「[the section called “ジョブ実行ステータス”](#)」を参照してください。

ジョブモニタリングダッシュボードにアクセスする

ジョブモニタリングダッシュボードにアクセスするには、AWS Glue ナビゲーションペインの [Monitoring] (モニタリング) リンクを選択します。

ジョブモニタリングダッシュボードの概要

ジョブモニタリングダッシュボードでは、ジョブの実行に関する全体的な概要と、ステータスが [Running] (実行中)、[Canceled] (キャンセル済み)、[Success] (成功)、または [Failed] (失敗) となっているジョブの合計を表示します。追加のタイルには、ジョブの実行の全体的な成功率、ジョブの DPU 使用率、ジョブタイプ、ワーカータイプ、日別のジョブステータスの内訳が表示されます。

タイル内のグラフはインタラクティブです。グラフ内の任意のブロックを選択して、ページの下部にある [Job runs] (ジョブの実行) テーブルで、それらのジョブのみを表示するフィルターを実行します。

[Date range] (日付範囲) セレクタを使用して、このページに表示される情報の日付範囲を変更できます。日付範囲を変更すると、現在の日付より前の指定した日数の値を表示するよう情報のタイルが調整されます。また、[Date range] (日付範囲) セレクタから [Custom] (カスタム) を選択して、特定の日付範囲を使用することもできます。

ジョブの実行ビュー

Note

ジョブ実行履歴には、ワークフローとジョブ実行のために 90 日間アクセスできます。

[Job runs] (ジョブの実行) リソースリストには、指定した日付範囲とフィルターのジョブが表示されます。

ステータス、ワーカータイプ、ジョブタイプ、ジョブ名などの追加の基準でジョブをフィルタリングできます。テーブルの上部にあるフィルターボックスに、フィルターとして使用するテキストを入力できます。テキストを入力すると、一致するテキストを含む行でテーブルの結果が更新されます。

ジョブモニタリングダッシュボードのグラフから要素を選択すると、ジョブのサブセットを表示できます。例えば、[Job runs summary] (ジョブの実行のサマリー) タイルで実行中のジョブの数を選択する場合、[Job runs] (ジョブの実行) リストには、その時点で Running のステータスにあるジョブのみが表示されます。[Worker type breakdown] (ワーカータイプの内訳) の棒グラフでいずれかのバーを選択する場合、ワーカータイプとステータスが一致するジョブの実行のみが [Job runs] (ジョブの実行) リストに表示されます。

[Job runs] (ジョブの実行) リソースリストには、ジョブの実行の詳細が表示されます。列見出しを選択して、テーブル内の行を並べ替えることができます。テーブルには次の情報が含まれます。

プロパティ	説明
ジョブ名	ジョブの名前。
タイプ	ジョブの環境のタイプ <ul style="list-style-type: none">• Glue ETL: AWS Glue が管理する Apache Spark 環境で実行します。• Glue Streaming: Apache Spark 環境で実行し、データストリームで ETL を実行します。• Python シェル: Python スクリプトをシェルとして実行します。
開始時間	このジョブ実行が開始された日付と時刻。
終了時間	このジョブ実行が完了した日付と時刻。
実行ステータス	現在のジョブ実行の状態。値は次のようになります。 <ul style="list-style-type: none">• STARTING• RUNNING• STOPPING• STOPPED• SUCCEEDED• FAILED• TIMEOUT
実行時間	ジョブの実行でリソースを消費した時間。
容量	このジョブの実行に割り当てられる AWS Glue データ処理ユニット (DPU) の数。容量計画の詳細については、AWS Glue デベロッパーガイドの「 Monitoring for DPU Capacity Planning 」を参照してください。

プロパティ	説明
ワーカータイプ	<p>ジョブの実行時に割り当てられた事前定義済みのワーカーのタイプ。値は G.1X、G.2X、G.4X または G.8X になります。</p> <ul style="list-style-type: none">• G.1X – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは、84 GB のディスク (約 34 GB の空き) を備えた 1 DPU (4 vCPU、16 GB のメモリ) にマッピングされます。メモリを大量に消費するジョブには、このワーカータイプをお勧めします。これが AWS Glue バージョン 2.0 以降のジョブの [Worker type] (ワーカータイプ) のデフォルトです。• G.2X – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは、128 GB のディスク (約 77 GB の空き) を備えた 2 DPU (8 vCPU、32 GB のメモリ) にマッピングされます。メモリを大量に消費するジョブには、機械学習変換を実行するこのワーカータイプをお勧めします。• G.4X – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは、256 GB のディスク (約 235 GB の空き) を備えた 4 DPU (16 vCPU、64 GB のメモリ) にマッピングされます。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、AWS Glue バージョン 3.0 以降の Spark ETL ジョブでのみ使用可能で、以下の AWS リージョンで使用でき

プロパティ	説明
	<p>ます。米国東部 (オハイオ)、米国東部 (バージニア北部)、米国西部 (オレゴン)、アジアパシフィック (シンガポール)、アジアパシフィック (シドニー)、アジアパシフィック (東京)、カナダ (中部)、欧州 (フランクフルト)、欧州 (アイルランド)、欧州 (ストックホルム)。</p> <ul style="list-style-type: none"> • G.8X – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは、512 GB のディスク (約 487 GB の空き) を備えた 8 DPU (32 vCPU、128 GB のメモリ) にマッピングされます。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、AWS Glue バージョン 3.0 以降の Spark ETL ジョブでのみ使用可能で、G.4X ワーカータイプでサポートされているのと同じ AWS リージョンで使用できます。
DPU 時間	<p>ジョブの実行に使用される DPU の推定数。DPU は、処理能力の相対的な尺度です。DPU は、ジョブの実行のコストを割り出すために使用されます。詳細については、「AWS Glue 料金表ページ」を参照してください。</p>

リストから任意のジョブの実行を選択し、追加の情報を表示できます。ジョブの実行を選択して、次のいずれかを実行します。

- [Actions] (アクション) メニューから [View job] (ジョブの表示) オプションを選択して、ビジュアルエディタでジョブを表示します。

- [Actions] (アクション) メニューから [Stop run] (実行の停止) オプションを使用して、現在行われているジョブ実行を停止します。
- [View CloudWatch logs] (CloudWatch ログを表示) ボタンをクリックして、そのジョブのジョブの実行ログを表示します。
- [詳細を表示] を選択して、ジョブの実行の詳細ページを表示します。

ジョブの実行ログの表示

ジョブのログは、さまざまな方法で表示できます。

- [Monitoring] (モニタリング) ページの [Job runs] (ジョブの実行) テーブルで、ジョブの実行を選択し、[View CloudWatch logs] (CloudWatch ログを表示) を選択します。
- ビジュアルジョブエディタの、ジョブの [Runs] (実行) タブで、ハイパーリンクを選択してログを表示できます。
- Logs (ログ) – ジョブの実行で継続的なログ記録が有効であるときに書き込まれる Apache Spark ジョブのログへのリンクです。このリンクを選択すると、/aws-glue/jobs/logs-v2 ロググループ内の Amazon CloudWatch ログに移動します。デフォルトでは、不要な Apache Hadoop YARN ハートビート、Apache Spark ドライバー、エグゼキューターログメッセージはログから除外されています。継続的なログ記録の詳細については、AWS Glue デベロッパーガイドの「[Continuous Logging for AWS Glue Jobs](#)」を参照してください。
- Error logs (エラーログ) – このジョブの実行で stderr に書き込まれるログへのリンクです。このリンクを選択すると、/aws-glue/jobs/error ロググループ内の Amazon CloudWatch ログに移動します。これらのログを使用して、ジョブの実行中に発生したエラーに関する詳細を表示できます。
- Output logs (出力ログ) – このジョブの実行で stdout に書き込まれるログへのリンクです。このリンクを選択すると、/aws-glue/jobs/output ロググループ内の Amazon CloudWatch ログに移動します。これらのログを使用して、AWS Glue Data Catalog で作成されたテーブルに関する詳細と、発生したエラーをすべて確認することができます。

ジョブの実行の詳細を表示する

[Monitoring] (モニタリング) ページの [Job runs] (ジョブの実行) リストでジョブを選択し、[View run details] (実行の詳細を表示する) をクリックして、ジョブの実行に関する詳細情報を表示できます。

ジョブの実行の詳細ページには、次の情報が表示されます。

プロパティ	説明
ジョブ名	ジョブの名前。
実行ステータス	現在のジョブ実行の状態。値は次のようになります。 <ul style="list-style-type: none">• STARTING• RUNNING• STOPPING• STOPPED• SUCCEEDED• FAILED• TIMEOUT
Glue バージョン	ジョブ実行に使用される AWS Glue バージョン
最近の試行	ジョブ実行時の自動再試行回数
開始時間	このジョブ実行が開始された日付と時刻。
終了時間	このジョブ実行が完了した日付と時刻。
起動時間	ジョブ実行の準備にかかった時間
実行時間	ジョブスクリプトの実行にかかった時間
トリガー名	ジョブに関連付けられているトリガーの名前
最終更新日	ジョブが最後に変更された日付
セキュリティ設定	ジョブのセキュリティ設定。これには、Amazon S3 の暗号化、CloudWatch の暗号化、ジョブブックマークの暗号化設定が含まれます。
タイムアウト	ジョブ実行のタイムアウトのしきい値

プロパティ	説明
割り当てられた容量	このジョブの実行に割り当てられる AWS Glue データ処理ユニット (DPU) の数。容量計画の詳細については、AWS Glue デベロッパガイドの「 Monitoring for DPU Capacity Planning 」を参照してください。
最大容量	ジョブの実行に使用可能な最大の容量。
ワーカー数	ジョブ実行に使用されるワーカーの数
ワーカータイプ	<p>ジョブの実行に割り当てられる定義済みのワーカータイプ。値は、G.1X または G.2X になります。</p> <ul style="list-style-type: none"> G.1X – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは 1 DPU (4 vCPU、16 GB のメモリ、64 GB のディスク) にマッピングされており、ワーカーごとに 1 個のエグゼキューターがあります。メモリを大量に消費するジョブには、このワーカータイプをお勧めします。これが AWS Glue バージョン 2.0 以降のジョブの [Worker type] (ワーカータイプ) のデフォルトです。 G.2X – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは 2 DPU (8 vCPU、32 GB のメモリ、128 GB のディスク) にマッピングされており、ワーカーごとに 1 個のエグゼキューターがあります。メモリを大量に消費するジョブには、機械学習変換を実行するこのワーカータイプをお勧めします。

プロパティ	説明
ログ	継続的なログ記録のためのジョブログへのリンク (/aws-glue/jobs/logs-v2)。
出力ログ	ジョブの出力ログファイルへのリンク (/aws-glue/jobs/output)。
エラーログ	ジョブのエラーログファイルへのリンク (/aws-glue/jobs/error)。

次の追加項目も表示できます。これらは、最近のジョブ実行の情報を表示する際に表示できます。詳細については、[the section called “最近のジョブの実行の情報を表示する”](#) を参照してください。

- [入力引数]
- 連続ログ
- [メトリクス] – 基本的なメトリクスのビジュアライゼーションを確認できます。含まれるメトリクスの詳細については、「[the section called “Spark ジョブ実行の Amazon CloudWatch メトリクスを表示”](#)」を参照してください。
- [Spark UI] – Spark UI でジョブの Spark ログを視覚化できます。Spark Web UI の使用の詳細については、「[the section called “Spark UI を使用したモニタリング”](#)」を参照してください。[the section called “ジョブ用の Spark UI の有効化”](#) のステップに従って、この機能を有効にします。

Spark ジョブ実行の Amazon CloudWatch メトリクスを表示

ジョブの実行の詳細ページの [Run details] (実行の詳細) セクションでは、ジョブのメトリクスを表示できます。AWS Glue Studio はジョブ実行のたびにジョブメトリクスを Amazon CloudWatch に送信します。

AWS Glue では、30 秒ごとに Amazon CloudWatch にメトリクスが報告されます。AWS Glue メトリクスは、以前に報告された値からデルタ値を表示します。必要に応じて、メトリクスのダッシュボードにより 30 秒の値が集計 (合計) され、直近 1 分間分の値を取得できます。一方、AWS Glue によって Amazon CloudWatch に渡される Apache Spark メトリクスは、一般的に報告された時点での状態を表す絶対値です。

Note

Amazon CloudWatch にアクセスできるようにアカウントを設定する必要があります。

メトリクスには、次のようなジョブの実行に関する情報が表示されます。

- ETL Data Movement (ETL データ移動) – Amazon S3 から読み取られたバイト数またはそこに書き込まれたバイト数。
- Memory Profile: Heap used (プロファイル: ヒープ使用) – Java 仮想マシン (JVM) ヒープによって使用されたメモリのバイト数。
- Memory Profile: heap usage (メモリプロファイル: ヒープ使用量) – JVM ヒープによって使用されたメモリの割合 (スケール: 0~1)。百分率で示されます。
- CPU Load (CPU 負荷) – 使用された CPU システムの負荷の割合 (スケール: 0~1)。百分率で示されます。

Ray ジョブ実行の Amazon CloudWatch メトリクスを表示

ジョブの実行の詳細ページの [Run details] (実行の詳細) セクションでは、ジョブのメトリクスを表示できます。AWS Glue Studio はジョブ実行のたびにジョブメトリクスを Amazon CloudWatch に送信します。

AWS Glue では、30 秒ごとに Amazon CloudWatch にメトリクスが報告されます。AWS Glue メトリクスは、以前に報告された値からデルタ値を表示します。必要に応じて、メトリクスのダッシュボードにより 30 秒の値が集計 (合計) され、直近 1 分間分の値を取得できます。一方、AWS Glue によって Amazon CloudWatch に渡される Apache Spark メトリクスは、一般的に報告された時点での状態を表す絶対値です。

Note

次に示されているように、Amazon CloudWatch にアクセスできるようにアカウントを設定する必要があります。

Ray ジョブでは、次の集約メトリクスグラフを表示できます。これらを使用すると、クラスターとタスクのプロファイルを作成し、各ノードの詳細情報にアクセスできます。これらのグラフを裏付ける時系列データは、CloudWatch でさらに分析できます。

タスクプロファイル: タスクの状態

システム内の Ray タスクの数を表示します。各タスクのライフサイクルには、独自の時系列が割り当てられます。

タスクプロファイル: タスク名

システム内の Ray タスクの数を表示します。保留中のタスクとアクティブなタスクのみが表示されます。タスクのタイプ (名前別) ごとに独自の時系列が割り当てられます。

クラスタープロファイル: 使用中の CPU

使用中の CPU コアの数を表示します。ノードごとに独自の時系列が割り当てられます。ノードは IP アドレスで識別されます。IP アドレスは一時的なものであり、識別にのみ使用されます。

クラスタープロファイル: オブジェクトストアのメモリ使用量

Ray オブジェクトキャッシュによるメモリ使用量を表示します。各メモリの場所 (物理メモリ、ディスクへのキャッシュ、Amazon S3 への流出) には、独自の時系列が割り当てられます。オブジェクトストアは、クラスター内に存在するすべてのノードのデータストレージを管理します。詳細については、Ray ドキュメントの「[Objects](#)」を参照してください。

クラスタープロファイル: ノード数

クラスターにプロビジョニングされたノード数を表示します。

ノードの詳細: CPU 使用量

各ノードの CPU 使用率をパーセンテージで表示します。各シリーズは、ノード上に存在する全コアの CPU 使用率の合計パーセンテージを表示します。

ノードの詳細: メモリ使用量

各ノードのメモリ使用量を GB 単位で表示します。各シリーズは、Ray タスクと Plasma ストア プロセスを含む、ノード上のすべてのプロセス間で集約されたメモリを表示します。これには、ディスクに保存されたオブジェクトや Amazon S3 に流出したオブジェクトは反映されません。

ノードの詳細: ディスク使用量

各ノードのディスク使用量を GB 単位で表示します。

ノードの詳細: ディスク I/O 速度

各ノードのディスク I/O を KB/秒単位で表示します。

ノードの詳細: ネットワーク I/O スループット

各ノードのネットワーク I/O を KB/秒単位で表示します。

ノードの詳細: Ray コンポーネントによる CPU 使用量

CPU 使用量をコアの分数で表示します。各ノードの Ray コンポーネントごとに独自の時系列が割り当てられます。

ノードの詳細: Ray コンポーネントによるメモリ使用量

メモリ使用量を GiB 単位で表示します。各ノードの Ray コンポーネントごとに独自の時系列が割り当てられます。

機密データを検出して処理する

Detect PII transform は、データソース内の個人識別情報 (PII) を識別します。エンティティを選択し、データのスキャン方法、Detect PII transform によって識別されてきた PII エンティティで何を行うかを識別します。

Detect PII transform は、定義したエンティティ、または AWS によって事前定義されたエンティティを検出、マスク、削除する機能を提供します。これにより、コンプライアンスを高め、責任を軽減できます。たとえば、読み取り可能な個人を特定できる情報がデータに含まれていないことを確認したい場合や、社会保障番号を固定文字列 (例: xxx-xx-xxxx)、電話番号、住所で隠したい場合があります。

AWS Glue Studio の外部で機密データを使用するには、「[AWS Glue Studio 外でのセンシティブデータ検出の使用](#)」を参照してください。

トピック

- [データのスキャン方法の選択](#)
- [検出する PII エンティティの選択](#)
- [検出感度のレベルの指定](#)
- [特定された PII データによる対処方法の選択](#)
- [詳細なアクションオーバーライドの追加](#)

データのスキャン方法の選択

データセットをスキャンして個人を特定できる情報 (PII) などの機密データを探す際に、各行で PII を検出するか、または PII データを含む列を検出するかを選択できます。

Detect PII in each cell
Scan the entire data set, and act on each occurrence individually.

Detect fields containing PII
To reduce costs and improve performance, sample only a portion of the data and act on fields across all records.

Sample portion

The percentage of rows to sample out of the entire data set.

100 %

Between 1 and 100.

Detection threshold

To consider a field as containing PII, set the minimum percentage of detected rows out of the sampled rows.

10 %

Between 1 and 100.

Detect PII in each cell を選択する場合、データソース内のすべての行をスキャンすることを選択しています。これは、PII エンティティを識別するための包括的なスキャンです。

Detect fields containing PII を選択する場合、PII エンティティの行のサンプルをスキャンすることを選択しています。これは、PII エンティティが見つかったフィールドを特定しながら、コストとリソースを低く抑える方法です。

PII を含むフィールドを検出することを選択した場合、行の一部をサンプリングすることで、コストを削減し、パフォーマンスを向上させることができます。このオプションを選択すると、追加のオプションを指定できます。

- Sample portion: これにより、サンプリングする行の割合を指定できます。例えば、50 と入力すると、PII エンティティのためにスキャンされた行の 50% を指定したいことになります。
- Detection threshold: これにより、列全体が PII エンティティを持つものとして識別されるように、PII エンティティを含む行の割合を指定できます。例えば、10 と入力した場合、フィールドに PII エンティティである米国電話機があると識別されるためには、スキャンされる行の PII エンティティの US Phone の数が 10% 以上になるように指定します。PII エンティティを含む行の割合が 10% 未満の場合、そのフィールドに PII エンティティ (US Phone) が含まれているというラベル付けできません。

検出する PII エンティティの選択

[Detect PII in each cell] (各セルの PII を検出する) を選択した場合は、次の 3 つのオプションのいずれかを選択できます。

- 利用可能なすべての PII パターン (これにはエンティティも含まれます)。AWS
- カテゴリを選択する - カテゴリを選択すると、PII パターンには、選択したカテゴリのパターンが自動的に含まれます。
- 特定のパターンを選択する - 選択したパターンのみが検出されます。

マネージド機密データタイプの全リストについては、「[Managed data types](#)」を参照してください。

利用可能なすべての PII パターンから選択する

「使用可能なすべての PII パターン」を選択した場合は、によって事前定義されたエンティティを選択します。AWSエンティティは、1 つでも、複数でも、すべてでも選択できます。

Select entities to detect



Available entities (19)



Select all

Clear all

Create new

Manage

All categories ▼

< 1 >

<input type="checkbox"/>	Entity name ▼	Category ▲
<input type="checkbox"/>	Person's name	Universal, HIPAA
<input type="checkbox"/>	Email (General)	Universal
<input type="checkbox"/>	Credit Card	Universal
<input type="checkbox"/>	IP Address	Networking
<input type="checkbox"/>	MAC Address	Networking
<input type="checkbox"/>	US Phone	United States, HIPAA
<input type="checkbox"/>	US Passport	United States
<input type="checkbox"/>	Social Security Number (SSN)	United States, HIPAA
<input type="checkbox"/>	US Individual Taxpayer Identification Number (ITIN)	United States, HIPAA
<input type="checkbox"/>	US/Canada bank account	United States, HIPAA
<input type="checkbox"/>	US driving license	HIPAA
<input type="checkbox"/>	Healthcare Common Procedure Coding System (HCPCS) code	HIPAA
<input type="checkbox"/>	National Drug Code (NDC)	HIPAA
<input type="checkbox"/>	National Provider Identifier (NPI)	HIPAA
<input type="checkbox"/>	Drug Enforcement Agency (DEA) Registration Number	HIPAA
<input type="checkbox"/>	Health Insurance Claim Number (HICN)	HIPAA
<input type="checkbox"/>	Medicare Beneficiary Identifier	HIPAA

カテゴリを選択する

PII パターンとして [Select categories] (カテゴリを選択する) を選択して検出する場合は、ドロップダウンメニューのオプションから選択します。一部のエンティティは複数のカテゴリに属する場合がありますことに注意してください。例えば、[Person's name] (人名) は、[Universal] (ユニバーサル) および [HIPAA] カテゴリに属するエンティティです。

- [Universal] (ユニバーサル) (例: [Email] (E メール)、[Credit Card] (クレジットカード))
- [HIPAA] (例: [US driving license] (米国の運転免許証)、[Healthcare Common Procedure Coding System (HCPCS) code] (Healthcare Common Procedure Coding System (HCPCS) コード))
- [Networking] (ネットワーキング) (例: [IP Address] (IP アドレス)、[MAC Address] (MAC アドレス))
- アルゼンチン
- オーストラリア
- オーストリア
- ベルギー
- ボスニア
- ブルガリア
- カナダ
- チリ
- コロンビア
- クロアチア
- キプロス
- チェコ共和国
- デンマーク
- エストニア
- フィンランド
- フランス
- ドイツ
- ギリシャ
- ハンガリー
- アイルランド

- 韓国
- 日本
- メキシコ
- オランダ
- ニュージーランド
- ノルウェー
- ポルトガル
- ルーマニア
- シンガポール
- スロバキア
- スロベニア
- スペイン
- スウェーデン
- スイス
- トルコ
- ウクライナ
- アメリカ
- 英国
- ベネズエラ

特定のパターンを選択する

検出する PII パターンとして [Select specific patterns] (特定のパターンを選択する) を選択すると、作成済みのパターンのリストから検索や参照したり、新しい検出エンティティパターンを作成したりできます。

次のステップでは、機密データを検出する新しいカスタムパターンを作成する方法について説明します。カスタムパターンの名前を入力して、カスタムパターンを作成し、正規表現を追加して、オプションでコンテキスト単語を定義します。

1. 新しいパターンを作成するには、[Create new] (新規作成) ボタンをクリックします。

Select patterns

[Browse](#)[Create new !\[\]\(ad56ff270e1b2c76650301ffd30abae8_img.jpg\)](#)

2. [Create detection entity] (検出エンティティの作成) ページで、エンティティ名と正規表現を入力します。正規表現 (Regex) は、AWS Glue がエンティティを照合するために使用するものです。
3. [Validate] (検証) をクリックします。検証が成功すると、文字列が有効な正規表現であることを示す確認メッセージが表示されます。検証に失敗した場合は、文字列が適切なフォーマット、および許容される文字リテラル、演算子、構文のいずれかに準拠していないことを示すメッセージが表示されます。
4. 正規表現にコンテキスト単語を追加することもできます。コンテキスト単語によって一致する可能性が高くなることがあります。コンテキスト単語は、フィールド名がエンティティを説明していない場合に有効です。例えば、社会保障番号には「SSN」または「SS」という名前が付けられます。これらのコンテキスト単語を追加すると、エンティティの照合に役立ちます。
5. [Create] (作成) をクリックして、検出エンティティを作成します。作成されたエンティティは、AWS Glue Studio コンソールに表示されます。左側のナビゲーションメニューの [Detection entities] (検出エンティティ) をクリックします。

[Detection entities] (検出エンティティ) ページから、検出エンティティの編集、削除、作成ができます。検索フィールドを使用してパターンを検索することもできます。

検出感度のレベルの指定

機密データの検出を使用する場合の感度レベルを設定できます。

- [高] – (デフォルト) より高いレベルの感度が必要なユースケースのために、より多くのエンティティを検出します。2023 年 11 月よりも後に作成されたすべての AWS Glue ジョブは、この設定を自動的にオプトインします。
- [低] - 検出するエンティティの数を減らし、誤検知を減らします。

Select global detection sensitivity

Choose the level of detection sensitivity to apply to your data set.

- High (default)**
Detects more entities for use cases that require a higher level of sensitivity.
- Low**
Detects fewer entities and reduces false positives.

特定された PII データによる対処方法の選択

データソース全体で PII を検出することを選択した場合は、適用するグローバルアクションを選択できます。

- **Enrich data with detection results:** 各セルで Detect PII を選択した場合、検出されたエンティティを新しい列に保存できます。
- **Redact detected text:** 検出された PII 値を、オプションの置換テキスト入力フィールド中に指定した文字列に置き換えることができます。文字列を指定しない場合、検出された PII エンティティは '*****' に指定されます。
- **[検出されたテキストを部分的にマスキング]:** 検出された PII の値の一部を、選択した文字列に置き換えることができます。可能なオプションは 2 つあります。すなわち、端をマスキングしないままにするか、または明示的な正規表現パターンを指定してマスキングするかのいずれかです。AWS Glue 2.0 では、この機能は使用できません。
- **Apply cryptographic hash:** 検出された PII 値を SHA-256 暗号化ハッシュ関数に渡し、その値を関数からの出力に置き換えることができます。

Select global action (required)

Choose an action to take on detected entities.

- DETECT. Enrich data with detection results.**
Create a new column that will contain any entity type detected in that row.
- REDACT. Redact detected text.**
Replace detected entity with a string you choose.
- PARTIAL_REDACT. Partially redact detected text.**
Replace part of a detected entity with a string you choose.
- SHA256_HASH. Apply cryptographic hash.**
Apply a SHA-256 cryptographic hash function to the input string.

AWS Glue バージョン 2.0 と 3.0 以降の相違点

AWS Glue 2.0 件のジョブでは、DataFrame 各列で検出された PII 情報を補足列に含めた新しい情報が返されます。マスキングまたはハッシュ作業は、ビジュアルタブの AWS Glue スクリプト内に表示されます。

AWS Glue 3.0 と 4.0 のジョブでは、DataFrame 同じ補足列を含む新しい情報が返されます。「actionUsed」の新しいキーが存在し、DETECT、REDACT、PARTIAL_REDACT、または SHA256_HASH のいずれかになります。マスキングアクションを選択すると、DataFrame 機密データがマスクされたデータが返されます。

詳細なアクションオーバーライドの追加

追加の検出およびアクションの設定を、詳細なアクションオーバーライドテーブルに追加できます。これにより、次のことが可能になります。

- [検出から特定の列を包含または除外] – データソース上の推論されたスキーマによって、使用可能な列がテーブルに入力されます。
- [グローバルアクションを使用するよりも詳細な特定の設定を指定] – 例えば、エンティティタイプごとに異なるマスキングテキストの設定を指定できます。
- [グローバルアクションとは異なるアクションを指定] – 異なる機密データタイプに異なるアクションを適用する場合は、ここで実行できます。同じ列に対して 2 edit-in-place つの異なるアクション (リダクションとハッシュ) を使用することはできませんが、detect はいつでも使用できることに注意してください。

Fine grained actions (overrides) (0)

[Edit as JSON](#) [Delete](#) [Edit](#) [Add](#)

Select entities to add a fine grained action different from the global action above.

<input type="checkbox"/>	Entity type ▲	Action ▼	Action options	Columns
No overrides				

AWS Glue Studio による ETL ジョブの管理

AWS Glue Studio ではシンプルでグラフィカルなインターフェースを使用して、ETL ジョブを管理できます。ナビゲーションメニューで、[Jobs] (ジョブ) を選択し、[Jobs] (ジョブ) ページを表示します。このページには、AWS Glue Studio または AWS Glue コンソールで作成したすべてのジョブが表示されます。このページでは、ジョブを表示、管理、および実行できます。

またこのページでは、次のタスクを実行できます。

- [ジョブの実行の開始](#)
- [ジョブの実行のスケジュール](#)
- [ジョブスケジュールの管理](#)
- [ジョブの実行の停止](#)
- [ジョブの表示](#)
- [最近のジョブの実行の情報を表示する](#)
- [ジョブスクリプトの表示](#)
- [ジョブのプロパティを変更する](#)
- [ジョブの保存](#)
- [ジョブのクローンを作成する](#)
- [ジョブの削除](#)

ジョブの実行の開始

AWS Glue Studio ではオンデマンドでジョブを実行できます。ジョブは複数回実行でき、ジョブを実行するたびに、AWS Glue によりジョブのアクティビティとパフォーマンスに関する情報が収集されます。この情報は [job run] (ジョブの実行) として参照され、ジョブの実行 ID によって識別されます。

AWS Glue Studio では次の方法でジョブの実行を開始できます。

- [Jobs] (ジョブ) ページで、開始するジョブを選択してから、[Run job] (ジョブの実行) ボタンをクリックします。
- ビジュアルエディタでジョブを表示していて、そのジョブが保存されている場合は、[Run] (実行) ボタンをクリックして、ジョブの実行を開始します。

ジョブの実行についての詳細は、AWS Glue デベロッパーガイドの「[Working with Jobs on the AWS Glue Console](#)」を参照してください。

ジョブの実行のスケジュール

AWS Glue Studio ではスケジュールを作成し、特定の時間にジョブを実行できます。ジョブが実行される回数、曜日、時間などの制約を指定できます。これらの制約は cron に基づいており、cron と同じ制限があります。例えば、毎月 31 日にジョブを実行することを選択した場合、31 日がない月があることに注意してください。cron の詳細については、AWS Glue デベロッパーガイドの「[cron 式](#)」を参照してください。

スケジュールに従ってジョブを実行するには

1. 次のいずれかの方法を使用して、ジョブのスケジュールを作成します。
 - [Jobs] (ジョブ) ページで、スケジュールを作成するジョブを選択し、[Actions] (アクション) を選択してから、[Schedule job] (ジョブのスケジュール) を設定します。
 - ビジュアルエディタでジョブを表示していて、そのジョブが保存されている場合は、[Schedules] (スケジュール) タブを選択します。次に [Create Schedule] (スケジュールの作成) を選択します。
2. [Schedule job run] (ジョブの実行のスケジュール) ページで、次の情報を入力します。
 - Name (名前): ジョブのスケジュールの名前を入力します。
 - Frequency (頻度): ジョブのスケジュールの頻度を入力します。以下のオプションを選択できます。
 - Hourly (毎時): ジョブは 1 時間ごとに実行され、特定の時間に始まります。ジョブを実行する時間の分を指定できます。デフォルトでは、[hourly] (毎時) を選択すると、ジョブは時間の初め (0 分) に実行されます。
 - Daily (毎日): ジョブは毎日実行され、一度に開始します。ジョブを実行する時間の分とジョブの開始時間を指定できます。時間は 23 時間制で指定され、午後の時間には 13 から 23 の数字を使用します。分と時間のデフォルトの値は 0 です。つまり、[Daily] (毎日) を選択している場合、ジョブはデフォルトで午前 0 時に実行されます。
 - Weekly (毎週): ジョブは毎週 1 日以上実行されます。前述の [Daily] (毎日) と同じ設定に加えて、ジョブを実行する曜日を選択できます。1 日以上選択できます。
 - Monthly (毎月): ジョブは毎月、特定の日に実行されます。前述の [Daily] (毎日) と同じ設定に加えて、その月のジョブを実行する日を選択できます。その日を 1 から 31 の数値で指定

します。2月30^日などの、1か月で存在しない日を選択した場合、その月にジョブは実行されません。

- Custom (カスタム): cron 構文を使用して、ジョブのスケジュールの式を入力します。cron 式を使用すると、その月の特定の日ではなく最終日、または3か月ごとに7^日と21^日などのより複雑なスケジュールを作成できます。

AWS Glue デベロッパーガイドの「[cron 式](#)」を参照してください。

- Description (説明): オプションで、ジョブのスケジュールについての説明を入力できます。複数のジョブに同じスケジュールを使用する場合、説明があることでジョブのスケジュールがどうなっているのかを簡単に確認できます。
3. [Create schedule] (スケジュールの作成) をクリックして、ジョブのスケジュールを保存します。
 4. スケジュールを作成すると、コンソールページの上部に成功のメッセージが表示されます。このバナーの [Job details] (ジョブの詳細) をクリックして、ジョブの詳細を表示できます。これにより、[Schedules] (スケジュール) タブが選択されているビジュアルジョブエディタのページが開きます。

ジョブスケジュールの管理

ジョブのスケジュールを作成したら、ビジュアルエディタでジョブを開き、[Schedules] (スケジュール) タブをクリックして、スケジュールを管理できます。

[Schedules] (スケジュール) タブでは、次のタスクを実行できます。

- 新しいスケジュールの作成。

[Create Schedule] (スケジュールの作成) を選択し、[the section called “ジョブの実行のスケジュール”](#) に示されているようにスケジュールの情報を入力します。

- 既存のスケジュールの編集。

編集するスケジュールを選択して、[Action] (アクション)、[Edit schedule] (スケジュールの編集) を選択します。既存のスケジュールの編集を選択すると、[Frequency] (頻度) は [Custom] (カスタム)、スケジュールは cron 式と表示されます。cron 式を変更するか、[Frequency] (頻度) ボタンを使用して新しいスケジュールを指定できます。変更が完了したら、[Update schedule] (スケジュールの更新) を選択します。

- アクティブなスケジュールの一時停止。

アクティブなスケジュールを選択して、[Action] (アクション)、[Pause schedule] (スケジュールの一時停止) を選択します。スケジュールは、すぐに非アクティブ化されます。更新されたジョブスケジュールのステータスを表示するには、[refresh (reload)] (更新 (再ロード)) ボタンをクリックします。

- 一時停止したスケジュールの再開。

非アクティブ化したスケジュールを選択して、[Action] (アクション)、[Resume schedule] (スケジュールの再開) を選択します。スケジュールはすぐにアクティブになります。更新されたジョブスケジュールのステータスを表示するには、[refresh (reload)] (更新 (再ロード)) ボタンをクリックします。

- スケジュールの削除。

削除するスケジュールを選択して、[Action] (アクション)、[Delete schedule] (スケジュールの削除) を選択します。スケジュールはすぐに削除されます。更新したジョブスケジュールのリストを表示するには、[refresh (reload)] (更新 (再ロード)) ボタンを選択します。スケジュールには、完全に削除されるまで [Deleting] (削除中) のステータスが表示されます。

ジョブの実行の停止

ジョブの実行が完了する前に、ジョブを停止できます。ジョブが正しく設定されていないことがわかっている場合、またはジョブの完了に時間がかかりすぎる場合は、このオプションを選択できます。

[Monitoring] (モニタリング) ページの [Job runs] (ジョブの実行) リストで停止するジョブを選択し、[Action] (アクション) を選択してから、[Stop run] (実行の停止) をクリックします。

ジョブの表示

すべてのジョブは、[Jobs] (ジョブ) ページで確認できます。このページにアクセスするには、ナビゲーションペインで [Jobs] (ジョブ) を選択します。

[Jobs] (ジョブ) ページでは、アカウントで作成されたすべてのジョブを確認できます。[Your jobs] (ジョブ) リストには、ジョブ名、タイプ、ジョブの前の実行のステータス、およびジョブが作成されて最後に変更された日付が表示されます。ジョブの名前を選択して、そのジョブの詳細情報を表示できます。

また、モニタリングダッシュボードを使用して、すべてのジョブを表示することもできます。ダッシュボードにアクセスするには、ナビゲーションペインで、[Monitoring] (モニタリング) を選択します。

ジョブの表示のカスタマイズ

[Jobs] (ジョブ) ページの [Your jobs] (ジョブ) セクションで、ジョブの表示方法をカスタマイズできます。また、検索テキストフィールドにテキストを入力して、そのテキストを含む名前のジョブのみを表示することもできます。

[Your jobs] (ジョブ) セクションの設定アイコン

() を選択すると、AWS Glue Studio でテーブルに情報を表示する方法をカスタマイズできます。表示内のテキストの行を折り返したり、ページに表示されるジョブの数を変更したり、表示する列を指定できます。

最近のジョブの実行の情報を表示する

新しいデータをソースの場所に追加する際、ジョブは複数回実行されます。ジョブを実行するたびに、一意の ID が割り当てられ、それに関する情報が収集されます。この情報は、次の方法で表示できます。

- [Runs] (実行) タブをクリックして、現在表示されているジョブの実行の情報を表示します。

[Recent job runs] (最近のジョブの実行) ページの [Runs] (実行) タブには、ジョブの実行ごとにカードがあります。[Runs] (実行) タブで表示される情報には、次の情報が含まれます。

- ジョブの実行 ID
 - このジョブの実行の試行回数
 - ジョブの実行のステータス
 - ジョブの実行開始および終了時刻
 - ジョブの実行のランタイム
 - ジョブのログファイルへのリンク
 - ジョブのエラーログファイルへのリンク
 - 失敗したジョブに返されたエラー
- ジョブ実行を選択すると、以下のようにジョブに関する追加情報を表示できます。

- [\[入力引数\]](#)

- 連続ログ
- [メトリクス] – 基本的なメトリクスのビジュアライゼーションを確認できます。含まれるメトリクスの詳細については、「[the section called “Spark ジョブ実行の Amazon CloudWatch メトリクスを表示”](#)」を参照してください。
- [Spark UI] – Spark UI でジョブの Spark ログを視覚化できます。Spark Web UI の使用の詳細については、「[the section called “Spark UI を使用したモニタリング”](#)」を参照してください。[the section called “ジョブ用の Spark UI の有効化”](#) のステップに従って、この機能を有効にします。

[詳細を表示] を選択すると、ジョブ実行の詳細ページに同じような情報が表示されます。あるいは、[モニタリング] ページからジョブ実行の詳細ページに移動することもできます。ナビゲーションペインで、[モニタリング] を選択します。[Job runs] (ジョブの実行) リストまで下にスクロールします。ジョブを選択し、[View run details] (実行の詳細を表示する) を選択します。内容は、[ジョブの実行の詳細を表示する](#) に記述されます。

このジョブのログについての詳細は、「[ジョブの実行ログの表示](#)」を参照してください。

ジョブスクリプトの表示

ジョブ内のすべてのノードの情報を入力した後、AWS Glue Studio により、スクリプトが生成されます。このスクリプトは、ジョブによるソースからのデータの読み取り、変換、ターゲットの場所への書き込みに使用されます。ジョブを保存すると、このスクリプトをいつでも表示できます。

ジョブ用に生成されたスクリプトを表示するには

1. ナビゲーションペインで [Jobs] (ジョブ) を選択します。
2. [Jobs] (ジョブ) ページの [Your Jobs] (ジョブ) リストで、レビューするジョブの名前を選択します。または、リストからジョブを選択し、[Actions] (アクション) メニューを選択して、[Edit job] (ジョブの編集) を選択することもできます。
3. ビジュアルエディタで、ページ上部の [Script] (スクリプト) タブを選択して、ジョブスクリプトを表示します。

ジョブスクリプトを編集する場合は、「[AWS Glue プログラミングガイド](#)」を参照してください。

ジョブのプロパティを変更する

ジョブで実行されるアクションは、ジョブ図のノードにより定義されますが、お客様がジョブに設定できるプロパティもいくつかあります。これらのプロパティでは、ジョブが実行される環境、使用するリソース、しきい値の設定、セキュリティ設定などを決定できます。

ジョブの実行環境をカスタマイズするには

1. ナビゲーションペインで [Jobs] (ジョブ) を選択します。
2. [Jobs] (ジョブ) ページの [Your Jobs] (ジョブ) リストで、レビューするジョブの名前を選択します。
3. ビジュアルエディタのページで、ジョブ編集ペインの上部にある [Job details] (ジョブの詳細) タブを選択します。
4. 必要に応じて、ジョブのプロパティを変更します。

ジョブのプロパティの詳細については、AWS Glue デベロッパーガイドの「[Defining Job Properties](#)」を参照してください。

5. 次のようなジョブのプロパティを指定する必要がある場合、[Advanced properties] (詳細プロパティ) セクションを展開します。
 - Script filename (スクリプトファイル名) – Amazon S3 でジョブスクリプトを保存するファイルの名前です。
 - Script path (スクリプトパス) – ジョブスクリプトが保存される Amazon S3 の場所です。
 - Job metrics (ジョブのメトリクス) – (Python シェルジョブでは使用不可) ジョブの実行時に Amazon CloudWatch メトリクスの作成をオンにします。
 - Continuous logging (継続的なログ記録) – (Python シェルジョブでは使用不可) CloudWatch への連続的なログ記録をオンにし、ジョブが完了する前にログを表示できるようにします。
 - Spark UI および Spark UI logs path (Spark UI ログパス) – (Python シェルジョブでは使用不可) このジョブをモニタリングするための Spark UI の使用をオンにし、Spark UI ログの場所を指定します。
 - Maximum concurrency (同時実行の最大数) – このジョブで許可される同時実行の最大数を設定します。
 - Temporary path (一時的なパス) – AWS Glue でジョブスクリプトが実行されるときに一時的な中間結果が書き込まれる Amazon S3 の作業ディレクトリの場所です。

- Delay notification threshold (minutes) (遅延通知のしきい値 (分)) – ジョブの遅延のしきい値を指定します。ジョブがしきい値で指定された時間よりも長い時間実行された場合、AWS Glue によりジョブの遅延通知が CloudWatch に送信されます。
 - Security configuration (セキュリティ設定) および Server-side encryption (サーバー側の暗号化) – これらのフィールドを使用して、ジョブの暗号化オプションを選択します。
 - Use Glue Data Catalog as the Hive metastore (Glue Data Catalog を Hive メタストアとして使用する) – Apache Hive Metastore の代わりに AWS Glue Data Catalog を使用する場合は、このオプションを選択します。
 - Additional network connection (追加のネットワーク接続) – VPC 内のデータソース向けに、タイプ Network の接続を指定して、ジョブが VPC 経由でデータにアクセスできるようにします。
 - Python library path (Python ライブラリパス)、Dependent jars path (依存 JARS パス) (Python シェルジョブでは使用不可)、Referenced files path (参照されるファイルパス) – これらのフィールドを使用して、スクリプトの実行時にジョブが使用する追加のファイルの場所を指定できます。
 - Job Parameters (ジョブパラメータ) – スクリプトに名前付きパラメータとして渡される一連のキーと値のペアを追加できます。AWS Glue API の Python 呼び出しでは、明示的に名前パラメータを渡すことが最善です。ジョブスクリプトでのパラメータの使用の詳細については、AWS Glue デベロッパーガイドの「[Passing and Accessing Python Parameters in AWS Glue](#)」を参照してください。
 - Tags (タグ) – ジョブにタグを追加することで、ジョブの整理および識別が容易になります。
6. ジョブのプロパティを変更した後、ジョブを保存します。

Spark のシャッフルファイルを Amazon S3 に保存する

一部の ETL ジョブでは、複数のパーティションから情報を読み込んで結合する必要があります。例えば、結合変換を使用する場合などです。この操作はシャッフリングと呼ばれます。シャッフル中、データはディスクに書き込まれ、ネットワーク経由で転送されます。AWS Glue バージョン 3.0 では、これらのファイルの保存場所として Amazon S3 を設定できます。AWS Glue では、Amazon S3 との間でシャッフルファイルの書き込みおよび読み込みを行うシャッフルマネージャーを利用できます。Amazon S3 からのシャッフルファイルの書き込みおよび読み込みは、ローカルディスク (または Amazon EC2 用に最適化された Amazon EBS) に比べて 5% ~ 20% 遅くなります。ただし、Amazon S3 には無制限のストレージ容量があるため、ジョブを実行する際に「No space left on device」エラーについて心配する必要はありません。

ファイルのシャッフルに Amazon S3 を使用するようにジョブを設定するには

1. [Jobs] (ジョブ) ページの [Your Jobs] (ジョブ) リストで、変更するジョブの名前を選択します。
2. ビジュアルエディタのページで、ジョブ編集ペインの上部にある [Job details] (ジョブの詳細) タブを選択します。

[Job parameters] (ジョブパラメータ) セクションまで下にスクロールします。

3. 次のキーと値のペアを指定します。

- `--write-shuffle-files-to-s3 — true`

これは、AWS Glue でシャッフルマネージャーを構成する主要なパラメータです。シャッフルデータの書き込みおよび読み取りには Amazon S3 バケットを使用します。デフォルトでは、このパラメータの値は `false` です。

- (オプション) `--write-shuffle-spills-to-s3 — true`

このパラメータを使用すると、流出ファイルを Amazon S3 バケットにオフロードできます。これにより、AWS Glue での Spark ジョブの耐障害性が強化されます。これは、大量のデータをディスクに退避させる大規模なワークロードにのみ必要です。デフォルトでは、このパラメータの値は `false` です。

- (オプション) `--conf spark.shuffle.glue.s3ShuffleBucket — S3://<shuffle-bucket>`

このパラメータを使用すると、シャッフルファイルを書き込む際に使用する Amazon S3 バケットを指定できます。このパラメータを設定しない場合、その場所は一時パス (`--TempDir`) 用に指定された場所にある `shuffle-data` フォルダになります。

Note

シャッフルバケットの場所がジョブが実行されるのと同じ AWS リージョンにあることを確認します。

また、ジョブの実行の終了後、ファイルはシャッフルサービスによりクリーンアップされないため、シャッフルバケットの場所で Amazon S3 ストレージライフサイクルポリシーを設定する必要があります。詳細については、Amazon S3 ユーザーガイドの「[ストレージのライフサイクルの管理](#)」をご参照ください。

ジョブの保存

ジョブが保存されるまで、[Save] (保存) ボタンの左側に、赤い [Job has not been saved] (ジョブは保存されていません) という吹き出しが表示されます。



ジョブを保存するには

1. [Visual] (ビジュアル) および [Job details] (ジョブの詳細) タブで、必要なすべての情報を入力します。
2. [保存] ボタンを選択します。

ジョブを保存すると、「保存されていません」という吹き出しが、ジョブが最後に保存された日時の表示に変わります。

ジョブを保存する前に AWS Glue Studio を終了した場合、次回の AWS Glue Studio の起動時に通知が表示されます。通知では、未保存のジョブがあり、それを復元するかどうかを尋ねられます。ジョブの復元を選択すると、そのジョブの編集を続行できます。

ジョブ保存時のエラーのトラブルシューティング

[Save] (保存) ボタンをクリックすると、ジョブに必要な情報が欠落している場合は、情報が欠落しているタブに赤い吹き出しが表示されます。吹き出しの数字は、欠落している検出されたフィールドの数を示します。

Untitled job 

Visual  | Script | Job details  | Runs | Schedules

- ビジュアルエディタのノードが正しく設定されていない場合、[Visual] (ビジュアル) タブに赤い吹き出しが表示され、エラーのあるノードに警告記号



が表示されます。

1. ノードを選択します。ノードの詳細パネルで、不足または誤った情報があるタブに赤い吹き出しが表示されます。

2. ノードの詳細パネルで赤い吹き出しが表示されるタブを選択し、ハイライトされている問題のフィールドを見つけます。フィールドの下のエラーメッセージには、問題に関する追加情報が表示されます。

The screenshot shows the AWS Glue console interface for an 'Untitled job'. At the top, there are buttons for 'Save' and 'Run', and a red notification 'Job has not been saved'. The main navigation bar includes 'Visual 2', 'Script', 'Job details 1', 'Runs', and 'Schedules'. Below this is a toolbar with icons for Source, Transform, Target, Undo, Redo, and Remove. The central workspace contains a single node labeled 'Data source - S3 bucket' with a red warning icon. To the right, the 'Node properties' panel is open to 'Data source properties - S3 2'. It shows the 'S3 source type' set to 'Data Catalog table'. The 'Database' dropdown menu is open, showing 'Choose a database' and a red warning message: 'Database is required.'. Below it, the 'Table' dropdown menu is also open, showing 'Choose a table' and a red warning message: 'Table is required.'. The 'Partition predicate - optional' field is empty.

- ジョブのプロパティに問題がある場合、[Job details] (ジョブの詳細) タブには、赤い吹き出しが表示されます。そのタブを選択し、ハイライトされている問題のフィールドを見つけます。フィールドの下のエラーメッセージにより、問題に関する追加情報が表示されます。

Untitled job Visual **2** | Script | **Job details 1** | Runs | SchedulesBasic properties [Info](#)

Name

Description - *optional*

Descriptions can be up to 2048 characters long.

IAM Role

Role assumed by the job with permission to access your data stores. Ensure that this role has permission to your Amazon S3 sources, targets, temporary directory, scripts, and any libraries used by the job.

 IAM Role is required.

Type

The type of ETL job. This is set automatically based on the types of data sources you have selected.

ジョブのクローンを作成する

[Clone job] (ジョブのクローン) アクションを使用して、既存のジョブを新しいジョブにコピーします。

既存のジョブをコピーして新しいジョブを作成するには

1. [Jobs] (ジョブ) ページの [Your Jobs] (ジョブ) リストで、複製するジョブを選択します。
2. [Actions] (アクション) メニューから、[Clone job] (ジョブのクローン) を選択します。
3. 新しいジョブの名前を入力します。その後、ジョブを保存または編集できます。

ジョブの削除

不要になったジョブを削除できます。1回の操作で1つ以上のジョブを削除できます。

AWS Glue Studio からジョブを削除する方法

1. [Jobs] (ジョブ) ページの [Your Jobs] (ジョブ) リストで、削除するジョブを選択します。
2. [Actions] (アクション) メニューから、[Delete job] (ジョブの削除) を選択します。
3. 「**delete**」を入力して、ジョブを削除することを確認します。

ビジュアルエディタで、そのジョブの [Job details] (ジョブの詳細) タブを表示して、保存したジョブを削除することもできます。

AWS Glue でのジョブの使用

以下のセクションでは、AWS Glueの ETL ジョブと Ray ジョブについて説明します。

トピック

- [AWS Glue バージョン](#)
- [での Spark ジョブの操作 AWS Glue](#)
- [AWS Glue での Ray ジョブの使用](#)
- [AWS Glue での Python シェルジョブに関するジョブプロパティの設定](#)
- [AWS Glue のモニタリング](#)
- [AWS Glue ジョブ実行ステータス](#)

AWS Glue バージョン

AWS Glue バージョンパラメータは、ジョブを追加または更新するときに設定できます。AWS Glue バージョンにより、AWS Glueがサポートする Apache Spark と Python のバージョンが決定されます。Python バージョンは、Spark タイプのジョブでサポートされるバージョンを示します。次の表に、使用可能な AWS Glue のバージョン、対応する Spark と Python のバージョン、および機能上のその他の変更点を示します。

AWS Glue バージョン

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
AWS Glue4.0	Spark 環境バージョン <ul style="list-style-type: none"> • Spark 3.3.0 • Python 3.10 	Java 8	AWS Glue 4.0 が AWS Glue の最新バージョンです。この AWS Glue リリースには、次のようないくつかの最適化とアップグレードが組み込まれています。

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
			<ul style="list-style-type: none"> • Spark 3.1 から Spark 3.3 になり、Spark の機能の多くがアップグレードされています。 • Pandas と組み合わせたときに改善されるいくつかの機能。詳細については、「What's New in Spark 3.3」(Spark 3.3 の新機能)を参照してください。 • Amazon EMR で開発されたその他の最適化。 • EMR ファイルシステム (EMRFS) 2.53 へのアップグレード。 • Log4j 1.x から Log4j 2 への移行 • Boto のアップグレード版など、いくつかの Python モジュールの AWS Glue 3.0 からの更新。

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
			<ul style="list-style-type: none"> • デフォルトの Amazon Redshift コネクタを含むいくつかのコネクタのアップグレード。 付録 C: コネクタのアップグレード を参照してください。 • いくつかの JDBC ドライバーのアップグレード。 付録 B: JDBC ドライバーのアップグレード を参照してください。 • 新しい Amazon Redshift コネクタと JDBC ドライバーで更新されました。 • Apache Hudi、Delta Lake、および Apache Iceberg によるオープンデータレイクフレームワークのネイティブサポート。 • Amazon S3 ベースのクラウドシャッフルストレージプラグイン (Apache Spark プラグイン)

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
			<p>ン) のネイティブサポート。これにより、Amazon S3 によるシャッフルと伸縮自在なストレージ容量が使用できるようになります。</p> <p>制限事項</p> <p>AWS Glue 4.0 での制限事項を次に示します。</p> <ul style="list-style-type: none"> • AWS Glue 機械学習変換と個人を特定できる情報 (PII) 変換は AWS Glue 4.0 ではまだ使用できません。 <p>AWS Glue バージョン 4.0 への移行に関する詳細については、「Spark ジョブの AWS Glue の AWS Glue バージョン 4.0 への移行」を参照してください。</p>

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
	Ray 環境バージョン <ul style="list-style-type: none"> • Ray 2.4.0 Python 3.9 	該当なし	<p>AWS Glue for Ray を使用して分散型 Python アプリケーションをビルドして実行します。</p> <ul style="list-style-type: none"> • Python 3.9 で Ray-2.4.0 データ分散 (<code>ray[data]</code>) をサポートします。この Ray リリースの詳細については、Ray リポジトリの Ray-2.4.0 を参照してください。GitHub • Ray2.4 ランタイム環境への追加の Python ライブラリのインストールをサポートします。詳細については、「the section called “Ray ジョブ用の Python モジュールを追加する”」を参照してください。 • Ray ジョブのログとメトリクスを Amazon CloudWatch と統合します。詳細については、

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
			<p>「the section called “Ray エラーをトラブルシューティングする”」および「the section called “Ray ジョブメトリクス”」を参照してください。</p> <ul style="list-style-type: none"> 各ジョブ実行ページにある AWS Glue Studio Ray ジョブのメトリクスを集約して視覚化します。 クラスター内の各作業ディレクトリへのファイルの配布、Ray オブジェクトストアから Amazon S3 へのオブジェクトのスピル、Ray ジョブに割り当てられるワーカーノードの最小数の制御をサポートします。詳細については、「the section called “Ray ジョブのパラメータ”」を参照してください。

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
			<p>AWS Glue 4.0 での Ray ジョブの制限事項</p> <ul style="list-style-type: none"> • AWS Glue Ray のインタラクティブセッションは、このリリースでもプレビュー段階にあります。 • AWS Glue フォーレイと Amazon VPC との統合は現在ご利用いただけません。の VPC AWS 内のリソースには、パブリックルートがないとアクセスできません。Amazon VPC AWS Glue との併用に関する詳細については、を参照してくださいthe section called “VPC エンドポイント (AWS PrivateLink)”。 • AWS Glue for Ray は、米国東部 (バージニア北部)、米国東部 (オハイオ)、米国西部 (オレゴン)、

AWS Glue バージョン	サポートされている ランタイム環境バー ジョン	サポートされている Java バージョン	機能の変更
			アジアパシフィッ ク (東京)、ヨーロッ パ (アイルランド) でご利用いただけ ます。

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
AWS Glue3.0	<ul style="list-style-type: none"> Spark 3.1.1 「Python 3.7」 	Java 8	<p>Spark エンジン 3.0 へのアップグレードに加えて、この AWS Glue リリースでは、以下の最適化とアップグレードが行われています。</p> <ul style="list-style-type: none"> Spark のメジャーリリースである Spark 3.0 に対する AWS Glue ETL ライブラリの構築。 AWS Glue 3.0 によるストリーミングジョブのサポート パフォーマンスと信頼性のための、以下の新しい AWS Glue Spark ランタイム最適化。 <ul style="list-style-type: none"> CSV データ読み取りのための、より高速化された Apache Arrow ベースのインメモリ列指向処理。 CSV データに関するベクトル化された読み取り

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
			<p>のための SIMD ベースの実行。</p> <ul style="list-style-type: none"> Spark のアップグレードには、Amazon EMR で開発された他の最適化も含まれます。 EMRFS 2.38 から 2.46 へのアップグレードとともに、Amazon S3 アクセスのための新機能追加とバグ修正。 新しい Spark バージョンに必要ないくつかの依存関係のアップグレード。付録 A: 注目すべき依存関係のアップグレード を参照してください。 ネイティブにサポートされているデータソース用の、JDBC ドライバのアップグレード。付録 B: JDBC ドライバのアップ

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
			<p>アップグレード を参照してください。</p> <p>制限事項</p> <p>以下に、AWS Glue 3.0での制限事項を示します。</p> <ul style="list-style-type: none">• 現在、AWS Glue の機械学習変換は AWS Glue 3.0 で使用できません。• Spark 2.4 に依存しており、Spark 3.1との互換性がない一部のカスタム Spark コネクタは、AWS Glue 3.0 では動作しません。 <p>AWS Glue バージョン 3.0 への移行に関する詳細については、「Spark ジョブの AWS Glue の AWS Glue バージョン 3.0 への移行」を参照してください。</p>

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
AWS Glue 2.0 (非推奨、サポート終了)	<ul style="list-style-type: none"> Spark 2.4.3 「Python 3.7」 	該当なし	<p>AWS Glue バージョン 2.0 では、AWS Glue バージョン 1.0 で提供される機能に加えて次の機能も提供されます。</p> <ul style="list-style-type: none"> AWS Glue で Apache Spark ETL ジョブを実行するための、アップグレードされたインフラストラクチャにより、起動時間が大幅に短縮されています。 デフォルトのログ記録機能はリアルタイムになり、ストリームがドライバー用とエグゼキューター用、出力用とエラー用の間でそれぞれ分離されています。 ジョブレベルで追加の Python モジュールまたは異なるバージョンを指定するためのサポート。

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
			<div data-bbox="1187 302 1507 1430" style="border: 1px solid #ccc; border-radius: 10px; padding: 10px; background-color: #e6f2ff;"> <p> Note</p> <p>基礎となるアーキテクチャの変更に伴い、AWS Glue バージョン 2.0 での一部の依存関係とバージョンが、AWS Glue バージョン 1.0 から変更されています。AWS Glue のメジャーバージョンリリース間での移行の前に、AWS Glue ジョブの検証を行ってください。</p> </div> <p>AWS Glue バージョン 2.0 の機能と制限事項の詳細については、「起動時間を短縮した Spark ETL ジョブの実行」 を参照してください。</p>

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
AWS Glue 1.0 (非推奨、サポート終了)	<ul style="list-style-type: none"> Spark 2.4.3 Python 2.7 Python 3.6 	該当なし	<p>AWS Glue ETL ジョブ (AWS Glue バージョン 1.0 を使用) には、Parquet 形式と ORC 形式のジョブのブックマークを維持できます。これまで AWS Glue ETL ジョブでは、JSON、CSV、Apache Avro、XML などの一般的な Amazon S3 ソース形式のみのブックマークが可能でした。</p> <p>ETL 入力および出力の形式オプションを設定する際に、Apache Avro のリーダー/ライター形式 1.8 を使用して Avro 論理型の読み取りと書き込みをサポートするように指定できます (AWS Glue バージョン 1.0 を使用)。以前は、バージョン 1.7 Avro リーダー/ライター形式のみがサポートされていました。</p>

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
			<p>DynamoDB 接続タイプは、(AWS Glue バージョン 1.0 を使用する) ライターオプションをサポートします。</p> <p>制限事項</p> <p>AWS Glue 1.0 での制限事項を次に示します。</p> <ul style="list-style-type: none"> • AWS Glue バージョン 0.9 と 1.0 は、今後、アジアパシフィック (ジャカルタ) (ap-southeast-3)、中東 (UAE) (me-central-1)、その他の新しい地域ではご利用いただけなくなります。

AWS Glue バージョン	サポートされているランタイム環境バージョン	サポートされている Java バージョン	機能の変更
AWS Glue 0.9 (非推奨、サポート終了)	<ul style="list-style-type: none"> Spark 2.2.1 Python 2.7 	該当なし	<p>AWS Glue バージョンを指定せずに作成されたジョブは、デフォルトで AWS Glue 0.9 に設定されます。</p> <p>制限事項</p> <p>AWS Glue 0.9 での制限事項を次に示します。</p> <ul style="list-style-type: none"> AWS Glue バージョン 0.9 と 1.0 は、今後、アジアパシフィック (ジャカルタ) (ap-southeast-3)、中東 (UAE) (me-central-1)、その他の新しい地域ではご利用いただけなくなります。

起動時間を短縮した Spark ETL ジョブの実行

AWS Glue バージョン 2.0 以降には、スタートアップ時間を短縮して AWS Glue で Apache Spark ETL (抽出、変換、ロード) ジョブを実行するためのアップグレード済みインフラストラクチャが用意されています。待機時間の短縮により、データエンジニアの生産性が向上し、AWS Glue とのやり取りが増加します。ジョブ開始時間の差異を低減することは、データが分析に利用できるようになる SLA の基準以上にすることに役立ちます。

AWS Glue ETL ジョブでこの特徴を使用するには、ジョブ作成時に **2.0** または Glue version の以降のバージョンを選択してください。

トピック

- [サポートされる新機能](#)
- [ログ記録の動作。](#)
- [サポートされていない機能](#)

サポートされる新機能

このセクションでは、AWS Glue バージョン 2.0 以降でサポートされる新機能について説明します。

ジョブレベルでの追加の Python モジュールの指定をサポート

AWS Glue バージョン 2.0 以降では、ジョブレベルで追加の Python モジュールや異なるバージョンを指定することもできます。--additional-python-modules オプションでコンマ区切りの Python モジュールのリストを指定することで、新しいモジュールを追加したり、既存のモジュールのバージョンを変更したりできます。

例えば、更新したり新しい scikit-learn モジュールを追加したりするには、"--additional-python-modules", "scikit-learn==0.21.3" のキー/値を使用します。

また、--additional-python-modules オプションの中で、Python ホイールモジュールへの Amazon S3 パスを指定できます。例:

```
--additional-python-modules s3://aws-glue-native-spark/tests/j4.2/ephem-3.7.7.1-cp37-cp37m-linux_x86_64.whl,s3://aws-glue-native-spark/tests/j4.2/fbprophet-0.6-py3-none-any.whl,scikit-learn==0.21.3
```

AWS Glue は Python パッケージインストーラ (pip3) を使用して追加のモジュールをインストールします。python-modules-installer-option で指定した追加オプションを pip3 に渡して、モジュールをインストールできます。pip3 からの非互換性や制限が適用されます。

AWS Glue バージョン 2.0 で既に提供されている Python モジュール

AWS Glue バージョン 2.0 では、すぐに使用できる次の Python モジュールがサポートされています。

- `setuptools-45.2.0`

- subprocess32–3.5.4
- ptvsd–4.3.2
- pydevd–1.9.0
- PyMySQL–0.9.3
- docutils–0.15.2
- jmespath–0.9.4
- six–1.14.0
- python_dateutil–2.8.1
- urllib3–1.25.8
- botocore–1.15.4
- s3transfer–0.3.3
- boto3–1.12.4
- certifi–2019.11.28
- chardet–3.0.4
- idna–2.9
- requests–2.23.0
- pyparsing–2.4.6
- enum34–1.1.9
- pytz–2019.3
- numpy–1.18.1
- cycloper–0.10.0
- kiwisolver–1.1.0
- scipy–1.4.1
- pandas–1.0.1
- pyarrow–0.16.0
- matplotlib–3.1.3
- pyhocon–0.3.54
- mpmath–1.1.0

- sympy-1.5.1
- patsy-0.5.1
- statsmodels-0.11.1
- fsspec-0.6.2
- s3fs-0.4.0
- Cython-0.29.15
- joblib-0.14.1
- pmdarima-1.5.3
- scikit-learn-0.22.1
- tbats-1.0.9

ログ記録の動作。

AWS Glue バージョン 2.0 では、さまざまなデフォルトのログ記録動作がサポートされています。相違点は次のとおりです。

- ログ記録はリアルタイムで行われます。
- ドライバーとエグゼキューターには別々のストリームがあります。
- ドライバーとエグゼキューターごとに、出カストリームとエラーストリームの 2 つのストリームがあります。

ドライバーとエグゼキューターのストリーム

ドライバーストリームは、ジョブ実行 ID で識別されます。エグゼキューターストリームは、ジョブ `<##ID>_<##### ID>` で識別されます。例:

- "logStreamName":
"jr_8255308b426fff1b4e09e00e0bd5612b1b4ec848d7884cebe61ed33a31789..._g-f65f617bd31d54bd94482af755b6cdf464542..."

出カストリームとエラーストリーム

出カストリームには、コードからの標準出力 (stdout) があります。エラーストリームには、コード/ライブラリからのログ記録メッセージがあります。

- ログストリーム:
 - ドライバーログストリームには `<jr>` があります。`<jr>` はジョブ実行 ID です。
 - エグゼキュターログストリームには `<jr>_<g>` があります。`<g>` はエグゼキュターのタスク ID です。ドライバーエラーログでエグゼキュタータスク ID を検索できます。

AWS Glue バージョン 2.0 のデフォルトのロググループは、以下のとおりです。

- 出力用の `/aws-glue/jobs/logs/output`
- エラー用の `/aws-glue/jobs/logs/error`

セキュリティ設定が指定されると、ロググループ名が次のように変更されます。

- `/aws-glue/jobs/<security configuration>-role/<Role Name>/output`
- `/aws-glue/jobs/<security configuration>-role/<Role Name>/error`

コンソールで、[Logs] (ログ) リンクが出力ロググループを指し、[Error] (エラー) リンクはエラーロググループを指します。連続ログ記録が有効な場合、[Logs] (ログ) リンクは連続ロググループを指し、[Output] (出力) リンクが出力ロググループを指します。

ログ記録ルール

Note

連続ログ記録のデフォルトのロググループ名は `/aws-glue/jobs/logs-v2` です。

AWS Glue バージョン 2.0 以降では、次のように連続ログ記録の動作は AWS Glue バージョン 1.0 の場合と同じです。

- デフォルトのロググループ: `/aws-glue/jobs/logs-v2`
- ドライバーログストリーム: `<jr>-driver`
- エグゼキュターログストリーム: `<jr>-<##### ID>`

ロググループ名は、`--continuous-log-logGroupName` の設定で変更できます。

ログストリーム名には、`--continuous-log-logStreamPrefix` の設定でプレフィックスを付加できます。

サポートされていない機能

以下の AWS Glue の機能はサポートされていません。

- 開発エンドポイント
- AWS Glue バージョン 2.0 以降は Apache YARN では実行されないため、YARN の設定は適用されません
- AWS Glue バージョン 2.0 以降は Hadoop Distributed File System (HDFS) がありません
- AWS Glue バージョン 2.0 以降では動的割り当てが使用されないため、ExecutorAllocationManager メトリクスは使用できません
- AWS Glue バージョン 2.0 以降のジョブでは、ワーカーの数とワーカータイプを指定しますが、maxCapacity は指定しません。
- AWS Glue バージョン 2.0 以降は s3n をサポートしていません。s3 または s3a を使用することをお勧めします。ジョブが何らかの理由で s3n を使用する必要がある場合、次の追加引数を渡すことができます。

```
--conf spark.hadoop.fs.s3n.impl=com.amazon.ws.emr.hadoop.fs.EmrFileSystem
```

Spark ジョブの AWS Glue の AWS Glue バージョン 3.0 への移行

このトピックでは、AWS Glue バージョン 0.9、1.0、2.0、および 3.0 の間の変更について説明し、Spark アプリケーションと ETL ジョブを AWS Glue 3.0 に移行できるようにします。

この機能を AWS Glue ETL ジョブに使用するには、ジョブを作成するとき、Glue version として **3.0** を選択します。

トピック

- [サポートされる新機能](#)
- [AWS Glue 3.0 に移行するためのアクション](#)
- [移行チェックリスト](#)
- [AWS Glue 0.9 から AWS Glue 3.0 への移行](#)
- [AWS Glue 1.0 から AWS Glue 3.0 への移行](#)
- [AWS Glue 2.0 から AWS Glue 3.0 への移行](#)
- [付録 A: 注目すべき依存関係のアップグレード](#)
- [付録 B: JDBC ドライバーのアップグレード](#)

サポートされる新機能

このセクションでは、AWS Glue バージョン 3.0 の新機能と利点について説明します。

- Apache Spark 3.1.1 に基づいており、オープンソースの Spark より最適化されています。また、AWS Glue および EMR サービス (アダプティブクエリ実行、ベクトル化されたリーダー、最適化されたシャッフルとパーティション結合など) によって開発されています。
- MySQL、Microsoft SQL Server、Oracle、PostgreSQL、MongoDB、アップグレードされた Spark ライブラリなど、すべての Glue ネイティブソース用に JDBC ドライバーをアップグレードしました。依存関係は Spark 3.1.1 によって取り込まれています。
- EMRFS をアップグレードして Amazon S3 アクセスを最適化し、デフォルトで Amazon S3 最適化出カコミッターを有効にしました。
- パーティションインデックス、プッシュダウン述語、パーティションリスト、アップグレードされた Hive メタストアクライアントにより、Data Catalog アクセスを最適化しました。
- セルレベルのフィルタリングとデータレイクランザクションを備えた管理されたカタログテーブル用の Lake Formation と統合しました。
- Spark 3.1.1 では、Spark エグゼキューターのメモリメトリクスと Spark 構造化ストリーミングメトリクスを含む新しい Spark UI が使用できます。
- スタートアップレイテンシーが短縮され、AWS Glue 2.0 と同様に、ジョブ全体の完了時間と対話性が改善されます。
- Spark ジョブは、1 秒単位で課金され、最小課金時間は 1/10 になります。例えば、AWS Glue 2.0 と同様、最小 10 分が最小 1 分になります。

AWS Glue 3.0 に移行するためのアクション

既存のジョブについては、ジョブ設定で、Glue version を以前のバージョンから Glue 3.0 に変更します。

- コンソールでは、Glue version で Spark 3.1, Python 3 (Glue Version 3.0) or Spark 3.1, Scala 2 (Glue Version 3.0) を選択します。
- AWS Glue Studio では、Glue version で Glue 3.0 - Supports spark 3.1, Scala 2, Python 3 を選択します。
- API では、[UpdateJob](#) API の GlueVersion パラメータで **3.0** を選択します。

新しいジョブについては、ジョブを作成するときに Glue 3.0 を選択します。

- コンソールでは、Glue version で Spark 3.1, Python 3 (Glue Version 3.0) or Spark 3.1, Scala 2 (Glue Version 3.0) を選択します。
- AWS Glue Studio では、Glue version で Glue 3.0 - Supports spark 3.1, Scala 2, Python 3 を選択します。
- API では、[CreateJob](#) API の GlueVersion パラメータで **3.0** を選択します。

AWS Glue 3.0 Spark イベントログを表示するには、[CloudFormation](#) または [Docker](#) を使用して [Glue 3.0 用にアップグレードされた Spark 履歴サーバーを起動します。](#)

移行チェックリスト

移行については、このチェックリストを確認してください。

- ジョブが HDFS に依存していますか。該当する場合は、HDFS を S3 に置き換えてみてください。
 - ジョブスクリプトコードで、DFS パスとして `hdfs://` または `/` で始まるファイルシステムのパスを検索します。
 - デフォルトのファイルシステムが HDFS で設定されていないかどうかを確認します。明示的に設定されている場合は、`fs.defaultFS` の設定を削除する必要があります。
 - ジョブに `dfs.*` のパラメータが含まれているかどうかを確認します。含まれている場合は、パラメータを無効にしても問題がないことを確認する必要があります。
- ジョブが YARN に依存していますか。該当する場合は、ジョブに次のパラメータが含まれているかどうかを確認して、影響を確認します。含まれている場合は、パラメータを無効にしても問題がないことを確認する必要があります。
 - `spark.yarn.*`

例:

```
spark.yarn.executor.memoryOverhead
spark.yarn.driver.memoryOverhead
spark.yarn.scheduler.reporterThread.maxFailures
```

- `yarn.*`

例:

```
yarn.scheduler.maximum-allocation-mb
```

```
yarn.nodemanager.resource.memory-mb
```

- ジョブが Spark 2.2.1 または Spark 2.4.3 に依存していますか。該当する場合は、Spark 3.1.1 で変更された機能をジョブで使用しているかどうかを確認して、影響を確認します。

- <https://spark.apache.org/docs/latest/sql-migration-guide.html#upgrading-from-spark-sql-22-to-23>

例えば、`percentile_approx` 関数、または既存の `SparkContext` が存在する場合、`SparkSession.builder.getOrCreate()` を使用した `SparkSession`。

- <https://spark.apache.org/docs/latest/sql-migration-guide.html#upgrading-from-spark-sql-23-to-24>

例えば、`array_contains` 関数、または `spark.sql.caseSensitive=true` の場合の `CURRENT_DATE`、`CURRENT_TIMESTAMP` 関数。

- ジョブの追加の jar に Glue 3.0 で競合がありますか。
 - AWS Glue 0.9/1.0 から: 既存の AWS Glue 0.9/1.0 ジョブに追加の jar が使用されていると、Glue 3.0 で利用可能なアップグレードまたは新しい依存関係により、クラスパスの競合が発生する可能性があります。--user-jars-first AWS Glue ジョブパラメータを使用するか、依存関係をシェーディングすることによって、AWS Glue 3.0 でのクラスパスの競合を避けることができます。
 - AWS Glue 2.0 から: この場合も、--user-jars-first AWS Glue ジョブパラメータを使用するか、依存関係をシェーディングすることによって、AWS Glue 3.0 でのクラスパスの競合を避けることができます。
- ジョブが Scala 2.11 に依存していますか。
 - AWS Glue 3.0 は Scala 2.12 を使用するため、ライブラリが Scala 2.11 に依存している場合は、Scala 2.12 でライブラリを再構築する必要があります。
- ジョブの外部 Python ライブラリが Python 2.7/3.6 に依存していますか。
 - Python ライブラリパスに egg/wheel/zip ファイルを設定するのではなく、--additional-python-modules パラメータを使用します。
 - Spark 3.1.1 では Python 2.7 のサポートが削除されたため、依存ライブラリを Python 2.7/3.6 から Python 3.7 に更新してください。

AWS Glue 0.9 から AWS Glue 3.0 への移行

移行時には、次の変更点に注意してください。

- AWS Glue 0.9 では、オープンソースの Spark 2.2.1 を使用していますが、AWS Glue 3.0 では、EMR で最適化された Spark 3.1.1 を使用しています。
- 一部の Spark の変更のみについてですが、削除された機能が参照されないように、スクリプトの修正が必要になる場合があります。
- 例えば、Spark 3.1.1 では Scala 型指定されていない UDF は有効になりませんが、Spark 2.2 では有効になります。
- AWS Glue 3.0 のすべてのジョブの実行におけるスタートアップ時間は、大幅に改善されています。スタートアップのレイテンシーが最大 10 分から最大 1 分になるため、Spark ジョブは 1 秒単位で課金され、最小課金時間は 1/10 になります。
- AWS Glue 2.0 以降、ログ記録の動作が変更されています。
- 依存関係の更新については、[付録 A: 注目すべき依存関係のアップグレード](#) を参照してください。
- Scala も 2.11 から 2.12 に更新されており、Scala 2.12 には Scala 2.11 との下位互換性がありません。
- また、AWS Glue 0.9 は Python 2 しか利用していませんでしたが、Python スクリプトで使用されるデフォルトのバージョンは Python 3.7 になりました。
 - Python 2.7 は、Spark 3.1.1 ではサポートされていません。
 - 追加の Python モジュールをインストールする新しいメカニズムが利用可能です。
- AWS Glue 3.0 は Apache YARN では実行されないため、YARN の設定は適用されません。
- AWS Glue 3.0 には、Hadoop Distributed File System (HDFS) はありません。
- 既存の AWS Glue 0.9 ジョブで使用されている追加の jar は、3.0 でいくつかの依存関係が 0.9 からアップグレードされたため、依存関係の競合を引き起こす可能性があります。--user-jars-first AWS Glue ジョブパラメータによって、AWS Glue 3.0 でのクラスパスの競合を避けることができます。
- AWS Glue 3.0 はまだ動的割り当てをサポートしていないため、ExecutorAllocationManager メトリクスは使用できません。
- AWS Glue バージョン 3.0 ジョブでは、ワーカーの数とワーカータイプを指定しますが、maxCapacity は指定しません。
- AWS Glue 3.0 はまだ機械学習変換をサポートしていません。
- AWS Glue 3.0 はまだ開発エンドポイントをサポートしていません。

Spark の移行ドキュメントを参照してください。

- 「[Upgrading from Spark SQL 2.2 to 2.3](#)」を参照してください。

- 「[Upgrading from Spark SQL 2.3 to 2.4](#)」を参照してください。
- 「[Upgrading from Spark SQL 2.4 to 3.0](#)」を参照してください。
- 「[Upgrading from Spark SQL 3.0 to 3.1](#)」を参照してください。
- 「[Changes in Datetime behavior to be expected since Spark 3.0.](#)」を参照してください。

AWS Glue 1.0 から AWS Glue 3.0 への移行

移行時には、次の変更点に注意してください。

- AWS Glue 1.0 はオープンソースの Spark 2.4 を使用していますが、AWS Glue 3.0 では、EMR で最適化された Spark 3.1.1 を使用しています。
 - 一部の Spark の変更のみについてですが、削除された機能が参照されないように、スクリプトの修正が必要になる場合があります。
 - 例えば、Spark 3.1.1 では Scala 型指定されていない UDF は有効になりませんが、Spark 2.4 では有効になります。
- AWS Glue 3.0 のすべてのジョブの実行におけるスタートアップ時間は、大幅に改善されています。スタートアップのレイテンシーが最大 10 分から最大 1 分になるため、Spark ジョブは 1 秒単位で課金され、最小課金時間は 1/10 になります。
- AWS Glue 2.0 以降、ログ記録の動作が変更されています。
- いくつかの依存関係の更新の要点については、次を参照してください。
- Scala も 2.11 から 2.12 に更新されており、Scala 2.12 には Scala 2.11 との下位互換性がありません。
- また、AWS Glue 0.9 は Python 2 しか利用していませんでしたが、Python スクリプトで使用されるデフォルトのバージョンは Python 3.7 になりました。
 - Python 2.7 は、Spark 3.1.1 ではサポートされていません。
 - 追加の Python モジュールをインストールする新しいメカニズムが利用可能です。
- AWS Glue 3.0 は Apache YARN では実行されないため、YARN の設定は適用されません。
- AWS Glue 3.0 には、Hadoop Distributed File System (HDFS) はありません。
- 既存の AWS Glue 1.0 ジョブで使用されている追加の jar は、3.0 でいくつかの依存関係が 1.0 からアップグレードされたため、依存関係の競合を引き起こす可能性があります。--user-jars-first AWS Glue ジョブパラメータによって、AWS Glue 3.0 でのクラスパスの競合を避けることができます。
- AWS Glue 3.0 はまだ動的割り当てをサポートしていないため、ExecutorAllocationManager メトリクスは使用できません。

- AWS Glue バージョン 3.0 ジョブでは、ワーカーの数とワーカータイプを指定しますが、maxCapacity は指定しません。
- AWS Glue 3.0 はまだ機械学習変換をサポートしていません。
- AWS Glue 3.0 はまだ開発エンドポイントをサポートしていません。

Spark の移行ドキュメントを参照してください。

- 「[Upgrading from Spark SQL 2.4 to 3.0](#)」を参照してください。
- 「[Changes in Datetime behavior to be expected since Spark 3.0.](#)」を参照してください。

AWS Glue 2.0 から AWS Glue 3.0 への移行

移行時には、次の変更点に注意してください。

- AWS Glue 2.0 に存在する既存のすべてのジョブパラメータと主要な機能は、AWS Glue 3.0 に存在します。
 - Parquet データを Amazon S3 に書き込むための EMRFS S3 最適化コミッターは、AWS Glue 3.0 ではデフォルトで有効になっています。ただし、`--enable-s3-parquet-optimized-committer` を `false` にすることで無効化することができます。
- AWS Glue 2.0 はオープンソースの Spark 2.4 を使用していますが、AWS Glue 3.0 では、EMR で最適化された Spark 3.1.1 を使用しています。
 - 一部の Spark の変更のみについてですが、削除された機能が参照されないように、スクリプトの修正が必要になる場合があります。
 - 例えば、Spark 3.1.1 では Scala 型指定されていない UDF は有効になりませんが、Spark 2.4 では有効になります。
- また、AWS Glue 3.0 の特徴として、EMRFS への更新や JDBC ドライバーの更新が行われ、AWS Glue による Spark 自体への追加の最適化も組み込まれています。
- AWS Glue 3.0 のすべてのジョブの実行におけるスタートアップ時間は、大幅に改善されています。スタートアップのレイテンシーが最大 10 分から最大 1 分になるため、Spark ジョブは 1 秒単位で課金され、最小課金時間は 1/10 になります。
- Python 2.7 は、Spark 3.1.1 ではサポートされていません。
- 依存関係の更新については、[付録 A: 注目すべき依存関係のアップグレード](#) を参照してください。
- Scala も 2.11 から 2.12 に更新されており、Scala 2.12 には Scala 2.11 との下位互換性がありません。

- 既存の AWS Glue 2.0 ジョブで使用されている追加の jar は、3.0 でいくつかの依存関係が 2.0 からアップグレードされたため、依存関係の競合を引き起こす可能性があります。--user-jars-first AWS Glue ジョブパラメータによって、AWS Glue 3.0 でのクラスパスの競合を避けることができます。
- AWS Glue 3.0 は AWS Glue 2.0 と比較すると、ドライバー/エグゼキューターの設定に対して Spark タスクの並列処理が異なり、パフォーマンス向上と利用可能なリソースをより有効に活用できます。spark.driver.cores と spark.executor.cores の両方とも、AWS Glue 3.0 のコア数 (スタンダードで 4、および G.1X ワーカー、G.2X ワーカーで 8) に設定されています。これらの構成は、AWS Glue ジョブのためのワーカータイプやハードウェアを変更するものではありません。これらの構成を利用して、Spark アプリケーションの Spark タスク並列性と一致するパーティションまたは分割の数を計算できます。

一般的に、ジョブのパフォーマンスは AWS Glue 2.0 と同等かそれ以上になります。ジョブの実行速度が遅い場合は、次のジョブ引数を渡すことでタスクの並列処理を増やすことができます。

- キー: --executor-cores 値: <#####>
- この値は、ワーカータイプの vCPU の数の 2 倍を超えないようにします。G.1X では 8、G.2X では 16、G.4X では 32、G.8X では 64 です。並列処理が増えるとメモリとディスクに負荷がかかり、ソースシステムとターゲットシステムがスロットリングされ、ジョブのパフォーマンスに影響する可能性があるため、この構成を更新する場合は注意が必要です。
- AWS Glue 3.0 は Spark 3.1 を使用しており、パーケットファイルから、またはパーケットファイルへのタイムスタンプの読み込み/保存に、動作が変更されています。詳細は「[Upgrading from Spark SQL 3.0 to 3.1](#)」を参照してください。

タイムスタンプ列を含むパーケットデータを読み書きする場合は、次のパラメータを設定することをお勧めします。これらのパラメータを設定すると、Spark 2 から Spark 3 へのアップグレード中に発生するカレンダーの非互換性の問題を、AWS Glue ダイナミックフレームとパークデータフレームの両方で解決することができます。CORRECTED オプションは、datetime 値をそのまま読み込むために使用し、LEGACY オプションは読み込む際にカレンダーの違いを考慮して datetime 値をリベースするために使用します。

```
- Key: --conf
- Value: spark.sql.legacy.parquet.int96RebaseModeInRead=[CORRECTED|LEGACY] --conf spark.sql.legacy.parquet.int96RebaseModeInWrite=[CORRECTED|LEGACY] --conf spark.sql.legacy.parquet.datetimeRebaseModeInRead=[CORRECTED|LEGACY]
```

Spark の移行ドキュメントを参照してください。

- 「[Upgrading from Spark SQL 2.4 to 3.0](#)」を参照してください。
- 「[Changes in Datetime behavior to be expected since Spark 3.0.](#)」を参照してください。

付録 A: 注目すべき依存関係のアップグレード

依存関係のアップグレードは次のとおりです。

依存関係	AWS Glue 0.9 でのバージョン	AWS Glue 1.0 でのバージョン	AWS Glue 2.0 でのバージョン	AWS Glue 3.0 でのバージョン
Spark	2.2.1	2.4.3	2.4.3	3.1.1-amzn-0
Hadoop	2.7.3-amzn-6	2.8.5-amzn-1	2.8.5-amzn-5	3.2.1-amzn-3
Scala	2.11	2.11	2.11	2.12
Jackson	2.7.x	2.7.x	2.7.x	2.10.x
[Hive]	1.2	1.2	1.2	2.3.7-amzn-4
EMRFS	2.20.0	2.30.0	2.38.0	2.46.0
Json4s	3.2.x	3.5.x	3.5.x	3.6.6
Arrow	該当なし	0.10.0	0.10.0	2.0.0
AWS Glue カタログクライアント	該当なし	該当なし	1.10.0	3.0.0

付録 B: JDBC ドライバーのアップグレード

JDBC ドライバーのアップグレードは次のとおりです。

ドライバー	過去の AWS Glue バージョンでの JDBC ドライバーのバージョン	AWS Glue 3.0 での JDBC ドライバーのバージョン
MySQL	5.1	8.0.23
Microsoft SQL Server	6.1.0	7.0.0
Oracle Database	11.2	21.1
PostgreSQL	42.1.0	42.2.18
MongoDB	2.0.0	4.0.0

Spark ジョブの AWS Glue の AWS Glue バージョン 4.0 への移行

このトピックでは、AWS Glue バージョン 0.9、1.0、2.0、および 3.0 の間の変更について説明し、Spark アプリケーションと ETL ジョブを AWS Glue 4.0 に移行できるようにします。また、AWS Glue 4.0 の機能とそれを使用する利点についても説明します。

この機能を AWS Glue ETL ジョブに使用するには、ジョブを作成するとき、Glue version として **4.0** を選択します。

トピック

- [サポートされる新機能](#)
- [AWS Glue 4.0 に移行するためのアクション](#)
- [移行チェックリスト](#)
- [AWS Glue 3.0 から AWS Glue 4.0 への移行](#)
- [AWS Glue 2.0 から AWS Glue 4.0 への移行](#)
- [AWS Glue 1.0 から AWS Glue 4.0 への移行](#)
- [AWS Glue 0.9 から AWS Glue 4.0 への移行](#)
- [AWS Glue 4.0 向けのコネクタと JDBC ドライバーの移行](#)
- [付録 A: 注目すべき依存関係のアップグレード](#)
- [付録 B: JDBC ドライバーのアップグレード](#)
- [付録 C: コネクタのアップグレード](#)

サポートされる新機能

このセクションでは、AWS Glue バージョン 4.0 の新機能と利点について説明します。

- これは Apache Spark 3.3.0 に基づいていますが、アダプティブクエリ実行、ベクトル化されたリーダー、最適化されたシャッフルとパーティション結合など、AWS Glue での最適化や Amazon EMR が含まれています。
- MySQL、Microsoft SQL Server、Oracle、PostgreSQL、MongoDB、アップグレードされた Spark ライブラリなど、すべての AWS Glue ネイティブソース用に JDBC ドライバーをアップグレードしました。依存関係は Spark 3.3.0 によって取り込まれています。
- 新しい Amazon Redshift コネクタと JDBC ドライバーで更新されました。
- EMR File System (EMRFS) をアップグレードして Amazon S3 アクセスを最適化し、デフォルトで Amazon S3 に最適化された出力コミッターが有効になりました。
- パーティションインデックス、プッシュダウン述語、パーティションリスト、アップグレードされた Hive メタストアクライアントにより、データカタログアクセスを最適化しました。
- セルレベルのフィルタリングとデータレイクトランザクションを備えた管理されたカタログテーブル用の Lake Formation と統合しました。
- スタートアップレイテンシーが短縮され、ジョブ全体の完了時間と対話性が改善されます。
- Spark ジョブは、1 秒単位で課金され、最小課金時間は 1/10 になります。つまり、最小 10 分が最小 1 分になります。
- Apache Hudi、Delta Lake、および Apache Iceberg によるオープンデータレイクフレームワークのネイティブサポート。
- Amazon S3 ベースのクラウドシャッフルストレージプラグイン (Apache Spark プラグイン) のネイティブサポート。これにより、Amazon S3 によるシャッフルと伸縮自在なストレージ容量が使用できるようになります。

Spark 3.1.1 から Spark 3.3.0 への主な機能強化

次の機能強化があります。

- 行レベルのランタイムフィルタリング ([SPARK-32268](#))。
- ANSI 拡張機能 ([SPARK-38860](#))。
- エラーメッセージの改善 ([SPARK-38781](#))。
- Parquet ベクトル化リーダーの複合型のサポート ([SPARK-34863](#))。
- Spark SQL の隠しファイルメタデータのサポート ([SPARK-37273](#))。

- Python/Pandas UDF 用のプロファイラーを提供 ([SPARK-37443](#))。
- Trigger.Once のようなストリーミングクエリを複数のバッチで実行できるように、Trigger.AvailableNow を導入 ([SPARK-36533](#))。
- より包括的な Datasource V2 プッシュダウン機能 ([SPARK-38788](#))。
- log4j 1 から log4j 2 への移行 ([SPARK-37814](#))。

その他の注目すべき変更

次の変更があります。

- **重要な変更**
 - Python 3.6 サポートへの参照を docs と Python/docs から削除 ([SPARK-36977](#))。
 - 組み込みの pickle を cloudpickle に置き換えることで、名前付きタプルハックを削除 ([SPARK-32079](#))。
 - Pandas の最小バージョンを 1.0.5 に引き上げ ([SPARK-37465](#))。

AWS Glue 4.0 に移行するためのアクション

既存のジョブについては、ジョブ設定で、Glue version を以前のバージョンから Glue 4.0 に変更します。

- AWS Glue Studio では、Glue version で Glue 4.0 - Supports Spark 3.3, Scala 2, Python 3 を選択します。
- API では、[UpdateJob](#) API オペレーションの GlueVersion パラメータで **4.0** を選択します。

新しいジョブについては、ジョブを作成するときに Glue 4.0 を選択します。

- コンソールでは、Glue version で Spark 3.3, Python 3 (Glue Version 4.0) or Spark 3.3, Scala 2 (Glue Version 3.0) を選択します。
- AWS Glue Studio では、Glue version で Glue 4.0 - Supports Spark 3.3, Scala 2, Python 3 を選択します。
- API では、[CreateJob](#) API オペレーションの GlueVersion パラメータで **4.0** を選択します。

AWS Glue 2.0 以前で生成された AWS Glue 4.0 Spark イベントログを表示するには、[AWS CloudFormation](#) または [Docker](#) を使用して AWS Glue 4.0 用にアップグレードされた Spark 履歴サーバーを起動します。

移行チェックリスト

移行については、次のチェックリストを確認してください。

Note

AWS Glue 3.0 に関連するチェックリスト項目については、「[移行チェックリスト](#)」を参照してください。

- ジョブの外部 Python ライブラリが Python 2.7/3.6 に依存していますか。
- Spark 3.3.0 では Python 2.7 と 3.6 のサポートが完全に削除されたため、依存ライブラリを Python 2.7/3.6 から Python 3.10 に更新してください。

AWS Glue 3.0 から AWS Glue 4.0 への移行

移行時には、次の変更点に注意してください。

- AWS Glue 3.0 に存在する既存のすべてのジョブパラメータと主要な機能は、AWS Glue 4.0 に存在します。
- AWS Glue 3.0 では、Amazon EMR に最適化された Spark 3.1.1 を使用していますが、AWS Glue 4.0 では、Amazon EMR に最適化された Spark 3.3.0 を使用しています。

一部の Spark の変更のみについてですが、削除された機能が参照されないように、スクリプトの修正が必要になる場合があります。

- AWS Glue 4.0 には EMRFS と Hadoop のアップデートも含まれています。特定のバージョンについては、「[付録 A: 注目すべき依存関係のアップグレード](#)」を参照してください。
- ETL ジョブで提供される AWS SDK が 1.11 から 1.12 にアップグレードされました。
- すべての Python ジョブは Python バージョン 3.10 を使用します。以前、AWS Glue 3.0 では、Python 3.7 が使用されていました。

その結果、AWS Glue に用意されている pymodule の一部はアップグレードされます。

- Log4j は Log4j2 にアップグレードされました。

- Log4j2 マイグレーションパスについては、[Log4j のドキュメント](#)を参照してください。
- カスタム log4j.properties ファイルの名前を log4j2.properties ファイルに変更し、適切な log4j2 プロパティを使用する必要があります。
- 特定のコネクタの移行については、「[AWS Glue 4.0 向けのコネクタと JDBC ドライバーの移行](#)」を参照してください。
- AWS Encryption SDK が 1.x から 2.x にアップグレードされました。AWS Glue セキュリティ設定を使用する AWS Glue ジョブや、AWS Encryption SDK の実行時に提供される依存関係に依存するジョブが影響を受けます。AWS Glue ジョブの移行については、次の説明を参照してください。

AWS Glue 2.0/3.0 には既に AWS Encryption SDK ブリッジバージョンが含まれているため、AWS Glue 2.0/3.0 ジョブを AWS Glue 4.0 ジョブに安全にアップグレードできます。

Spark の移行ドキュメントを参照してください。

- 「[Upgrading from Spark SQL 3.1 to 3.2](#)」(Spark SQL 3.1 から 3.2 へのアップグレード)
- 「[Upgrading from Spark SQL 3.2 to 3.3](#)」(Spark SQL 3.2 から 3.3 へのアップグレード)

AWS Glue 2.0 から AWS Glue 4.0 への移行

移行時には、次の変更点に注意してください。

Note

AWS Glue 3.0 に関連する移行手順については、「[AWS Glue 3.0 から AWS Glue 4.0 への移行](#)」を参照してください。

- AWS Glue 2.0 に存在する既存のすべてのジョブパラメータと主要な機能は、AWS Glue 4.0 に存在します。
- Parquet データを Amazon S3 に書き込むための EMRFS S3 最適化コミッターは、AWS Glue 3.0 以降デフォルトで有効になっています。ただし、`--enable-s3-parquet-optimized-commmitter` を `false` にすることで無効化することができます。
- AWS Glue 2.0 では、オープンソースの Spark 2.4 を使用していますが、AWS Glue 4.0 では、Amazon EMR で最適化された Spark 3.3.0 を使用しています。
- 一部の Spark の変更のみについてですが、削除された機能が参照されないように、スクリプトの修正が必要になる場合があります。

- 例えば、Spark 3.3.0 では Scala 型指定されていない UDF は有効になりませんが、Spark 2.4 では有効になります。
- ETL ジョブで提供される AWS SDK が 1.11 から 1.12 にアップグレードされました。
- また、AWS Glue 4.0 の特徴として、EMRFS への更新や JDBC ドライバーの更新が行われ、AWS Glue により Spark 自体に追加された最適化も組み込まれています。
- Scala は 2.11 から 2.12 に更新されており、Scala 2.12 には Scala 2.11 との下位互換性がありません。
- AWS Glue 2.0 は Python 3.7 と 2.7 しか利用していませんでしたが、Python スクリプトで使用されるデフォルトのバージョンは Python 3.10 になりました。
 - Python 2.7 は、Spark 3.3.0 ではサポートされていません。ジョブ設定で Python 2 をリクエストするジョブはすべて、`IllegalArgumentException` が発生して失敗します。
 - AWS Glue 2.0 以降、追加の Python モジュールをインストールする新しいメカニズムが利用可能です。
- 依存関係の更新については、[付録 A: 注目すべき依存関係のアップグレード](#) を参照してください。
- 既存の AWS Glue 2.0 ジョブで使用されている追加の JAR ファイルは、4.0 で一部の依存関係が 2.0 からアップグレードされたため、依存関係の競合を引き起こす可能性があります。--user-jars-first AWS Glue ジョブパラメータによって、AWS Glue 4.0 でのクラスパスの競合を避けることができます。
- AWS Glue 4.0 は Spark 3.3 を使用します。Spark 3.1 以降、parquet ファイルから、または parquet ファイルへのタイムスタンプのロード/保存の動作が変更されています。詳細は「[Upgrading from Spark SQL 3.0 to 3.1](#)」を参照してください。

タイムスタンプ列を含むパーケットデータを読み書きする場合は、次のパラメータを設定することをお勧めします。これらのパラメータを設定すると、Spark 2 から Spark 3 へのアップグレード中に発生するカレンダーの非互換性の問題を、AWS Glue ダイナミックフレームとパークデータフレームの両方で解決することができます。CORRECTED オプションは、datetime 値をそのまま読み込むために使用し、LEGACY オプションは読み込む際にカレンダーの違いを考慮して datetime 値をリベースするために使用します。

```
- Key: --conf
- Value: spark.sql.legacy.parquet.int96RebaseModeInRead=[CORRECTED|LEGACY] --
conf spark.sql.legacy.parquet.int96RebaseModeInWrite=[CORRECTED|LEGACY] --conf
spark.sql.legacy.parquet.datetimeRebaseModeInRead=[CORRECTED|LEGACY]
```

- 特定のコネクタの移行については、「[AWS Glue 4.0 向けのコネクタと JDBC ドライバーの移行](#)」を参照してください。

- AWS Encryption SDK が 1.x から 2.x にアップグレードされました。AWS Glue セキュリティ設定を使用する AWS Glue ジョブや、AWS Encryption SDK の実行時に提供される依存関係に依存するジョブが影響を受けます。AWS Glue ジョブの移行については、次の説明を参照してください。
- AWS Glue 2.0 には既に AWS Encryption SDK ブリッジバージョンが含まれているため、AWS Glue 2.0 ジョブを AWS Glue 4.0 ジョブに安全にアップグレードできます。

Spark の移行ドキュメントを参照してください。

- 「[Upgrading from Spark SQL 2.4 to 3.0](#)」 (Spark SQL 2.4 から 3.0 へのアップグレード)
- 「[Upgrading from Spark SQL 3.1 to 3.2](#)」 (Spark SQL 3.1 から 3.2 へのアップグレード)
- 「[Upgrading from Spark SQL 3.2 to 3.3](#)」 (Spark SQL 3.2 から 3.3 へのアップグレード)
- 「[Changes in Datetime behavior to be expected since Spark 3.0](#)」 (Spark 3.0 以降で予想される日時に関する動作の変化)

AWS Glue 1.0 から AWS Glue 4.0 への移行

移行時には、次の変更点に注意してください。

- AWS Glue 1.0 では、オープンソースの Spark 2.4 を使用していますが、AWS Glue 4.0 では、Amazon EMR で最適化された Spark 3.3.0 を使用しています。
 - 一部の Spark の変更のみについてですが、削除された機能が参照されないように、スクリプトの修正が必要になる場合があります。
 - 例えば、Spark 3.3.0 では Scala 型指定されていない UDF は有効になりませんが、Spark 2.4 では有効になります。
- AWS Glue 4.0 のすべてのジョブの実行におけるスタートアップ時間は、大幅に改善されています。スタートアップのレイテンシーが最大 10 分から最大 1 分になるため、Spark ジョブは 1 秒単位で課金され、最小課金時間は 1/10 になります。
- AWS Glue 4.0 でログ記録の動作は大幅に変更されました。Spark 3.3.0 では最低でも Log4j2 が必須です。
 - 一部の依存関係の更新があり、付録に記載されています。
- Scala も 2.11 から 2.12 に更新されており、Scala 2.12 には Scala 2.11 との下位互換性はありません。
- また、AWS Glue 0.9 は Python 2 しか使用していませんでしたが、Python スクリプトで使用されるデフォルトのバージョンは Python 3.10 になりました。

Python 2.7 は、Spark 3.3.0 ではサポートされていません。ジョブ設定で Python 2 をリクエストするジョブはすべて、`IllegalArgumentExcpion` が発生して失敗します。

- AWS Glue 2.0 以降、追加の Python モジュールを pip を通じてインストールする新しいメカニズムが利用可能です。詳細については、「[AWS Glue 2.0+ の pip を使用した追加 Python モジュールのインストール](#)」を参照してください。
- AWS Glue 4.0 は Apache YARN では実行されないため、YARN の設定は適用されません。
- AWS Glue 4.0 には、Hadoop 分散ファイルシステム (HDFS) はありません。
- 既存の AWS Glue 1.0 ジョブで使用されている追加の JAR ファイルは、4.0 で一部の依存関係が 1.0 からアップグレードされたため、依存関係の競合を引き起こす可能性があります。この問題を避けるため、デフォルトで `--user-jars-first` AWS Glue ジョブパラメータで AWS Glue 4.0 を有効にしています。
- AWS Glue 4.0 は Auto Scaling をサポートしています。そのため、Auto Scaling が有効になっていると、`ExecutorAllocationManager` メトリクスが使用可能になります。
- AWS Glue バージョン 4.0 ジョブでは、ワーカーの数とワーカータイプを指定しますが、`maxCapacity` は指定しません。
- AWS Glue 4.0 はまだ機械学習変換をサポートしていません。
- 特定のコネクタの移行については、「[AWS Glue 4.0 向けのコネクタと JDBC ドライバーの移行](#)」を参照してください。
- AWS Encryption SDK が 1.x から 2.x にアップグレードされました。AWS Glue セキュリティ設定を使用する AWS Glue ジョブや、AWS Encryption SDK の実行時に提供される依存関係に依存するジョブが影響を受けます。AWS Glue ジョブの移行については、次の説明を参照してください。
- AWS Glue 0.9/1.0 ジョブを AWS Glue 4.0 ジョブに直接移行することはできません。バージョン 2.x 以降に直接アップグレードしてすべての新機能をすぐに有効にすると、AWS Encryption SDK では、AWS Encryption SDK の古いバージョンで暗号化された暗号文の復号ができないからです。
- 安全にアップグレードするには、AWS Encryption SDK ブリッジバージョンが含まれている AWS Glue 2.0/3.0 ジョブにまず移行することをお勧めします。ジョブを 1 回実行して AWS Encryption SDK ブリッジバージョンを利用します。
- 完了したら、AWS Glue 2.0/3.0 ジョブを AWS Glue 4.0 に安全に移行できます。

Spark の移行ドキュメントを参照してください。

- 「[Upgrading from Spark SQL 2.4 to 3.0](#)」(Spark SQL 2.4 から 3.0 へのアップグレード)

- 「[Upgrading from Spark SQL 3.0 to 3.1](#)」 (Spark SQL 3.0 から 3.1 へのアップグレード)
- 「[Upgrading from Spark SQL 3.1 to 3.2](#)」 (Spark SQL 3.1 から 3.2 へのアップグレード)
- 「[Upgrading from Spark SQL 3.2 to 3.3](#)」 (Spark SQL 3.2 から 3.3 へのアップグレード)
- 「[Changes in Datetime behavior to be expected since Spark 3.0](#)」 (Spark 3.0 以降で予想される日時に関する動作の変化)

AWS Glue 0.9 から AWS Glue 4.0 への移行

移行時には、次の変更点に注意してください。

- AWS Glue 0.9 では、オープンソースの Spark 2.2.1 を使用していますが、AWS Glue 4.0 では、Amazon EMR で最適化された Spark 3.3.0 を使用しています。
 - 一部の Spark の変更のみについてですが、削除された機能が参照されないように、スクリプトの修正が必要になる場合があります。
 - 例えば、Spark 3.3.0 では Scala 型指定されていない UDF は有効になりませんが、Spark 2.2 では有効になります。
- AWS Glue 4.0 のすべてのジョブの実行におけるスタートアップ時間は、大幅に改善されています。スタートアップのレイテンシーが最大 10 分から最大 1 分になるため、Spark ジョブは 1 秒単位で課金され、最小課金時間は 1/10 になります。
- AWS Glue 4.0 以降、ログ記録の動作は大幅に変更されました。ここ (<https://spark.apache.org/docs/latest/core-migration-guide.html#upgrading-from-core-32-to-33>) に記載されているように、Spark 3.3.0 では最低でも Log4j2 が必須です。
- 一部の依存関係の更新があり、付録に記載されています。
- Scala も 2.11 から 2.12 に更新されており、Scala 2.12 には Scala 2.11 との下位互換性がありません。
- また、AWS Glue 0.9 は Python 2 しか使用していませんでしたが、Python スクリプトで使用されるデフォルトのバージョンは Python 3.10 になりました。
 - Python 2.7 は、Spark 3.3.0 ではサポートされていません。ジョブ設定で Python 2 をリクエストするジョブはすべて、IllegalArgumentException が発生して失敗します。
 - 追加の Python モジュールを pip を通じてインストールする新しいメカニズムが利用可能です。
- AWS Glue 4.0 は Apache YARN では実行されないため、YARN の設定は適用されません。
- AWS Glue 4.0 には、Hadoop 分散ファイルシステム (HDFS) はありません。
- 既存の AWS Glue 0.9 ジョブで使用されている追加の JAR ファイルは、3.0 で一部の依存関係が 0.9 からアップグレードされたため、依存関係の競合を引き起こす可能性があります。--user-

jars-first AWS Glue ジョブパラメータによって、AWS Glue 3.0 でのクラスパスの競合を避けることができます。

- AWS Glue 4.0 は Auto Scaling をサポートしています。そのため、Auto Scaling が有効になっていると、ExecutorAllocationManager メトリクスが使用可能になります。
- AWS Glue バージョン 4.0 ジョブでは、ワーカーの数とワーカータイプを指定しますが、maxCapacity は指定しません。
- AWS Glue 4.0 はまだ機械学習変換をサポートしていません。
- 特定のコネクタの移行については、「[AWS Glue 4.0 向けのコネクタと JDBC ドライバーの移行](#)」を参照してください。
- AWS Encryption SDK が 1.x から 2.x にアップグレードされました。AWS Glue セキュリティ設定を使用する AWS Glue ジョブや、AWS Encryption SDK の実行時に提供される依存関係に依存するジョブが影響を受けます。AWS Glue ジョブの移行については、次の説明を参照してください。
- AWS Glue 0.9/1.0 ジョブを AWS Glue 4.0 ジョブに直接移行することはできません。バージョン 2.x 以降に直接アップグレードしてすべての新機能をすぐに有効にすると、AWS Encryption SDK では、AWS Encryption SDK の古いバージョンで暗号化された暗号文の復号ができないからです。
- 安全にアップグレードするには、AWS Encryption SDK ブリッジバージョンが含まれている AWS Glue 2.0/3.0 ジョブにまず移行することをお勧めします。ジョブを 1 回実行して AWS Encryption SDK ブリッジバージョンを利用します。
- 完了したら、AWS Glue 2.0/3.0 ジョブを AWS Glue 4.0 に安全に移行できます。

Spark の移行ドキュメントを参照してください。

- 「[Upgrading from Spark SQL 2.2 to 2.3](#)」 (Spark SQL 2.2 から 2.3 へのアップグレード)
- 「[Upgrading from Spark SQL 2.3 to 2.4](#)」 (Spark SQL 2.3 から 2.4 へのアップグレード)
- 「[Upgrading from Spark SQL 2.4 to 3.0](#)」 (Spark SQL 2.4 から 3.0 へのアップグレード)
- 「[Upgrading from Spark SQL 3.0 to 3.1](#)」 (Spark SQL 3.0 から 3.1 へのアップグレード)
- 「[Upgrading from Spark SQL 3.1 to 3.2](#)」 (Spark SQL 3.1 から 3.2 へのアップグレード)
- 「[Upgrading from Spark SQL 3.2 to 3.3](#)」 (Spark SQL 3.2 から 3.3 へのアップグレード)
- 「[Changes in Datetime behavior to be expected since Spark 3.0](#)」 (Spark 3.0 以降で予想される日時に関する動作の変化)

AWS Glue 4.0 向けのコネクタと JDBC ドライバーの移行

アップグレードされた JDBC コネクタとデータレイクコネクタのバージョンについては、次を参照してください。

- [付録 B: JDBC ドライバーのアップグレード](#)
- [付録 C: コネクタのアップグレード](#)

Hudi

- Spark SQL サポートの強化:
 - Call Procedure コマンドにより、アップグレード、ダウングレード、ブートストラップ、クリーニング、および修復のサポートが追加されます。Spark SQL で Create/Drop/Show/Refresh Index 構文が可能です。
 - Spark DataSource を経由して使用する場合、Spark SQL に対してパフォーマンスのギャップがありました。過去には、データソースの書き込みは、SQL よりも高速でした。
 - 組み込みキージェネレーターはすべて、より高性能な Spark 固有の API 操作を実装しています。
 - SerDe の使用コストを削減するため、一括 insert 操作の UDF 変換を RDD 変換に置き換えました。
 - Hudi の Spark SQL では、SQL ステートメントの tblproperties またはオプションで指定される primaryKey が必要です。更新と削除のオペレーションには、preCombineField も必要です。
- バージョン 0.10.0 より前に primaryKey なしで作成された Hudi テーブルは、バージョン 0.10.0 以降、primaryKey フィールドつきで再作成する必要があります。

PostgreSQL

- いくつかの脆弱性 (CVE) に対応しました。
- Java 8 はネイティブにサポートされています。
- ジョブが配列の配列を使用している場合、バイト配列を除き、この状況は多次元配列として扱うことができます。

MongoDB

- 現在の MongoDB コネクタは Spark バージョン 3.1 以降と MongoDB バージョン 4.0 以降をサポートしています。
- コネクタのアップグレードにより、いくつかのプロパティ名が変更されました。例えば、URI プロパティ名が `connection.uri` に変更されました。現在のオプションの詳細については、「[MongoDB Spark Connector blog](#)」(MongoDB Spark コネクタブログ) を参照してください。
- Amazon DocumentDB によってホストされている MongoDB 4.0 を使用することには、いくつかの機能上の違いがあります。詳細については、以下のトピックを参照してください。
 - [機能の違い: Amazon DocumentDB と MongoDB](#)
 - [サポートされている MongoDB API、オペレーション、およびデータ型](#)
- 「partitioner」オプションは ShardedPartitioner、PaginateIntoPartitionsPartitioner、および SinglePartitionPartitioner に制限されています。ステージ演算子は MongoDB API をサポートしていないため、Amazon DocumentDB にはデフォルトの SamplePartitioner と PaginateBySizePartitioner を使用できません。詳細については、「[サポートされている MongoDB API、オペレーション、およびデータ型](#)」を参照してください。

Delta Lake

- Delta Lake は [SQL でのタイムトラベル](#) をサポートし、古いデータを簡単にクエリできるようになりました。今回の更新により、タイムトラベルは Spark SQL と DataFrame API の両方で利用できるようになりました。SQL の現在のバージョンの TIMESTAMP のサポートが追加されました。
- Spark 3.3 では、バッチクエリに対する `Trigger.Once` と同様に、ストリーミングクエリを実行するための [Trigger.AvailableNow](#) が導入されています。このサポートは、Delta テーブルをストリーミングソースとして使用する場合にも利用できます。
- テーブル内の列のリストを返す `SHOW COLUMNS` のサポート。
- Scala と Python の DeltaTable API での [DESCRIBE DETAIL](#) のサポート。DeltaTable API または Spark SQL のいずれかを使用して、Delta テーブルに関する詳細情報を取得します。
- SQL [Delete](#)、[Merge](#)、および [Update](#) コマンドからオペレーションメトリクスを返すことのサポート。以前は、これらの SQL コマンドは空の DataFrame を返していましたが、実行されたオペレーションに関する有用なメトリクスを含む DataFrame を返すようになりました。
- パフォーマンス向上の最適化:
 - Optimize コマンドの設定オプションを `spark.databricks.delta.optimize.repartition.enabled=true` に設定して

coalesce(1) ではなく repartition(1) を使用するようにすると、多数の小さなファイルを圧縮するときのパフォーマンスが向上します。

- キューベースのアプローチを使用して圧縮ジョブを並列化すると、[パフォーマンスが向上](#)します。
- その他の注目すべき変更点:
 - VACUUM および OPTIMIZE SQL コマンドでの[変数の使用をサポート](#)。
 - カタログテーブルによる CONVERT TO DELTA の改善:
 - パーティションスキーマが指定されていない場合は、カタログから[パーティションスキーマを自動入力](#)します。
 - ディレクトリ全体をスキャンせずに、カタログの[パーティション情報を使用](#)してコミットするデータファイルを検索します。テーブルディレクトリ内のすべてのデータファイルをコミットせずに、アクティブなパーティションのディレクトリにあるデータファイルのみがコミットされます。
 - DROP COLUMN と RENAME COLUMN が使用されていない場合、列マッピングが有効なテーブルでの[変更データフィールド \(CDF\) バッチ読み取りをサポート](#)。詳細については、[Delta Lake のドキュメント](#)を参照してください。
 - 最初のパスでスキーマのプルーニングを有効にすることで、[Update コマンドのパフォーマンスを改善](#)。

Apache Iceberg

- スキャンプランニングと Spark クエリの[パフォーマンス改善](#)をいくつか追加。
- 変更ベースのコミットを使用してサービス側のコミット競合を解決する一般的な REST カタログクライアントを追加。
- SQL タイムトラベルクエリの AS OF 構文をサポート。
- MERGE クエリと UPDATE クエリの読み取り時マージのサポートを追加。
- Z オーダーを使用してパーティションを書き換えるサポートを追加。
- [Theta スケッチ](#)やブルームフィルタのような大きな統計やインデックス blob 用の形式である Puffin の仕様と実装を追加。
- データを段階的に使用するための新しいインターフェイスを追加 (追加スキャンと変更ログスキャンの両方)。
- FileIO インターフェイスへの一括オペレーションと範囲読み取りのサポートを追加。
- メタデータツリーに削除ファイルを表示するメタデータテーブルをさらに追加。

- ドロップテーブルの動作を変更。Iceberg 0.13.1 では、DROP TABLE を実行するとカタログからテーブルが削除され、テーブルの内容も削除されます。Iceberg 1.0.0 では、DROP TABLE を実行してもカタログからテーブルが削除されるだけです。テーブルの内容を削除するには、DROP TABLE PURGE を使用します。
- Iceberg 1.0.0 では、Parquet のベクトル化された読み取りがデフォルトで有効になっています。ベクトル化された読み取りを無効にする場合は、`read.parquet.vectorization.enabled` を `false` に設定します。

Oracle

変更は軽微です。

MySQL

変更は軽微です。

Amazon Redshift

AWS Glue 4.0 には、新しい JDBC ドライバーを備えた新しい Amazon Redshift コネクタが搭載されています。機能強化と以前の AWS Glue バージョンからの移行方法については、「[the section called “Redshift 接続”](#)」を参照してください。

付録 A: 注目すべき依存関係のアップグレード

依存関係のアップグレードは次のとおりです。

依存関係	AWS Glue 4.0 でのバージョン	AWS Glue 3.0 でのバージョン	AWS Glue 2.0 でのバージョン	AWS Glue 1.0 でのバージョン
Spark	3.3.0-amzn-1	3.1.1-amzn-0	2.4.3	2.4.3
Hadoop	3.3.3-amzn-0	3.2.1-amzn-3	2.8.5-amzn-5	2.8.5-amzn-1
Scala	2.12	2.12	2.11	2.11
Jackson	2.13.3	2.10.x	2.7.x	2.7.x
[Hive]	2.3.9-amzn-2	2.3.7-amzn-4	1.2	1.2
EMRFS	2.54.0	2.46.0	2.38.0	2.30.0

依存関係	AWS Glue 4.0 でのバージョン	AWS Glue 3.0 でのバージョン	AWS Glue 2.0 でのバージョン	AWS Glue 1.0 でのバージョン
Json4s	3.7.0-M11	3.6.6	3.5.x	3.5.x
Arrow	7.0.0	2.0.0	0.10.0	0.10.0
AWS Glue データカタログクライアント	3.7.0	3.0.0	1.10.0	該当なし
Python	3.10	37	2.7 と 3.6	2.7 と 3.6
Boto	1.26	1.18	1.12	該当なし

付録 B: JDBC ドライバーのアップグレード

JDBC ドライバーのアップグレードは次のとおりです。

ドライバー	過去の AWS Glue バージョンでの JDBC ドライバーのバージョン	AWS Glue 3.0 での JDBC ドライバーのバージョン	AWS Glue 4.0 での JDBC ドライバーのバージョン
MySQL	5.1	8.0.23	8.0.23
Microsoft SQL Server	6.1.0	7.0.0	9.4.0
Oracle Database	11.2	21.1	21.7
PostgreSQL	42.1.0	42.2.18	42.3.6
MongoDB	2.0.0	4.0.0	4.7.2
Amazon Redshift	redshift-jdbc41-1.2.12.1017	redshift-jdbc41-1.2.12.1017	redshift-jdbc42-2.1.0.16

付録 C: コネクタのアップグレード

コネクタのアップグレードは次のとおりです。

ドライバー	AWS Glue 3.0 のコネクタバージョン	AWS Glue 4.0 のコネクタバージョン
MongoDB	3.0.0	10.0.4
Hudi	0.10.1	0.12.1
Delta Lake	1.0.0	2.1.0
Iceberg	0.13.1	1.0.0
DynamoDB	1.11	1.12

AWS Glue for Ray のプレビューから Ray2.4 ランタイム環境への移行

Warning

AWS Glue for Ray (プレビュー) ジョブを AWS Glue Studio で保存すると、自動的に Ray2.4 ランタイムにアップグレードされます。スクリプトに互換性の問題が発生した場合は、サポートにお問い合わせください。

AWS Glue for Ray (プレビュー) で作成された AWS Glue ジョブは AWS Glue for Ray に移行する必要があります。これには、ジョブ設定に同時にいくつかの変更を加える必要があります。

- [Runtime] フィールドには、Ray2.4 ランタイム値を入力します。これにより、基盤となる Ray バージョンが 2.0.0 から 2.4.0 にアップグレードされます。
- プレビューにデフォルトで含まれている特定の Python ライブラリは提供されなくなりました。ジョブで pandas (aws wrangler)、dask、modin、または pymars 用の AWS SDK を活用する場合は、これらを追加のライブラリとして組み込む必要があります。追加の Python ライブラリを含める方法の詳細については、「[the section called “Ray ジョブ用の Python モジュールを追加する”](#)」を参照してください。
- `--additional-python-modules` パラメータを使用している場合、このワークフローをサポートするために使用されるパラメータは `--pip-install` と `--s3-py-modules` に分割されてい

ます。これらのパラメータの詳細については、「[the section called “Ray ジョブ用の Python モジュールを追加する”](#)」を参照してください。

- `--auto-scaling-ray-min-workers` パラメータを使用している場合、名前が `--min-workers` に変更されています。

AWS Glue バージョンサポートポリシー

AWS Glue は、分析、機械学習、アプリケーション開発のためのデータの検出、準備、結合を容易にするサーバーレスデータ統合サービスです。AWS Glue ジョブでは、AWS Glue でのデータ統合作業を実行するビジネスロジックが含まれています。AWS Glue のジョブには、Spark (バッチとストリーミング)、Ray および Python シェルの 3 つのタイプがあります。ジョブを定義するときは、基盤となるランタイム環境の Spark、Ray および Python のバージョンを構成する AWS Glue バージョンを指定します。例: AWS Glue バージョン 2.0 Spark のジョブでは、Spark 2.4.3 と Python 3.7 がサポートされています。

サポートポリシー

AWS Glue では、古い AWS Glue バージョンのサポートを打ち切ることがあります。ただし、推奨していないバージョンで実行されているジョブは、テクニカルサポートの対象外となります。AWS Glue は、廃止されたバージョンに対して、セキュリティパッチやその他のアップデートは適用しません。また、推奨していないバージョンでジョブが実行される場合、AWS Glue は SLA を遵守しません。

AWS Glue のバージョン 2.0 以降のサポートが終了すると、ジョブは作成できなくなります。行えるのはジョブの編集と実行のみです。

次の AWS Glue バージョンは、サポートが終了したか、またはサポートの終了が予定されています。サポートの終了は、指定日の午前 0 時 (太平洋標準時) に開始となります。

タイプ	Glue バージョン	サポート終了
Spark	Spark 2.2、Scala 2 (Glue バージョン 0.9)	2022 年 6 月 1 日
Spark	Spark 2.2、Python 2 (Glue バージョン 0.9)	2022 年 6 月 1 日

タイプ	Glue バージョン	サポート終了
Spark	Spark 2.4、Python 2 (Glue バージョン 1.0)	2022 年 6 月 1 日
Spark	Spark 2.4、Python 3 (Glue バージョン 1.0)	2022 年 9 月 30 日
Spark	Spark 2.4、Scala 2 (Glue バージョン 1.0)	2022 年 9 月 30 日
Spark	Glue バージョン 2.0	2024 年 1 月 31 日
タイプ	Python バージョン	サポート終了
Python シェル	Python 2 (Glue バージョン 1.0)	2022 年 6 月 1 日
タイプ	ノートブックのバージョン	サポート終了
開発エンドポイント	Zeppelin ノートブック	2022 年 9 月 30 日

AWS は、ジョブをサポートしているバージョンに移行することを強く推奨します。

Spark ジョブを最新の AWS Glue バージョンへ移行する方法の詳細については、「[AWS Glue ジョブの AWS Glue バージョン 4.0 への移行](#)」を参照してください。

Python シェルのジョブを最新のものに移行するため AWS Glue バージョン:

- コンソールで、[Python 3 (Glue Version 4.0)] を選択します。
- [CreateJob/UpdateJob](#) API で、GlueVersion パラメータを 2.0 に、Command パラメータで PythonVersion を 3 に設定します。GlueVersion 設定は Python シェルジョブの動作に影響しないため、GlueVersion をインクリメントするメリットはありません。
- ジョブスクリプトを Python 3 と互換性を持たせる必要があります。

Note

2022年8月にインドネシアのジャカルタ (ap-southeast-3) リージョンが開設される前に開設されたすべてのAWSリージョンには、AWS Glueバージョン0.9/1.0のジョブを実行できるお客様の許可リストがあります。これらの古いリージョンでは、NULL値を使用してジョブを作成でき、デフォルトのバージョンはリージョンによって0.9/1.0になります。それ以降に開設されたAWSリージョンについては、APIでAWS Glueバージョンを明示的に設定する必要があります。AWS GlueはNULLパラメータを受け付けなくなりました。パラメータで0.9または1.0を渡すと、「Glueバージョン0.9(または)1.0はサポートされていません」というエラーが表示されます。

での Spark ジョブの操作 AWS Glue

Spark ETL AWS Glue ジョブに関する情報を提供します。

トピック

- [AWS Glue ジョブのパラメータ](#)
- [AWS Glue スパークとジョブズ PySpark](#)
- [AWS Glue でのストリーミング ETL ジョブ](#)
- [AWS Lake Formation FindMatches によるレコードのマッチング](#)
- [Apache Spark プログラムを AWS Glue に移行する](#)

AWS Glue ジョブのパラメータ

AWS Glue ジョブを作成するときは、Roleやなどの標準フィールドを設定しますWorkerType。[Argument] フィールド (コンソールの [ジョブパラメータ]) で追加の設定情報を指定できます。これらのフィールドでは、このトピックに記載されている引数 (パラメーター) AWS をGlue eジョブに提供できます。AWS Glue Job API の詳細については、を参照してください[the section called “ジョブ”](#)。

ジョブのパラメータを設定する

コンソールの [Job details] (ジョブの詳細) タブの [Job Parameters] (ジョブパラメータ) の見出しの下でジョブを設定できます。または、NonOverridableArgumentsジョブ上で設定したり、AWS

CLI `DefaultArguments` ジョブ実行時に設定したりして、`Arguments` ジョブを設定することもできます。ジョブに設定された引数は、ジョブが実行されるたびに渡されますが、ジョブの実行に設定された引数は、その個々の実行に対してのみ渡されます。

次の例は、ジョブのパラメータを設定するために `--arguments` を使用してジョブを実行する構文です。

```
$ aws glue start-job-run --job-name "CSV to CSV" --arguments='--scriptLocation="s3://my_glue/libraries/test_lib.py"'
```

ジョブのパラメータへのアクセス

AWS Glue スクリプトを作成する場合、ジョブのパラメータ値にアクセスして独自のコードの動作を変更したい場合があります。そのためのヘルパーメソッドをライブラリに用意しています。これらのメソッドは、ジョブパラメータ値をオーバーライドするジョブ実行パラメータ値を解決します。複数の場所で設定されたパラメータを解決する場合、ジョブの `NonOverridableArguments` はジョブ実行の `Arguments` をオーバーライドしますが、ジョブの `DefaultArguments` はオーバーライドされません。

Python では:

Python ジョブでは、`getResolvedParameters` という名前の関数を用意しています。詳細については、「[the section called “getResolvedOptions”](#)」を参照してください。ジョブパラメータは `sys.argv` 変数から取得できます。

Scala では:

Scala ジョブでは、`GlueArgParser` という名前のオブジェクトを用意しています。詳細については、「[the section called “GlueArgParser”](#)」を参照してください。ジョブパラメータは `sysArgs` 変数から取得できます。

ジョブパラメータリファレンス

AWS Glue は、ジョブおよびジョブ実行のスクリプト環境をセットアップするために使用できる以下の引数名を認識します。

--additional-python-modules

インストールする Python パッケージのセットを表すカンマ区切りのリスト。パッケージのインストールは、PyPI からでも、カスタムディストリビューションを提供することでも行

えます。PyPI パッケージのエントリは、ターゲットパッケージの PyPI 名とバージョンを含む、`package==version` のような形式になります。カスタムディストリビューションエントリとは S3 のパスであり、ディストリビューションを指しています。

エントリは Python のバージョンマッチングを使用して、パッケージとバージョンの照合を行います。つまり、2 つの等号 (==) を使用する必要があります。バージョンマッチングには他の演算子もあります。詳細については「[PEP 440](#)」を参照してください。

モジュールインストールオプションを pip3 に渡すには、`--python-modules-installer-option` パラメータを使用します。

--auto-scale-within-microbatch

デフォルト値は false です。このパラメータは、ストリーミングデータを一連のマイクロバッチで処理する AWS Glue ストリーミングジョブにのみ使用でき、auto スケーリングを有効にする必要があります。この値を false に設定すると、完了したマイクロバッチのバッチ持続時間の指数平滑移動平均線を計算し、この値をウィンドウサイズと比較して、エグゼキューターの数スケールアップするかスケールダウンするかを決定します。スケーリングはマイクロバッチが完了したときにのみ行われます。この値を true に設定すると、マイクロバッチの実行中に Spark タスクの数が 30 秒間変わらない場合や、現在のバッチ処理がウィンドウサイズを超えると、スケールアップします。エグゼキューターが 60 秒以上アイドル状態であったり、バッチ持続時間の指数平滑移動平均線が短い場合、エグゼキューター数は減少します。

--class

Scala スクリプトのエントリポイントとなる Scala クラス。これは、`--job-language` を `scala` に設定した場合にのみ適用されます。

--continuous-log-conversionPattern

連続ログ記録を有効にしたジョブのカスタム変換ログパターンを指定します。変換パターンは、ドライバーログとエグゼキューターログにのみ適用されます。AWS Glue の進行状況バーには影響しません。

--continuous-log-logGroup

連続ロギングが有効になっているジョブのカスタム Amazon CloudWatch ロググループ名を指定します。

--continuous-log-logStreamPrefix

CloudWatch 連続ロギングが有効なジョブのカスタムログストリームプレフィックスを指定します。

--customer-driver-env-vars および --customer-executor-env-vars

これらのパラメータは、各ワーカー（ドライバーまたはエグゼキューター）のオペレーティングシステムの環境変数をそれぞれ設定します。これらのパラメータは、Glue 上でプラットフォームやカスタムフレームワークを構築する場合に使用して、ユーザーが AWS Glue 上でジョブを作成できるようにすることができます。これら 2 つのフラグを有効にすると、ジョブスクリプト自体に同じロジックを注入しなくても、ドライバーおよびエグゼキューターにそれぞれ異なる環境変数を設定できます。

使用例

次の内容は、これらのパラメータを使用する例を示します。

```
"--customer-driver-env-vars", "CUSTOMER_KEY1=VAL1,CUSTOMER_KEY2=\"val2,val2 val2\"",  
"--customer-executor-env-vars", "CUSTOMER_KEY3=VAL3,KEY4=VAL4"
```

これらをジョブ実行の引数に設定することは、次のコマンドを実行することと同じです。

ドライバーに含まれるもの:

- export CUSTOMER_KEY1=VAL1
- export CUSTOMER_KEY2="val2,val2 val2"

エグゼキューターに含まれるもの:

- export CUSTOMER_KEY3=VAL3

その後、ジョブスクリプト自体で `os.environ.get("CUSTOMER_KEY1")` または `System.getenv("CUSTOMER_KEY1")` を使用して環境変数を取得できます。

強制構文

環境変数を定義するとき、次の基準に従ってください。

- 各キーには、CUSTOMER_ prefix が必要です。

たとえば "CUSTOMER_KEY3=VAL3,KEY4=VAL4" の場合、KEY4=VAL4 は無視されて設定されません。

- キーおよび値のペアは、それぞれ 1 つのカンマで区切る必要があります。

例: "CUSTOMER_KEY3=VAL3,CUSTOMER_KEY4=VAL4"

- 「値」にスペースやカンマが含まれている場合、引用符で囲んで定義する必要があります。

例 : CUSTOMER_KEY2=\"val2,val2 val2\"

この構文は、bash 環境変数の設定基準を忠実に模倣しています。

--datalake-formats

AWS Glue 3.0 以降のバージョンでサポートされています。

使用するデータレイクフレームワークを指定します。AWS Glue は、指定したフレームワークに必要な JAR ファイルをに追加します。classpath 詳細については、「[AWS Glue ETL ジョブでのデータレイクフレームワークの使用](#)」を参照してください。

次の中から 1 つまたは複数の値をカンマで区切って指定できます。

- hudi
- delta
- iceberg

例えば、3 つのフレームワークすべてを指定するには、次の引数を渡します。

```
'--datalake-formats': 'hudi,delta,iceberg'
```

--disable-proxy-v2

サービスプロキシを無効にして、Amazon S3 CloudWatch、AWS へのサービス呼び出しと AWS Glue VPC 経由のスクリプトからのサービス呼び出しを許可します。詳細については、「[すべての AWS コールを VPC を経由するように設定する](#)」を参照してください。サービスプロキシを無効にするには、このパラメータの値を true に設定します。

--enable-auto-scaling

値を true に設定した場合、Auto Scaling およびワーカーごとの課金がオンになります。

--enable-continuous-cloudwatch-log

AWS Glue ジョブのリアルタイムの連続ログ記録を有効にします。Apache Spark ジョブのログインをリアルタイムで表示できます。CloudWatch

--enable-continuous-log-filter

連続ログ記録が有効であるジョブを作成または編集するときに、標準フィルタ (true) またはフィルタなし (false) を指定します。標準フィルタを選択すると、無用な Apache Spark ドライ

バー/エグゼキュータや Apache Hadoop YARN ハートビートのログメッセージは除外されます。フィルタなしを選択すると、すべてのログメッセージが表示されます。

--enable-glue-datacatalog

AWS Glue データカタログを Apache Spark Hive メタストアとして使用できるようにします。この機能を有効にするには、値を true に設定します。

--enable-job-insights

AWS Glue ジョブ実行インサイトによる追加のエラー分析監視を可能にします。詳細については、「[the section called “AWS Glue ジョブ実行インサイトでのモニタリング”](#)」を参照してください。デフォルトでは、値は true に設定されており、ジョブ実行インサイトは有効になっています。

このオプションは、AWS Glue バージョン 2.0 および 3.0 で使用できます。

--enable-metrics

このジョブの実行のジョブプロファイリングに関するメトリクスの収集を有効にします。これらのメトリックスは、AWS Glueコンソールと Amazon CloudWatch コンソールで利用できます。このパラメータの値は関係ありません。この機能を有効にするには、このパラメータに任意の値を指定できますが、わかりやすくするために true が推奨されます。この機能を無効にするには、このパラメータをジョブ設定から削除します。

--enable-observability-metrics

AWS Glueコンソールとコンソールの Job Runs Monitoring ページで、各ジョブ実行内で何が起きているかについてのインサイトを生成するためのオブザーバビリティメトリクスのセットを有効にします。Amazon CloudWatch この機能を有効にするには、パラメータの値を true に設定します。この機能を無効にするには、false に設定するか、ジョブ設定からこのパラメータを削除します。

--enable-rename-algorithm-v2

EMRFS 名前変更アルゴリズムのバージョンをバージョン 2 に設定します。Spark ジョブがパーティションの動的な上書きモードを使用している場合、パーティションが重複して作成される可能性があります。例えば、s3://bucket/table/location/p1=1/p1=1 のような重複したパーティション作成されます。ここで、P1 は上書きされているパーティションです。名前変更アルゴリズムのバージョン 2 では、この問題が修正されています。

このオプションは、AWS Glue バージョン 1.0 でのみ使用できます。

--enable-s3-parquet-optimized-committer

Parquet データを Amazon S3 に書き込むために、EMRFS S3 最適化コミッターを有効にします。AWS Glue ジョブを作成または更新するときに、AWS Glue コンソールからパラメータ/値のペアを指定できます。値を **true** に設定すると、コミッターが有効になります。デフォルトでは、このフラグは AWS Glue 3.0 ではオン、AWS Glue 2.0 ではオフになっています。

詳細については、「[EMRFS S3 向けに最適化されたコミッターの使用](#)」を参照してください。

--enable-spark-ui

true を設定した場合、Spark UI を使用して AWS Glue ETL ジョブのモニタリングとデバッグを行う機能が有効になります。

--executor-cores

並列に実行できる Spark タスクの数。このオプションは AWS Glue 3.0 以降でサポートされています。この値は、ワーカータイプの vCPU の数の 2 倍を超えないようにします。G.1X では 8、G.2X では 16、G.4X では 32、G.8X では 64 です。タスクの並列処理が増えるとメモリやディスクに負荷がかかるだけでなく、ソースシステムとターゲットシステムがスロットリングされる (例えば、Amazon RDS での同時接続数が増える) ため、この構成を更新するときは注意が必要です。

--extra-files

スクリプトを実行する前に、AWS Glue がスクリプトの作業ディレクトリにコピーする、設定ファイルなどの追加ファイルへの Amazon S3 のパス。複数の値はコンマ (,) で区切られた完全なパスでなければなりません。ディレクトリパスではなく、個別のファイルのみがサポートされています。このオプションは Python Shell ジョブタイプではサポートされていません。

--extra-jars

スクリプトを実行する前に、AWS Glue が Java クラスパスに追加する Java .jar ファイルへの Amazon S3 パス。複数の値はコンマ (,) で区切られた完全なパスでなければなりません。

--extra-py-files

スクリプトを実行する前に、AWS Glue が Python パスに追加する Python モジュールへの Amazon S3 パス。複数の値はコンマ (,) で区切られた完全なパスでなければなりません。ディレクトリパスではなく、個別のファイルのみがサポートされています。

--job-bookmark-option

ジョブブックマークの動作を制御します。次のオプション値を設定できます。

--job-bookmark-option 値	説明
job-bookmark-enable	以前に処理されたデータを追跡します。ジョブが実行されると、最後のチェックポイントから新しいデータを処理します。
job-bookmark-disable	常にデータセット全体を処理します。以前のジョブからの出力の管理は、ユーザーが行います。
job-bookmark-pause	<p>最後のブックマークの状態は更新せずに、最後に正常に実行された後の増分データ、または次のサブオプションで識別される範囲内のデータを処理します。以前のジョブからの出力の管理は、ユーザーが行います。2つのサブオプションは以下のとおりです。</p> <ul style="list-style-type: none"> • job-bookmark-from <from-value> は、指定された実行 ID を含む最後に成功した実行までに処理されたすべての入力を表す実行 ID です。対応する入力は無視されます。 • job-bookmark-to <to-value> は、指定された実行 ID を含む最後に成功した実行までに処理されたすべての入力を表す実行 ID です。<from-value> によって識別される入力を除く対応する入力は、ジョブによって処理されます。この入力より後の入力も処理対象から除外されます。 <p>このオプションが設定されている場合、ジョブのブックマークの状態は更新されません。</p> <p>サブオプションはオプションです。ただし、使用する場合は、両方のサブオプションを指定する必要があります。</p>

例えば、ジョブブックマークを有効にするには、以下の引数を渡します。

```
'--job-bookmark-option': 'job-bookmark-enable'
```

--job-language

スクリプトプログラミング言語。この値は `scala` または `python` のいずれかである必要があります。このパラメータが存在しない場合、デフォルトで `python` が使用されます。

--python-modules-installer-option

[--additional-python-modules](#) でモジュールをインストールする場合に、`pip3` に渡されるオプションを定義するプレーンテキストの文字列。コマンドラインの場合と同様に、スペースで区切り、先頭にダッシュを付けてオプションを指定します。使用方法の詳細については、「[the section called “AWS Glue 2.0 の pip を使用した追加 Python モジュールのインストール”](#)」(Python モジュールのサポートの追加) を参照してください。

Note

Python 3.9 を使用する場合、このオプションは AWS Glue ジョブではサポートされません。

--scriptLocation

ETL スクリプトが (`s3://path/to/my/script.py` 形式で) 置かれている Amazon Simple Storage Service (Amazon S3) の場所。このパラメータは、`JobCommand` オブジェクトで設定されているスクリプトの場所を上書きします。

--spark-event-logs-path

Amazon S3 パスを指定します。Spark UI モニタリング機能を使用する場合、AWS Glue は Spark イベントログを、30 秒ごとに、この Amazon S3 パスの (Spark UI イベントを保存するための一時ディレクトリとして使用可能な) バケットにフラッシュします。

--TempDir

ジョブの一時ディレクトリとして使用できるバケットへの Amazon S3 のパスを指定します。

例えば、一時ディレクトリを設定するには、以下の引数を渡します。

```
'--TempDir': 's3-path-to-directory'
```

Note

AWS Glue では、リージョンにバケットが存在していない場合、ジョブの一時バケットが作成されます。このバケットは、パブリックアクセスを許可する場合があります。

す。Amazon S3 のバケットは、パブリックアクセスブロックを設定するように変更するか、そのリージョンのすべてのジョブが完了した後に削除するかのいずれかが可能です。

--use-postgres-driver

この値を true に設定した場合、Amazon Redshift の JDBC ドライバーとの競合を避けるために、クラスパスにある Postgres JDBC ドライバーが優先されます。このオプションは、AWS Glue バージョン 2.0 でのみ使用可能です。

--user-jars-first

この値を true に設定した場合、クラスパス内にあるお客様の追加 JAR ファイルが優先されます。このオプションは、AWS Glue バージョン 2.0 以降でのみ使用可能です。

--conf

Spark の設定パラメータを制御します。高度なユースケース向けです。

--encryption-type

レガシーパラメータ。対応する動作は、セキュリティ設定を使用して設定する必要があります。セキュリティ設定の詳細については、「[the section called “AWS Glue によって書き込まれたデータの暗号化”](#)」を参照してください。

AWS Glue では内部で以下のような引数を使用するので、それは絶対に使用しないでください。

- --debug — AWS Glue 内部用。設定する必要はありません。
- --mode — AWS Glue 内部用。設定する必要はありません。
- --JOB_NAME — AWS Glue 内部用。設定する必要はありません。
- --endpoint — AWS Glue 内部用。設定する必要はありません。

AWS Glue はサイト固有のカスタマイズを実行するための sitecustomize を使用した Python の site モジュールによる環境のブートストラップをサポートします。独自の初期化関数のブートストラップは、高度なユースケースにのみ推奨され、AWS Glue 4.0 ではベストエフォートベースでサポートされています。

環境変数のプレフィックス GLUE_CUSTOMER は、お客様専用です。

AWS Glue スパークとジョブズ PySpark

以下のセクションでは、AWS Glue Spark PySpark とジョブについて説明します。

トピック

- [Spark ジョブと PySpark ジョブをAWS Glueに追加する](#)
- [ジョブのブックマークを使用した処理済みデータの追跡](#)
- [AWS Glue Spark シャッフルプラグインと Amazon S3](#)
- [モニタリングAWS GlueSpark ジョブ](#)

Spark ジョブと PySpark ジョブをAWS Glueに追加する

以下のセクションでは、Spark と PySpark のジョブをAWS Glueに追加する方法について説明します。

トピック

- [での Spark ジョブのジョブプロパティの設定 AWS Glue](#)
- [AWS Glueコンソールで Spark スクリプトの編集](#)
- [仕事 \(レガシー\)](#)

での Spark ジョブのジョブプロパティの設定 AWS Glue

AWS Glue ジョブには、ソースデータに接続して処理し、データターゲットに書き出すスクリプトがカプセル化されています。通常、ジョブは、抽出、変換、ロード (ETL) スクリプトを実行します。ジョブでは、汎用 Python スクリプト (Python シェルジョブ) を実行することもできます。AWS Glue トリガーでは、スケジュールまたはイベントに基づいて、またはオンデマンドでジョブを開始できます。ジョブ実行をモニタリングすると、完了ステータス、継続時間、開始時間などのランタイムメトリクスを知ることができます。

AWS Glue で生成されたスクリプトを使用することも、独自のスクリプトを使用することもできます。ソーススキーマとターゲットの場所またはスキーマを使用すると、AWS Glueコードジェネレーターは Apache Spark API (PySpark) スクリプトを自動的に作成できます。このスクリプトを出発点として使用し、目標に合わせて編集できます。

AWS Glue は、JSON、CSV、ORC (Optimized Row Columnar)、Apache Parquet、Apache Avro などのいくつかのデータ形式で出力ファイルを書き込むことができます。一部のデータ形式では、一般的な圧縮形式を記述できます。

AWS Glue のジョブには、Spark、ストリーミング ETL、Python シェルの 3 タイプがあります。

- Spark ジョブは、によって管理される Apache Spark 環境で実行されます AWS Glue。データはバッチで処理されます。
- ストリーミング ETL ジョブは Spark ジョブに似ていますが、データストリームで ETL を実行する点が異なります。このタイプのジョブでは Apache Spark 構造化ストリーミングフレームワークが使用されます。一部の Spark ジョブ機能は、ストリーミング ETL ジョブでは使用できません。
- Python シェルジョブは Python スクリプトをシェルとして実行し、使用している AWS Glue バージョンに応じた Python バージョンをサポートします。このタイプのジョブでは、Apache Spark 環境を必要としないタスクをスケジューリングして実行できます。

Spark ジョブのジョブプロパティの定義

AWS Glue コンソールでジョブを定義するときに、AWS Glue ランタイム環境をコントロールするためのプロパティの値を指定します。

次のリストは、Spark ジョブのプロパティについて説明しています。Python シェルジョブのプロパティについては、「[Python シェルジョブのジョブプロパティの定義](#)」を参照してください。ストリーミング ETL ジョブのプロパティについては、「[the section called “ストリーミング ETL ジョブのジョブプロパティの定義”](#)」を参照してください。

プロパティは、AWS Glue コンソールの [Add job] (ジョブの追加) ウィザードに表示された順に一覧表示されます。

名前

UTF-8 文字を 255 文字以内で入力します。

説明

オプションとして最大 2,048 文字の説明を提示します。

IAM ロール

ジョブ実行とデータストアへのアクセスに使用されるリソースの認証に使用する IAM ロールを指定します。AWS Glue でジョブを実行するためのアクセス権限の詳細については、[AWS Glue の Identity and Access Management](#) を参照してください。

タイプ

ETL ジョブの種類。これは、選択するデータソースの種類に基づいて自動的に設定されます。

- [Spark] は、`glueetl` のジョブコマンドで Apache Spark ETL スクリプトを実行します。
- [Spark ストリーミング] は、`gluestreaming` のジョブコマンドで Apache Spark ストリーミング ETL スクリプトを実行します。詳細については、「[the section called “ストリーミング ETLジョブ”](#)」を参照してください。
- [Python シェル] は、`pythonshell` のジョブコマンドで Python スクリプトを実行します。詳細については、「[AWS Glue での Python シェルジョブに関するジョブプロパティの設定](#)」を参照してください。

AWS Glue バージョン

AWS Glue のバージョンによって、ジョブで使用できる Apache Spark および Python のバージョンが次の表に指定されているように決まります。

AWS Glue バージョン	サポートされている Spark および Python のバージョン
4.0	<ul style="list-style-type: none"> • Spark 3.3.0 • Python 3.10
3.0	<ul style="list-style-type: none"> • Spark 3.1.1 • 「Python 3.7」
2.0	<ul style="list-style-type: none"> • Spark 2.4.3 • 「Python 3.7」
1.0	<ul style="list-style-type: none"> • Spark 2.4.3 • Python 2.7 • Python 3.6
0.9	<ul style="list-style-type: none"> • Spark 2.2.1 • Python 2.7

ワーカータイプ

以下のワーカータイプを使用できます。

AWS Glue ワーカーで使用できるリソースは DPU。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。

- **G.1X** – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは、84 GB のディスク (約 34 GB の空き) を備えた 1 DPU (4 vCPU、16 GB のメモリ) にマッピングされます。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケーラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- **G.2X** – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは、128 GB のディスク (約 77 GB の空き) を備えた 2 DPU (8 vCPU、32 GB のメモリ) にマッピングされます。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケーラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- **G.4X** – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは、256 GB のディスク (約 235 GB の空き) を備えた 4 DPU (16 vCPU、64 GB のメモリ) にマッピングされます。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、米国東部 (オハイオ)、米国東部 (バージニア北部)、米国西部 (オレゴン)、アジアパシフィック (シンガポール)、アジアパシフィック (シドニー)、アジアパシフィック (東京)、カナダ (中部)、欧州 (フランクフルト)、欧州 (アイルランド)、欧州 (ストックホルム) の各 AWS リージョンで、AWS Glueバージョン 3.0 以降の Spark ETL ジョブでのみ使用できます。
- **G.8X** – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは、512 GB のディスク (約 487 GB の空き) を備えた 8 DPU (32 vCPU、128 GB のメモリ) にマッピングされます。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、ワーカータイプでサポートされているのと同じ AWS リージョンで、AWS Glueバージョン 3.0 以降の Spark ETL G.4X ジョブでのみ使用できます。
- **G.025X** – このタイプを選択する場合は、[Number of workers] (ワーカー数) の値も指定します。各ワーカーは、84 GB のディスク (約 34 GB の空き) を備えた 0.25 DPU (2 vCPU、4 GB のメモリ) にマッピングされます。少量のストリーミングジョブには、このワーカータイプをお勧めします。このワーカータイプは、AWS Glueバージョン 3.0 のストリーミングジョブでのみ使用できます。

ETL ジョブの実行に使用された DPU の数に基づいて時間あたりの料金が請求されます。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

AWS Glue バージョン 1.0 以前のジョブでは、コンソールを使用してジョブを設定し、[Worker type] (ワーカータイプ) に [Standard] (標準) を指定する場合は、[Maximum capacity] (最大キャパシティー) が設定され、[Number of workers] (ワーカー数) は [Maximum capacity] (最大キャパシティー) から 1 を引いた値になります。AWS Command Line Interface (AWS CLI) または AWS

SDK を使用する場合は、最大容量パラメータを指定するか、ワーカータイプとワーカー数の両方を指定できます。

AWS Glue バージョン 2.0 以降のジョブでは、[最大キャパシティ] の指定はできません。代わりに、[Worker type] (ワーカータイプ) と [Number of workers] (ワーカー数) を指定します。

[言語]

ETL スクリプト内のコードでジョブのロジックを定義します。Python または Scala でスクリプトを記述できます。ジョブが実行するスクリプトを AWS Glue によって生成するのか、それとも自分で提供するのかを選択できます。Amazon Simple Storage Service (Amazon S3) でスクリプト名と場所を指定します。パスのスクリプトディレクトリと同じ名前のファイルが存在していないことを確認します。スクリプトの記述の詳細については、「[AWS Glue プログラミングガイド](#)」を参照してください。

リクエストしたワーカーの人数

ほとんどのワーカータイプで、ジョブの実行時に割り当てられるワーカーの数を指定する必要があります。

ジョブのブックマーク

ジョブ実行時に AWS Glue が状態情報を処理する方法を指定します。以前に処理されたデータの記憶、状態情報の更新、または状態情報の無視を指定できます。詳細については、「[the section called “ジョブのブックマークを使用した処理済みデータの追跡”](#)」を参照してください。

Flex 実行

AWS Studio または API を使用してジョブを設定する場合、標準または柔軟なジョブ実行クラスを指定できます。ジョブには、さまざまな度合いの優先順位や時間的な制約を設定することができます。標準の実行クラスは、素早くジョブを起動する必要があり、専用のリソースが必要な時間的な制約のあるワークロードに最適です。

柔軟な実行クラスは、実稼働前のジョブ、テスト、1 回限りのデータの読み込みなど、緊急性のないジョブに適しています。柔軟なジョブの実行は、AWS Glue バージョン 3.0 以降と G.1X または G.2X ワーカータイプを使用するジョブでサポートされています。

Flex ジョブの実行は、任意の時点で実行されているワーカーの数に基づいて課金されます。柔軟なジョブの実行では、ジョブ数を追加または削除することができます。Max Capacity * Execution Time という単純な計算で課金されるのではなく、各ワーカーがジョブ実行中に実行された時間が含まれます。請求は (Number of DPU's per worker * time each worker ran) の合計です。

詳細については、AWS Studio のヘルプパネル、または [ジョブ](#) および [ジョブ実行](#) を参照してください。

再試行回数

失敗した場合に AWS Glue がジョブを自動的に再起動する回数を 0～10 の間で指定します。タイムアウト制限に達したジョブは再起動されません。

ジョブのタイムアウト

最大の実行時間 (分) を設定します。バッチジョブのデフォルト値は 2880 分 (48 時間) です。ジョブ実行時間がこの制限値を超えると、ジョブ実行状態は TIMEOUT に変わります。

ストリーミングジョブのタイムアウト値は 7 日または 10080 分未満である必要があります。値を空白のままにすると、メンテナンスウィンドウを設定していない場合、ジョブは 7 日後に再開されます。メンテナンスウィンドウを設定している場合、メンテナンスウィンドウ中に 7 日後に再起動されます。

ジョブタイムアウトのベストプラクティス

ジョブは実行時間に基づいて課金されます。予期しない課金を避けるには、ジョブの予想実行時間に適切なタイムアウト値を設定してください。

詳細プロパティ

スクリプトのファイル名

ジョブの一意的スクリプト名。[タイトルがないジョブ] という名前を付けられません。

スクリプトパス

スクリプトの Amazon S3 ロケーション。パスは `s3://bucket/prefix/path/` の形式で指定する必要があります。末尾にスラッシュ (/) を付けてファイルは含めないでください。

ジョブのメトリクス

このジョブの実行時に Amazon CloudWatch メトリクスの作成を有効または無効にします。プロファイリングデータを表示するには、このオプションを有効にする必要があります。メトリクスを有効にして表示する方法の詳細については、「[ジョブのモニタリングとデバッグ](#)」を参照してください。

ジョブのオブザーバビリティメトリクス

このジョブの実行時に追加のオブザーバビリティ CloudWatch メトリクスの作成を有効にします。詳細については、「[the section called “AWS Glue オブザーバビリティメトリクスを使用したモニタリング”](#)」を参照してください。

連続ログ記録

Amazon への継続的なログ記録を有効にします CloudWatch。このオプションが有効になっていない場合、ログはジョブの完了後にのみ使用できます。詳細については、[the section called “AWS Glue ジョブの連続ログ記録”](#) を参照してください。

Spark UI

このジョブをモニタリングするために Spark UI の使用を有効にします。詳細については、「[AWS Glue ジョブ用の Apache Spark ウェブ UI の有効化](#)」を参照してください。

Spark UI ログパス

Spark UI が有効になっているときにログを書き込むパス。

Spark UI のログ記録とモニタリングの設定

以下のオプションのいずれかを選択します。

- 標準: AWS Glue ジョブ実行 ID をファイル名として使用してログを書き込みます。AWS Glue コンソールで Spark UI モニタリングを有効にします。
- レガシー: 「spark-application-{timestamp}」をファイル名として使用してログを書き込みます。Spark UI モニタリングをオンにしないでください。
- スタンダードとレガシー: スタンダードロケーションとレガシーロケーションの両方にログを書き込みます。AWS Glue コンソールで Spark UI モニタリングを有効にします。

最大同時実行数

このジョブで許可される同時実行の最大数を設定します。デフォルトは 1 です。このしきい値に達すると、エラーが返されます。指定できる最大値は、サービスの制限によってコントロールされます。たとえば、新しいインスタンスの開始時に前回のジョブがまだ実行されている場合、同じジョブの 2 つのインスタンスが同時に実行されないようにエラーを戻すことができます。

一時的なパス

AWS Glue がスクリプトを実行するときに一時的な中間結果が書き込まれる Amazon S3 の作業ディレクトリの場所を指定します。パスの一時ディレクトリと同じ名前のファイルが存在し

ていないことを確認します。このディレクトリは、Amazon Redshift に対して AWS Glue が読み取りと書き込みをするときに使用します。また、特定の AWS Glue 変換で使用します。

Note

AWS Glue では、リージョンにバケットが存在しない場合、ジョブの一時バケットが作成されます。このバケットは、パブリックアクセスを許可する場合があります。この Amazon S3 に置かれたバケットは、パブリックアクセスブロックを構成するために設定することができます。また、そのリージョンのすべてのジョブが完了した後に削除することが可能です。

遅延通知のしきい値 (分)

遅延通知を送信するまでのしきい値 (分単位) を設定します。RUNNING、STARTING、または STOPPING ジョブの実行が想定時間 (分単位) を超えると通知を送信するように、このしきい値を設定できます。

セキュリティ設定

リストからセキュリティ設定を選択します。セキュリティ設定では、Amazon S3 ターゲットのデータの暗号化方法として、暗号化なし、AWS KMSで管理されたキー (SSE-KMS) を使用したサーバー側の暗号化、または Amazon S3 で管理された暗号化キー (SSE-S3) を使用したサーバー側の暗号化を指定します。

サーバー側の暗号化

このオプションを選択すると、ETL ジョブが Amazon S3 に書き込むときに、データは SSE-S3 暗号化を使用して保管時に暗号化されます。Amazon S3 のデータターゲットと、Amazon S3 の一時ディレクトリに書き込まれるデータは、両方とも暗号化されています。このオプションはジョブパラメータとして渡されます。詳細については、Amazon Simple Storage Service ユーザーガイドの「[Amazon S3 が管理する暗号化キーによるサーバー側の暗号化 \(SSE-S3\) を使用したデータの保護](#)」を参照してください。

Important

セキュリティ設定を指定している場合、このオプションは無視されます。

Glue データカタログを Hive メタストアとして使用する

AWS Glue データカタログを Hive メタストアとして使用することを選択します。ジョブに使用される IAM ロールには、`glue:CreateDatabase` アクセス許可が必要です。「default」という名前のデータベースが存在しない場合はデータカタログに作成されます。

接続

VPC 設定を選択し、仮想プライベートクラウド (VPC) にある Amazon S3 データソースにアクセスします。AWS Glueでネットワーク接続を作成および管理できます。詳細については、「[データへの接続](#)」を参照してください。

[Libraries] (ライブラリ)

Python ライブラリパス、依存 JAR パス、参照されるファイルパス

これらのオプションはスクリプトで必要に応じて指定します。ジョブを定義するときに、これらのオプションのカンマ区切りの Amazon S3 パスを定義できます。ジョブ実行時にこれらのパスを上書きできます。詳細については、「[独自のカスタムスクリプトの提供](#)」を参照してください。

ジョブのパラメータ

スクリプトに名前付きパラメータとして渡される一連のキーと値のペア。これらは、スクリプトの実行時に使用されるデフォルト値ですが、トリガーまたはジョブの実行時に上書きできます。キー名の前に `--` を付ける必要があります (`--myKey` など)。を使用する場合、ジョブパラメータをマップとして渡します AWS Command Line Interface。

例については、「[AWS Glue の Python パラメータの受け渡しとアクセス](#)」で Python パラメータを参照してください。

タグ

[タグキー] とオプションの [タグ値] を使用してジョブにタグを付けます。作成されたタグキーは読み取り専用になります。リソースを整理、識別しやすいように、いくつかのリソースでタグを使用します。詳細については、「[AWS Glue のタグ](#)」を参照してください。

Lake Formation 管理対象テーブルにアクセスするジョブの制約事項

によって管理されるテーブルに対して読み書きするジョブを作成するときは、次の注意事項と制限に注意してください AWS Lake Formation。

- 次の機能は、セルレベルフィルタを使用してテーブルにアクセスするジョブではサポートされません。

- [ジョブのブックマークと境界実行の実行](#)
- [プッシュダウン述語](#)
- [サーバーサイドのカタログパーティション述語](#)
- [enableUpdateCatalog](#)

AWS Glueコンソールで Spark スクリプトの編集

スクリプトには、ソースからデータを抽出し、変換してターゲットに読み込むコードが含まれています。AWS Glue ジョブの開始時にスクリプトを実行します。

Python または Scala で AWS Glue ETL スクリプトを記述できます。Python スクリプトは、抽出、変換、読み込み (ETL) ジョブに PySpark Python 方言を拡張した言語を使用します。スクリプトには ETL 変換を処理する拡張構造が含まれます。自動でジョブのソースコードロジックを生成するときに、スクリプトが作成されます。このスクリプトを編集するか、または、独自のスクリプトを指定して ETL 作業を処理することができます。

AWS Glue のスクリプトの定義と編集の詳細については、「[AWS Glue プログラミングガイド](#)」を参照してください。

その他のライブラリまたはファイル

スクリプトに追加のライブラリやファイルが必要な場合は、次のように指定できます。

Python ライブラリパス

スクリプトで必要とされる Python ライブラリへのカンマ区切りの Amazon Simple Storage Service (Amazon S3) パス。

Note

純粋な Python ライブラリのみを使用できます。pandas Python データ解析ライブラリなど、C 拡張機能に依存するライブラリはまだサポートされていません。

依存 JARS パス

スクリプトで必要とされる JAR ファイルへのカンマ区切りの Amazon S3 パスです。

Note

現在、純粋な Java または Scala (2.11) ライブラリのみを使用できます。

参照されるファイルパス

スクリプトに必要な追加のファイル (例えば、設定ファイル) への、カンマ区切りの Amazon S3 パス。

仕事 (レガシー)

スクリプトには、抽出、変換、ロード (ETL) ワークを実行するコードが含まれます。独自のスクリプトを提供することもできますし、お客様のガイダンスで AWS Glue がスクリプトを生成することもできます。独自のスクリプトの作成については、「[独自のカスタムスクリプトの提供](#)」を参照してください。

スクリプトは、AWS Glue コンソールで編集できます。スクリプトを編集する場合、ソース、ターゲット、および変換を追加することができます。

スクリプトを編集するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。その後 [Jobs] (ジョブ) タブを選択します。
2. リストでジョブを選択し、次に [Action] (アクション)、[Edit script] (スクリプトの編集) を選択してスクリプトエディタを開きます。

ジョブ詳細ページからスクリプトエディタにアクセスすることもできます。[Script] (スクリプト) タブを選択し、次に [Edit script] (スクリプトの編集) を選択します。

スクリプトエディタ

AWS Glue スクリプトエディタを使用して、スクリプトのソース、ターゲット、変換を挿入、変更、および削除できます。スクリプトエディタにはスクリプトとダイアグラムの両方が表示され、データの流れを可視化しやすくなります。

スクリプトのダイアグラムを作成するには、[ダイアグラムの生成] を選択します。AWS Glue は ## で始まるスクリプト内の注釈行を使用してダイアグラムをレンダリングします。ダイアグラムでスク

リプトを正しく表すために、注釈のパラメータと Apache Spark コードのパラメータの同期を保つ必要があります。

スクリプトエディタを使用して、スクリプトのカーソルが置かれている任意の場所にコードテンプレートを追加することができます。エディタの上部で、次のオプションから選択します。

- ソーステーブルをスクリプトに追加するには、[Source] (ソース) を選択します。
- ターゲットテーブルをスクリプトに追加するには、[Target] (ターゲット) を選択します。
- ターゲット位置をスクリプトに追加するには、[Target location] (ターゲット位置) を選択します。
- 変換をスクリプトに追加するには、[Transform] (変換) を選択します。スクリプトで呼び出される関数については、「[で AWS Glue ETL スクリプトをプログラムする PySpark](#)」を参照してください。
- スピゴット変換をスクリプトに追加するには、[Spigot] (スピゴット) を選択します。

挿入されたコードで、注釈および Apache Spark コード両方の parameters を変更します。たとえば、スピゴット変換を追加したら、path が @args 注釈行および output コード行の両方で置き換えられていることを検証します。

[Logs] (ログ) タブでは、実行されるジョブに関連するログが表示されます。最新の 1,000 行が表示されます。

[Schema] (スキーマ) タブでは、Data Catalog で使用可能な場合、選択されたソースとターゲットのスキーマが表示されます。

ジョブのブックマークを使用した処理済みデータの追跡

AWS Glue ではジョブの実行による状態情報を保持することで、ETL ジョブの以前の実行中にすでに処理されたデータを追跡します。この継続状態の情報はジョブのブックマークと呼ばれています。ジョブのブックマークは、AWS Glue で状態情報を保持して、古いデータを再処理しないために役立ちます。ジョブのブックマークを使用すると、スケジュールされた間隔で再実行する際に新しいデータを処理できます。ジョブのブックマークは、ソース、変換、ターゲットなど、さまざまなジョブの要素で構成されています。例えば、ETL ジョブが Amazon S3 ファイルで新しいパーティションを読み込むとします。AWS Glue は、そのジョブにより正常に処理されたのはどのパーティションなのかを追跡し、処理の重複およびジョブのターゲットデータストアにデータが重複するのを防ぎます。

ジョブのブックマークは、JDBC データソース、Relationalize 変換、および一部の Amazon Simple Storage Service (Amazon S3) ソースに実装されています。次の表に、AWS Glue がジョブのブックマークに対してサポートする Amazon S3 ソース形式を示します。

AWS Glue バージョン	Amazon S3 ソース形式
バージョン 0.9	JSON、CSV、Apache Avro、XML
バージョン 1.0 以降	JSON、CSV、Apache Avro、XML、Parquet、ORC

AWS Glue バージョンの詳細に関しては、「[Spark ジョブのジョブプロパティの定義](#)」を参照してください。

ジョブのブックマーク機能には、AWS Glue スクリプトからアクセスする際の追加機能があります。生成されたスクリプトを参照すると、この機能に関連する変換コンテキストが表示される場合があります。詳細については、「[the section called “ジョブのブックマークを使用する”](#)」を参照してください。

トピック

- [AWS Glue でジョブのブックマークを使用する](#)
- [ジョブブックマーク機能の運用に関する詳細](#)

AWS Glue でジョブのブックマークを使用する

ジョブのブックマークオプションは、ジョブが開始したときにパラメータとして渡されます。AWS Glue コンソールでジョブのブックマークを設定するためのオプションを次の表に示します。

ジョブのブックマーク	説明
有効	ジョブの実行後に状態を更新させて以前に処理されたデータを追跡します。ジョブのブックマークをサポートしているソースのあるジョブの場合、ジョブは処理されたデータを追跡し、ジョブが実行されると、最後のチェックポイント以降の新しいデータを処理します。
[無効]	ジョブのブックマークは使用されず、ジョブは常にデータセット全体を処理します。以前のジョブからの出力の管理は、ユーザーが行います。これがデフォルトです。
[Pause] (一時停止)	最後のブックマークの状態は更新せずに、最後に正常に実行された後の増分データ、または次のサブオプションで識別される範囲内のデータを

ジョブのブックマーク	説明
	<p>処理します。以前のジョブからの出力の管理は、ユーザーが行います。次の 2 つのサブオプションがあります。</p> <ul style="list-style-type: none">• <code>job-bookmark-from <from-value></code> は、指定された実行 ID を含む最後に成功した実行までに処理されたすべての入力を表す実行 ID です。対応する入力は無視されます。• <code>job-bookmark-to <to-value></code> は、指定された実行 ID を含む最後に成功した実行までに処理されたすべての入力を表す実行 ID です。<code><from-value></code> によって識別される入力を除く対応する入力は、ジョブによって処理されます。この入力より後の入力も処理対象から除外されません。 <p>このオプションが設定されている場合、ジョブのブックマークの状態は更新されません。</p> <p>サブオプションはオプションですが、使用する場合、両方のサブオプションを提供する必要があります。</p>

コマンドラインでジョブに渡されるパラメータ、および特にジョブブックマークの詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。

Amazon S3 入力ソースの場合、AWS Glue ジョブのブックマークではオブジェクトの最終更新日時を確認して、どのオブジェクトを再処理する必要があるのかを確認します。入力ソースデータが最後のジョブ実行以降に変更されている場合、ジョブを再度実行すると、ファイルが再度処理されます。

JDBC ソースでは、次のルールが適用されます。

- AWS Glue では、テーブルごとに 1 つ以上の列をブックマークキーとして使用して、新しいデータおよび処理済みのデータを決定します。ブックマークキーは、結合して単一の複合キーを形成します。
- AWS Glue はデフォルトでプライマリキーをブックマークキーとして使用します。ただし、プライマリキーが連続して (ギャップなく) 増減することが条件です。

- AWS Glue スクリプトでブックマークキーとして使用する列を指定できます。AWS Glue スクリプトでのジョブブックマーク使用の詳細については、「[the section called “ジョブのブックマークを使用する”](#)」を参照してください。
- AWS Glue は、ジョブブックマークキーとして名前の大文字と小文字を区別する列の使用をサポートしていません。

AWS Glue Spark ETL ジョブのジョブブックマークを以前のジョブ実行に巻き戻すことができます。ジョブのブックマークを以前のジョブ実行に巻き戻すことで、データのバックフィルシナリオをより適切にサポートできます。その結果、後続のジョブ実行ではブックマークされたジョブ実行からのデータだけが再処理されます。

同じジョブを使用してすべてのデータを再処理する場合は、ジョブのブックマークをリセットします。ジョブのブックマークの状態をリセットするには、AWS Glue コンソール、[ResetJobBookmark アクション \(Python: reset_job_bookmark\)](#) API オペレーション、または AWS CLI を使用します。例えば、AWS CLI を使用して以下のコマンドを入力します。

```
aws glue reset-job-bookmark --job-name my-job-name
```

ブックマークを巻き戻したりリセットしたりしたとき、複数のターゲットが存在する可能性があり、ターゲットがジョブブックマークで追跡されないため、AWS Glue はターゲットファイルをクリーンアップしません。ジョブブックマークで追跡されるのは、ソースファイルだけです。ソースファイルの巻き戻しと再処理時に異なる出力ターゲットを作成して、出力内のデータが重複しないようにすることができます。

AWS Glue では、ジョブでジョブのブックマークを管理します。ジョブを削除すると、ジョブのブックマークは削除されます。

場合によっては、AWS Glue ジョブのブックマークを有効にしてあっても、以前の実行で処理したデータを ETL ジョブで再処理することがあります。このエラーの一般的な原因を解決する方法の詳細については、「[AWS Glue for Spark のエラーのトラブルシューティング](#)」を参照してください。

ジョブブックマーク機能の運用に関する詳細

このセクションでは、ジョブのブックマークを使用する運用の詳細を説明します。

ジョブのブックマークはジョブの状態を保存します。状態の各インスタンスは、ジョブ名とバージョン番号がキーになっています。スクリプトで呼び出された `job.init` では、その状態を取得し、常

に最新バージョンを取得します。1つの状態の中に複数の状態要素が存在します。要素は、スクリプト内のソース、変換、シンクインスタンスごとに固有です。これらの状態要素は、スクリプト内の対応する要素(ソース、変換、またはシンク)にアタッチされている変換コンテキストによって識別されます。状態要素はユーザースクリプトから `job.commit` が呼び出されたときにアトミックに保存されます。スクリプトは引数からジョブのブックマークのジョブ名およびコントロールオプションを取得します。

ジョブのブックマーク内の状態要素は、ソース、変換、またはシンク固有のデータです。例えば、アップストリームジョブまたはプロセスによって継続的に書き込まれている Amazon S3 の場所から増分データを読み取るとします。この場合、スクリプトではそれまでに処理されている部分を判別する必要があります。Amazon S3 ソースのジョブのブックマーク実装では情報を保存するため、再度ジョブを実行するときに、保存された情報を使用して新しいオブジェクトのみをフィルターでき、次にジョブを実行するための状態を再計算できます。新しいファイルをフィルタリングするためにタイムスタンプが使用されます。

ジョブのブックマークは、状態要素に加えて、実行番号、試行番号、バージョン番号を持ちます。実行番号はジョブの実行を追跡し、試行番号はジョブ実行の試行回数を記録します。ジョブ実行番号は正常な実行ごとに増分される、一定間隔で増加する番号です。試行番号は各実行の試行回数を追跡し、失敗した試行後にのみ実行がある場合にのみ増分されます。バージョン番号は、一定間隔で増加し、ジョブのブックマークの更新を追跡します。

AWS Glue サービスデータベースでは、すべての変換のブックマーク状態がキーと値のペアとして一緒に格納されます。

```
{
  "job_name" : ...,
  "run_id": ...,
  "run_number": ..,
  "attempt_number": ...
  "states": {
    "transformation_ctx1" : {
      bookmark_state1
    },
    "transformation_ctx2" : {
      bookmark_state2
    }
  }
}
```

ベストプラクティス

以下は、ジョブブックマークを使用するためのベストプラクティスです。

- ブックマークを有効にした状態でデータソースプロパティを変更しないでください。例えば、Amazon S3 入力パス A を指している `datasource0` があります。ジョブはブックマークを有効にして数ラウンド実行されているソースから読み出しています。`transformation_ctx` を変更せずに `datasource0` の入力パスを Amazon S3 パス B に変更すると、AWS Glue ジョブは保存されている古いブックマークの状態を使用します。これによって、入力パス B のファイルは、見つからないか、スキップされます。AWS Glue は、これらのファイルが以前の実行で処理されたと仮定します。
- パーティション管理を改善するために、ブックマーク付きのカタログテーブルを使用してください。ブックマークは、データカタログのデータソースとオプションの両方に対して機能します。ただし、`from` オプションのアプローチでは、新しいパーティションを削除/追加するのは困難です。クローラーでカタログテーブルを使用すると、新しく追加された [パーティション](#) を追跡するための自動化が向上し、[プッシュダウン述語](#) を使用して特定のパーティションを柔軟に選択できるようになります。
- 大規模なデータセット向けの [AWS Glue Amazon S3 ファイルリスター](#) の使用 ブックマークは、各入力パーティションの下すべてのファイルをリストし、ファイリングを実行します。したがって、単一のパーティションの下にファイルが多すぎると、ブックマークがドライバ OOM で実行される可能性があります。AWS Glue Amazon S3 ファイルリスターを使用して、メモリ内のすべてのファイルを一度にリストしないようにします。

AWS Glue Spark シャッフルプラグインと Amazon S3

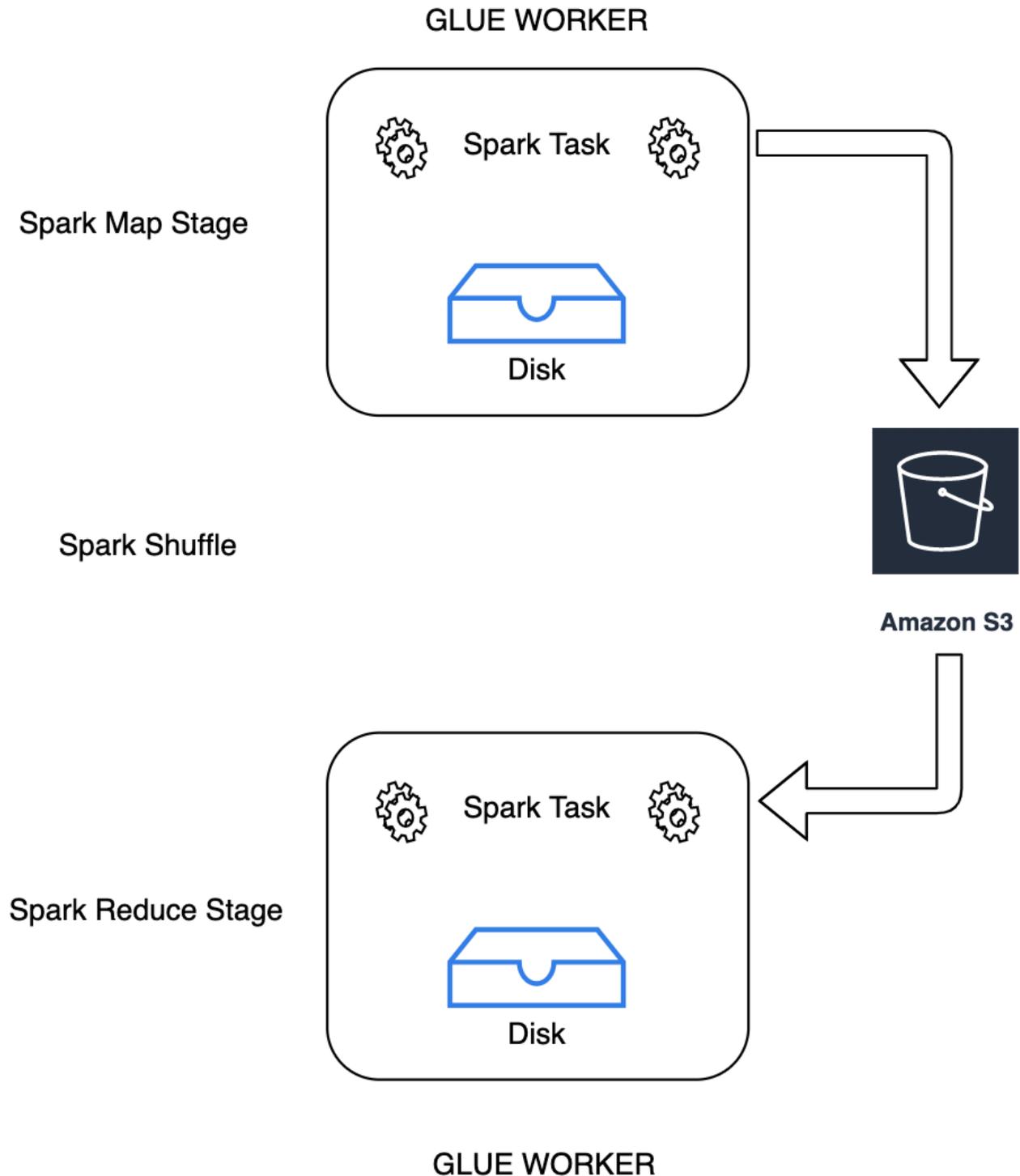
シャッフルは、データがパーティション間で再配置されるたびに行われる Spark ジョブの重要なステップです。これが必要なのは、`join`、`groupByKey`、`reduceByKey`、`repartition` などのさまざまな変換に、処理を完了するために他のパーティションからの情報が必要なためです。Spark は各パーティションから必要なデータを収集し、新しいパーティションに結合します。シャッフル中、データはディスクに書き込まれ、ネットワーク経由で転送されます。その結果、シャッフルオペレーションはローカルディスク容量に制約されます。Spark の `No space left on device` または `MetadataFetchFailedException` エラーは、エグゼキュターに十分なディスク領域がなく、リカバリがない場合に発生します。

Note

Amazon S3 での AWS Glue Spark シャッフルプラグインは、AWS Glue ETL ジョブについてのみサポートされています。

ソリューション

AWS Glue では、Amazon S3 を使用して Spark シャッフルデータを保存できるようになりました。Amazon S3 は、業界をリードするスケーラビリティ、データ可用性、セキュリティ、およびパフォーマンスを提供するオブジェクトストレージサービスです。このソリューションで、Spark ジョブ用のコンピューティングとストレージが非集約化され、完全な伸縮自在性と低コストのシャッフルストレージが得られ、シャッフル負荷の高いワークロードが高い信頼性で実行できるようになります。



Amazon S3 を使用するための Cloud Shuffle Storage Plugin for Apache Spark を導入します。大きなシャッフルオペレーションのためにローカルディスク容量の制約を受けることがわかっている場合、Amazon S3 シャッフルを有効にすると、AWS Glue ジョブを失敗することなく高い信頼性で実行できます。小さなパーティションが多数ある場合や、Amazon S3 に書き込まれたファイルを

シャッフルする場合、Amazon S3 へのシャッフルはローカルディスク (または EBS) よりもわずかに遅くなる場合があります。

クラウドシャッフルストレージプラグインを使用するための前提条件

AWS Glue ETL ジョブにクラウドシャッフルストレージプラグインを使用するには、以下が必要です。

- 途中のシャッフルデータやスピルしたデータを保存するための、ジョブが実行されているのと同じリージョンにある Amazon S3 バケット。シャッフルストレージの Amazon S3 プレフィックスは、`--conf spark.shuffle.glue.s3ShuffleBucket=s3://shuffle-bucket/prefix/`で次の例のように指定できます。

```
--conf spark.shuffle.glue.s3ShuffleBucket=s3://glue-shuffle-123456789-us-east-1/glue-shuffle-data/
```

- シャッフルマネージャーはジョブの終了後にファイルをクリーンアップしないため、プレフィックス (`glue-shuffle-data` など) に Amazon S3 ストレージライフサイクルポリシーを設定します。途中のシャッフルデータやスピルしたデータは、ジョブの終了後に削除する必要があります。ユーザーはプレフィックスに短いライフサイクルポリシーを設定できます。Amazon S3 ライフサイクルポリシーのセットアップ手順については、「Amazon Simple Storage Service ユーザーガイド」の「[バケットのライフサイクル設定の指定](#)」を参照してください。

AWS コンソールからの AWS Glue Spark シャッフルマネージャーの使用

ジョブを設定するときに AWS Glue コンソールまたは AWS Glue Studio を使用して AWS Glue Spark シャッフルマネージャーをセットアップするには、`--write-shuffle-files-to-s3` ジョブパラメータを選択して、ジョブの Amazon S3 シャッフルを有効にします。

Job parameters

Key	Value - optional	
<input type="text" value="Q --write-shuffle-files- X"/>	<input type="text" value="Q"/>	<input type="button" value="Remove"/>

You can add 49 more parameters.

AWS Glue Spark シャッフルプラグインの使用

次のジョブパラメータで、AWS Glue シャッフルマネージャーが有効になり、チューニングされます。これらのパラメータはフラグなので、指定された値は考慮されません。

- `--write-shuffle-files-to-s3` – メインフラグです。Amazon S3 バケットを使用してシャッフルデータの書き込みと読み込みを行う AWS Glue Spark シャッフルマネージャーが有効になります。フラグを指定しない場合、シャッフルマネージャーは使用されません。
- `--write-shuffle-spills-to-s3` – (AWS Glue バージョン 2.0 でのみサポートされています)。オプションのフラグです。スピルファイルを Amazon S3 バケットにオフロードできます。これにより、Spark ジョブの耐障害性が強化されます。これは、大量のデータをディスクに退避させる大規模なワークロードにのみ必要です。フラグが指定されていない場合、中間スピルファイルは書き込まれません。
- `--conf spark.shuffle.glue.s3ShuffleBucket=s3://<shuffle-bucket>` – 別のオプションのフラグです。シャッフルファイルを書き込む Amazon S3 バケットを指定します。デフォルトでは、`--TempDir/shuffle data` です。AWS Glue 3.0 以降では、`--conf spark.shuffle.glue.s3ShuffleBucket=s3://shuffle-bucket-1/prefix,s3://shuffle-bucket-2/prefix` のようにバケットをカンマ区切りで指定することで、シャッフルファイルを複数のバケットに書き込むことができます。複数のバケットを使用するとパフォーマンスが向上します。

シャッフルデータの保存時の暗号化を有効にするには、セキュリティ構成設定を提供する必要があります。セキュリティの構成の詳細については、「[the section called “暗号化のセットアップ”](#)」を参照してください。AWS Glue は、Spark が提供する他のすべてのシャッフル関連の構成をサポートします。

クラウドシャッフルストレージプラグインのソフトウェアバイナリ

また、Apache 2.0 ライセンスで Cloud Shuffle Storage Plugin for Apache Spark のソフトウェアバイナリをダウンロードして、任意の Spark 環境で実行することもできます。新しいプラグインは Amazon S3 をそのまま使えるようにサポートしており、[Google クラウドストレージ](#)や [Microsoft Azure Blob ストレージ](#)などの他の形式のクラウドストレージを使用するように簡単に設定することもできます。詳細については、「[Cloud Shuffle Storage Plugin for Apache Spark](#)」を参照してください。

注意事項と制限事項

AWS Glue シャッフルマネージャーには次の注意事項や制限事項があります。

- AWS Glue シャッフルマネージャーは、ジョブの完了後に Amazon S3 バケットに保存されている (一時的な) シャッフルデータファイルを自動的に削除しません。データを確実に保護するには、Cloud Shuffle ストレージプラグインを有効にする前に、[クラウドシャッフルストレージプラグインを使用するための前提条件](#) の手順に従ってください。
- データに偏りがある場合、この機能が使用できます。

Cloud Shuffle Storage Plugin for Apache Spark

Cloud Shuffle ストレージプラグインは [ShuffleDataIO API](#) と互換性のある Apache Spark プラグインで、シャッフルデータをクラウドストレージシステム (Amazon S3 など) に保存することを可能にします。この機能では、Spark アプリケーションにおいて join、reduceByKey、groupByKey、および repartition などの変換によってトリガーされることが多い、大規模なシャッフル操作のローカルディスクストレージ容量を補充または置換できます。このため、サーバーレスデータ分析ジョブやパイプラインでの、一般的な障害やコストパフォーマンスの不整合を軽減するのに役立ちます。

AWS Glue

AWS Glue バージョン 3.0 と 4.0 には、このプラグインが事前インストールされており、特別な手順なしで Amazon S3 でシャッフルが利用可能になります。Spark アプリケーションでこの機能を有効にする方法については、「[AWS Glue シャッフルマネージャーと Amazon S3](#)」を参照してください。

その他の Spark 環境

他の Spark 環境では、このプラグインに以下の Spark 構成を設定する必要があります。

- `--conf spark.shuffle.sort.io.plugin.class=com.amazonaws.spark.shuffle.io.cloud.Chopper`
この Spark に対し、Shuffle IO 用としてこのプラグインを使用することを通知します。
- `--conf spark.shuffle.storage.path=s3://bucket-name/shuffle-file-dir`: シャッフルファイルが保存されるパスです。

Note

このプラグインは、Spark コアクラスの 1 つを上書きします。そのため、プラグインの jar ファイルは、Spark jar に先行してロードする必要があります。プラグインを AWS Glue の外

部で使用する場合は、オンプレミスの YARN 環境で `userClassPathFirst` を使用して、この処理を実行できます。

Spark アプリケーションへのプラグインのバンドル

Spark アプリケーションをローカルで開発する際、Maven `pom.xml` にプラグインの依存関係を追加することで、Spark アプリケーションと Spark ディストリビューション (バージョン 3.1 以降) にプラグインをバンドルできます。このプラグインと Spark バージョンの詳細については、「[プラグインのバージョン](#)」を参照してください。

```
<repositories>
  ...
  <repository>
    <id>aws-glue-etl-artifacts</id>
    <url>https://aws-glue-etl-artifacts.s3.amazonaws.com/release/ </url>
  </repository>
</repositories>
...
<dependency>
  <groupId>com.amazonaws</groupId>
  <artifactId>chopper-plugin</artifactId>
  <version>3.1-amzn-LATEST</version>
</dependency>
```

別の方法として、以下のように AWS Glue Maven アーティファクトからバイナリを直接ダウンロードして、Spark アプリケーションにインクルードすることも可能です。

```
#!/bin/bash
sudo wget -v https://aws-glue-etl-artifacts.s3.amazonaws.com/release/com.amazonaws/
chopper-plugin/3.1-amzn-LATEST/chopper-plugin-3.1-amzn-LATEST.jar -P /usr/lib/spark/
jars/
```

spark-submit の例

```
spark-submit --deploy-mode cluster \
--conf spark.shuffle.storage.s3.path=s3://<ShuffleBucket>/<shuffle-dir> \
--conf spark.driver.extraClassPath=<Path to plugin jar> \
--conf spark.executor.extraClassPath=<Path to plugin jar> \
--class <your test class name> s3://<ShuffleBucket>/<Your application jar> \
```

オプションの設定

これらは Amazon S3 シャッフルの動作を制御するための、オプションの設定値です。

- `spark.shuffle.storage.s3.enableServerSideEncryption`: ファイルのシャッフルおよびスピルのために、S3 SSE を有効化/無効化します。デフォルト値は `true` です。
- `spark.shuffle.storage.s3.serverSideEncryption.algorithm`: 使用される SSE アルゴリズム。デフォルト値は `AES256` です。
- `spark.shuffle.storage.s3.serverSideEncryption.kms.key`: `aws:kms` が有効になっている場合の KMS キーの ARN。

これらの構成に加え、ユースケースに応じて適切な暗号化が適用されるように `spark.hadoop.fs.s3.enableServerSideEncryption`、その他の環境固有の構成などの設定が必要になる場合があります。

プラグインのバージョン

このプラグインは、AWS Glue のバージョンごとに関連付けられている Spark バージョンで、それぞれサポートされています。次の表に、AWS Glue のバージョン、Spark バージョン、およびプラグインのソフトウェアバイナリ用の Amazon S3 ロケーションと関連付けたプラグインバージョンを示します。

AWS Glue バージョン	Spark バージョン	プラグインバージョン	Amazon S3 ロケーション
3.0	3.1	3.1-amzn-LATEST	s3://aws-glue-etl-artifacts/release/com/amazonaws/chopper-plugin/3.1-amzn-0/chopper-plugin-3.1-amzn-LATEST.jar
4.0	3.3	3.3-amzn-LATEST	s3://aws-glue-etl-artifacts/release/com/amazonaws/chopper-plugin/3.3-amzn-0/chopper-plugin-3.3-amzn-LATEST.jar

ライセンス

このプラグインのソフトウェアバイナリは、Apache-2.0 License の下でライセンスされています。

モニタリングAWS GlueSpark ジョブ

トピック

- [Spark メトリクスは以下で利用可能です AWS Glue Studio](#)
- [Apache Spark ウェブ UI を使用したジョブのモニタリング](#)
- [AWS Glue ジョブ実行インサイトでのモニタリング](#)
- [Amazon CloudWatch によるモニタリング](#)
- [ジョブのモニタリングとデバッグ](#)

Spark メトリクスは以下で利用可能です AWS Glue Studio

[Metrics] (メトリクス) タブには、ジョブが実行され実行とプロファイリングが有効になっているときに収集されるメトリクスが表示されます。以下のグラフが Spark ジョブに表示されます。

- ETL データ移動
- メモリプロファイル: ドライバーおよびエグゼキューター

[View additional metrics (追加のメトリクスの表示)] を選択し、以下のグラフを表示します。

- ETL データ移動
- メモリプロファイル: ドライバーおよびエグゼキューター
- エグゼキューター間のデータシャッフル
- CPU 負荷: ドライバーおよびエグゼキューター
- ジョブの実行: アクティブなエグゼキューター、完了したステージ、必要な最大エグゼキューター

ジョブがメトリクスを収集するように設定されている場合、CloudWatch これらのグラフのデータはメトリクスにプッシュされます。メトリクスを有効にしてグラフを解釈する方法の詳細については、「[ジョブのモニタリングとデバッグ](#)」を参照してください。

Example ETL データ移動グラフ

ETL データ移動グラフには、以下のメトリクスが表示されます。

- すべてのエグゼキュターにより Amazon S3 から読み取られたバイト数 – [glue.ALL.s3.filesystem.read_bytes](#)
- すべてのエグゼキュターにより Amazon S3 に書き込まれたバイト数 – [glue.ALL.s3.filesystem.write_bytes](#)

Jobs > e2e-straggler

Detailed job metrics

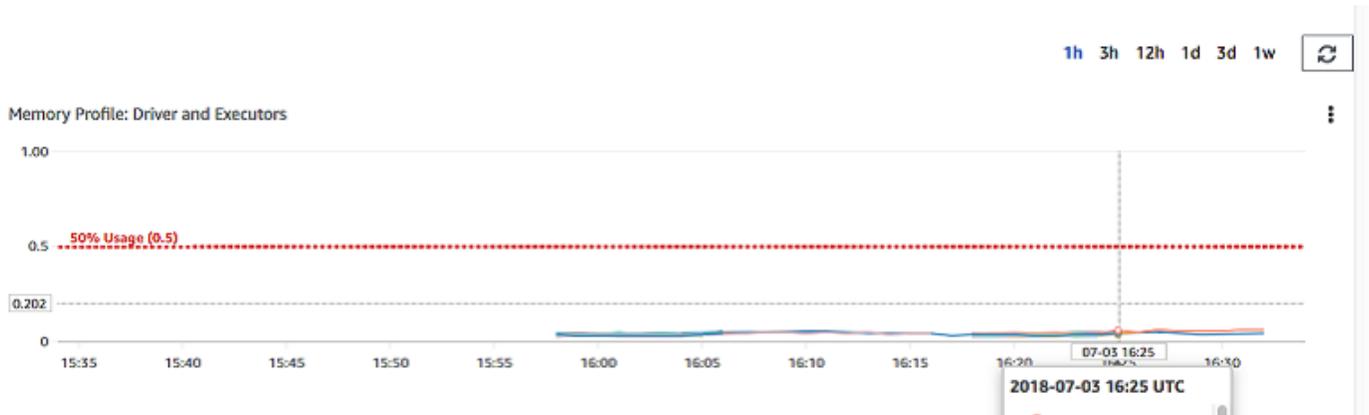
ETL Data Movement



Example メモリプロファイルグラフ

メモリプロファイルグラフには、以下のメトリクスが表示されます。

- このドライバーの JVM ヒープにより使用されるメモリの割合を、スケール 0~1 で、ドライバーごと、`executorId` により識別されるエグゼキュターごと、またはすべてのエグゼキュターごとに表示 –
 - [glue.driver.jvm.heap.usage](#)
 - [glue.executorId.jvm.heap.usage](#)
 - [glue.ALL.jvm.heap.usage](#)



Example エグゼキューター間のデータシャッフルグラフ

エグゼキューター間でのデータシャッフルグラフには、以下のメトリクスが表示されます。

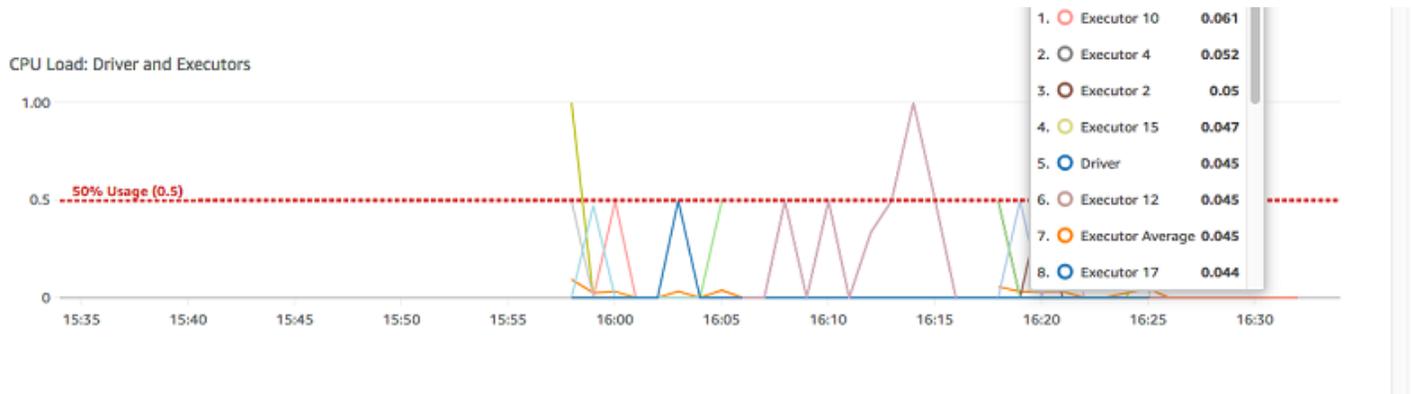
- すべてのエグゼキューター間でデータをシャッフルするためにそれらのエグゼキューターにより読み取られたバイト数 — [glue.driver.aggregate.shuffleLocalBytesRead](#)
- すべてのエグゼキューター間でデータをシャッフルするためにそれらのエグゼキューターにより書き込まれたバイト数 — [glue.driver.aggregate.shuffleBytesWritten](#)



Example CPU ロードグラフ

CPU ロードグラフには、以下のメトリクスが表示されます。

- 使用された CPU システムロードの割合を、スケール 0~1 で、ドライバーごと、executorId により識別されるエグゼキューターごと、またはすべてのエグゼキューターごとに表示 —
 - [glue.driver.system.cpuSystemLoad](#)
 - [glue.executorId.system.cpuSystemLoad](#)
 - [glue.ALL.system.cpuSystemLoad](#)



Example ジョブの実行グラフ

ジョブの実行グラフには、以下のメトリクスが表示されます。

- アクティブに実行されているエグゼキュターの数 — [glue.driver.ExecutorAllocationManager.executors.numberAllExecutors](#)
- 完了したステージの数 — [glue.aggregate.numCompletedStages](#)
- 必要なエグゼキュターの最大数 — [glue.driver.ExecutorAllocationManager.executors.numberMaxNeededExecutors](#)



Apache Spark ウェブ UI を使用したジョブのモニタリング

Apache Spark ウェブ UI を使用して、AWS Glue ジョブシステムで実行されている AWS Glue ETL ジョブと、AWS Glue 開発エンドポイントで実行されている Spark アプリケーションをモニタリングおよびデバッグできます。Spark UI では、ジョブごとに以下の項目を確認できます。

- 各 Spark ステージのイベントタイムライン
- ジョブの Directed Acyclic Graph (DAG)
- SparkSQL クエリの物理プランと論理プラン
- 各ジョブの基盤となる Spark 環境変数

Spark Web UI の使用の詳細については、Spark ドキュメントの「[Web UI](#)」を参照してください。Spark UI の結果を解釈してジョブのパフォーマンスを向上させる方法のガイダンスについては、AWS 「[規範ガイダンス](#)」の「[Apache Spark ジョブのパフォーマンスチューニング AWS Glue のベストプラクティス](#)」を参照してください。

AWS Glue コンソールで Spark UI を確認できます。これは、AWS Glue ジョブが AWS Glue 3.0 以降のバージョンで実行され、新しいジョブのデフォルトである標準 (レガシーではなく) 形式で生成されたログがある場合に使用できます。0.5 GB を超えるログファイルがある場合は、AWS Glue 4.0 以降のバージョンで実行されるジョブのローリングログサポートを有効にして、ログのアーカイブ、分析、トラブルシューティングを簡素化できます。

AWS Glue コンソールまたは AWS Command Line Interface () を使用して、Spark UI を有効にできます。AWS CLI。Spark UI を有効にすると、AWS Glue ETL ジョブと、AWS Glue 開発エンドポイントの Spark アプリケーションの Spark イベントログを Amazon Simple Storage Service (Amazon S3) の指定した場所にバックアップできます。Amazon S3 にバックアップされたイベントログは、ジョブの実行中 (リアルタイム) でも、ジョブの完了後でも Spark UI で使用できます。ログが Amazon S3 に残っている間は、AWS Glue コンソールの Spark UI でログを表示できます。

アクセス許可

AWS Glue コンソールで Spark UI を使用するには、すべての個々のサービス APIs を使用 `UseGlueStudio` または追加できます。Spark UI を完全に使用するには、すべての API が必要です。ただし、ユーザーはきめ細かくアクセスするために、IAM 許可にサービス API を追加して Spark UI の機能にアクセスできます。

`RequestLogParsing` はログの解析を行うため、最も重要です。残りの API は、それぞれの解析済みデータを読み取るためのものです。たとえば、`GetStages` は Spark ジョブのすべてのステージに関するデータへのアクセスを提供します。

`UseGlueStudio` にマッピングされた Spark UI サービス API のリストは、以下のサンプルポリシーに記載されています。以下のポリシーは、Spark UI 機能のみを使用するためのアクセスを提供します。Amazon S3 や IAM などのアクセス許可を追加するには、「[の カスタム IAM ポリシーの作成](#)」を参照してください [AWS Glue Studio](#)。

`UseGlueStudio` にマッピングされた Spark UI サービス API のリストは、以下のサンプルポリシーに記載されています。Spark UI サービス API を使用するとき、`glue:<ServiceAPI>` という名前空間を使用してください。

```
{
  "Version": "2012-10-17",
```

```
"Statement": [
  {
    "Sid": "AllowGlueStudioSparkUI",
    "Effect": "Allow",
    "Action": [
      "glue:RequestLogParsing",
      "glue:GetLogParsingStatus",
      "glue:GetEnvironment",
      "glue:GetJobs",
      "glue:GetJob",
      "glue:GetStage",
      "glue:GetStages",
      "glue:GetStageFiles",
      "glue:BatchGetStageFiles",
      "glue:GetStageAttempt",
      "glue:GetStageAttemptTaskList",
      "glue:GetStageAttemptTaskSummary",
      "glue:GetExecutors",
      "glue:GetExecutorsThreads",
      "glue:GetStorage",
      "glue:GetStorageUnit",
      "glue:GetQueries",
      "glue:GetQuery"
    ],
    "Resource": [
      "*"
    ]
  }
]
```

制限事項

- AWS Glue コンソールの Spark UI は、レガシーログ形式であるため、2023 年 11 月 20 日より前に発生したジョブ実行では使用できません。
- AWS Glue コンソールの Spark UI は、ストリーミングジョブでデフォルトで生成されるログなど、AWS Glue 4.0 のローリングログをサポートします。生成されたすべてのロールされたログファイルの最大合計は 2 GB です。ロールログがサポートされていない AWS Glue ジョブの場合、SparkUI でサポートされるログファイルの最大サイズは 0.5 GB です。
- Serverless Spark UI は、VPC のみがアクセスできる Amazon S3 バケットに保存されている Spark イベントログでは使用できません。

例: Apache Spark Web UI

この例では、Spark UI を使用してジョブのパフォーマンスを理解する方法を示します。 スクリーンショットは、自己管理型の Spark 履歴サーバーが提供する Spark ウェブ UI を示しています。AWS Glue コンソールの Spark UI も同様のビューを提供します。Spark Web UI の使用の詳細については、Spark ドキュメントの「[Web UI](#)」を参照してください。

次の例に示す Spark アプリケーションでは、2 つのデータソースから読み取り、結合変換を実行し、それを Parquet 形式で Amazon S3 に書き込みます。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from pyspark.sql.functions import count, when, expr, col, sum, isnull
from pyspark.sql.functions import countDistinct
from awsglue.dynamicframe import DynamicFrame

args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session

job = Job(glueContext)
job.init(args['JOB_NAME'])

df_persons = spark.read.json("s3://awsglue-datasets/examples/us-legislators/all/
persons.json")
df_memberships = spark.read.json("s3://awsglue-datasets/examples/us-legislators/all/
memberships.json")

df_joined = df_persons.join(df_memberships, df_persons.id == df_memberships.person_id,
'fullouter')
df_joined.write.parquet("s3://aws-glue-demo-sparkui/output/")

job.commit()
```

次の DAG 可視化は、この Spark ジョブのさまざまなステージを示しています。

APACHE **Spark** 2.2.1 tape-sparksql-jr_80b2f86d42bfb62... application UI

Jobs Stages Storage Environment Executors SQL

Details for Job 2

Status: SUCCEEDED
Completed Stages: 3

- ▶ Event Timeline
- ▼ DAG Visualization

▶ **Completed Stages (3)**

次に示すジョブのイベントタイムラインは、さまざまな Spark エグゼキュターの開始、実行、終了を示しています。



- Jobs
- Stages
- Storage
- Environment
- Executors
- SQL

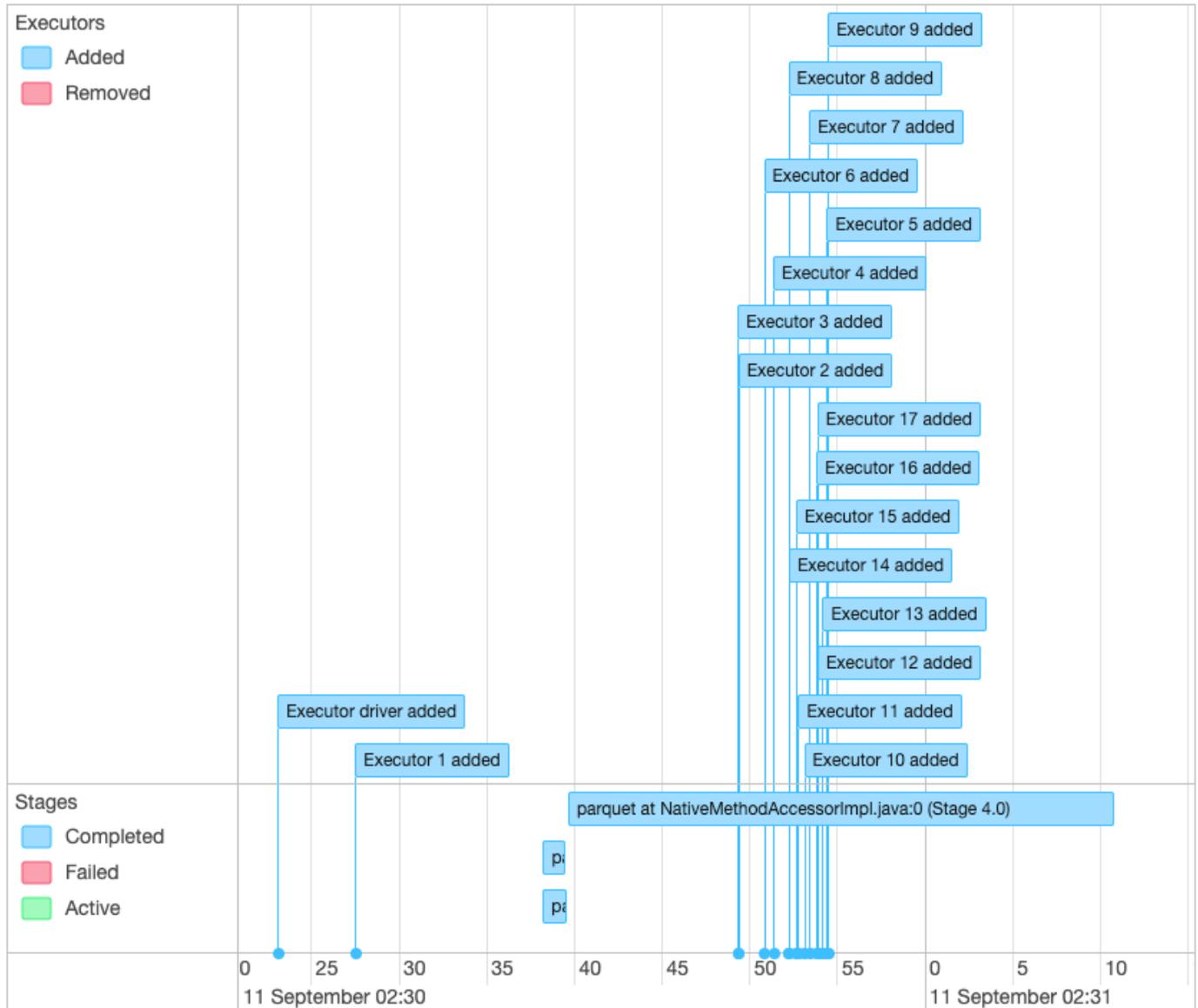
Details for Job 2

Status: SUCCEEDED

Completed Stages: 3

Event Timeline

Enable zooming



▶ DAG Visualization

▶ Completed Stages (3)

次の画面は、SparkSQL クエリプランの詳細を示しています。

- 解析された論理プラン
- 分析された論理プラン
- 最適化された論理プラン
- 実行のための物理プラン



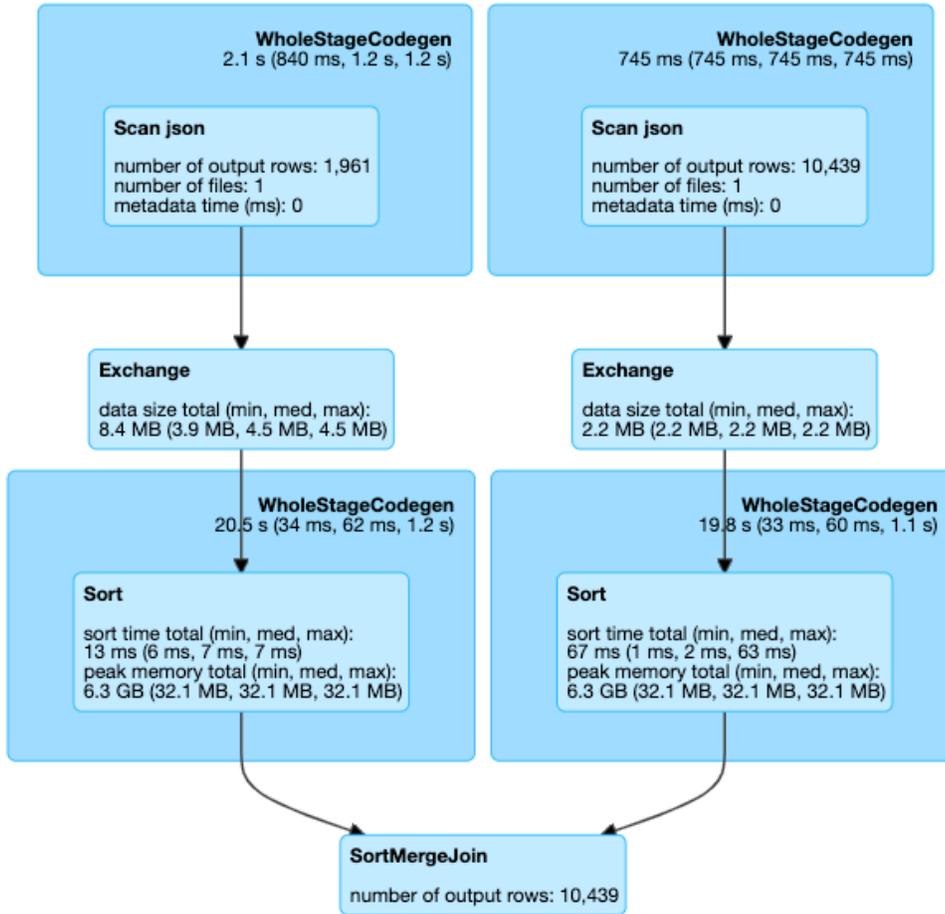
- Jobs
- Stages
- Storage
- Environment
- Executors
- SQL**

Details for Query 0

Submitted Time: 2019/09/11 02:30:37

Duration: 34 s

Succeeded Jobs: 2



▼ Details

```

== Parsed Logical Plan ==
Join FullOuter, (id#14 = person_id#50)
:-
Relation[birth_date#8,contact_details#9,death_date#10,family_name#11,gender#12,given_name#13,id#14,identifiers#15
,image#16,images#17,links#18,name#19,other_names#20,sort_name#21] json
+-
Relation[area_id#45,end_date#46,legislative_period_id#47,on_behalf_of_id#48,organization_id#49,person_id#50,role#
51,start_date#52] json

== Analyzed Logical Plan ==
birth_date: string, contact_details: array<struct<type:string,value:string>>, death_date: string, family_name:
string, gender: string, given_name: string, id: string, identifiers:
array<struct<identifier:string,scheme:string>>, image: string, images: array<struct<url:string>>, links:
array<struct<url:string>>, name: string, other_names:
array<struct<lang:string,name:string,note:string>>, sort_name: string, area_id: string, end_date: string,
legislative_period_id: string, on_behalf_of_id: string, organization_id: string, person_id: string, role: string,
start_date: string
Join FullOuter, (id#14 = person_id#50)

```

トピック

- [AWS Glue ジョブ用の Apache Spark ウェブ UI の有効化](#)
- [Spark 履歴サーバーの起動](#)

AWS Glue ジョブ用の Apache Spark ウェブ UI の有効化

Apache Spark ウェブ UI を使用して、AWS Glue ジョブシステムで実行されている AWS Glue ETL ジョブをモニタリングおよびデバッグできます。Spark UI は、AWS Glue コンソールまたは AWS Command Line Interface (AWS CLI) を使用して設定できます。

30 秒ごとに、AWS Glue は指定した Amazon S3 パスに Spark イベントログをバックアップします。

トピック

- [Spark UI の設定 \(コンソール\)](#)
- [Spark UI の設定 \(AWS CLI\)](#)
- [Notebook を使用するセッション用の Spark UI の設定](#)
- [ローリングログを有効にする](#)

Spark UI の設定 (コンソール)

以下のステップに従い、AWS Management Consoleを使用して Spark UI を設定します。AWS Glue ジョブを作成する場合、Spark UI はデフォルトで有効になっています。

ジョブを作成または編集するときに Spark UI をオンにするには

1. にサインイン AWS Management Console し、<https://console.aws.amazon.com/glue/> でAWS Glueコンソールを開きます。
2. ナビゲーションペインで ジョブを選択します。
3. [ジョブを追加] を選択するか、または既存のジョブを選択します。
4. [ジョブの詳細] で、[詳細プロパティ] を開きます。
5. [Spark UI] タブで、[Spark UI ログを Amazon S3 に書き込む] を選択します。
6. ジョブの Spark イベントログを保存するための Amazon S3 パスを指定します。ジョブでセキュリティ設定を使用する場合、暗号化は Spark UI ログファイルにも適用されることに注意してください。詳細については、「[AWS Glue によって書き込まれたデータの暗号化](#)」を参照してください。

7. [Spark UI のログ記録とモニタリングの設定] で、次のように操作します。

- AWS Glue コンソールで表示するログを生成する場合は、標準 を選択します。
- Spark 履歴サーバーで表示するログを生成している場合は、[レガシー] を選択します。
- 両方を生成することを選択することもできます。

Spark UI の設定 (AWS CLI)

Spark UI で表示するためのログを生成するには、AWS Glue コンソールで を使用して、次のジョブパラメータ AWS CLI をAWS Glueジョブに渡します。詳細については、「[the section called “ジョブのパラメータ”](#)」を参照してください。

```
'--enable-spark-ui': 'true',  
'--spark-event-logs-path': 's3://s3-event-log-path'
```

ログをレガシーの場所に配布するには、`--enable-spark-ui-legacy-path` パラメータを `"true"` に設定します。両方の形式でログを生成したくない場合は、`--enable-spark-ui` パラメータを削除します。

Notebook を使用するセッション用の Spark UI の設定

Warning

AWS Glue インタラクティブセッションは現在、コンソールで Spark UI をサポートしていません。Spark 履歴サーバーを設定します。

AWS Glue ノートブックを使用する場合は、セッションを開始する前に SparkUI 設定をセットアップします。そのためには、次のように `%%configure` セルマジックを使用します。

```
%%configure { "--enable-spark-ui": "true", "--spark-event-logs-path": "s3://path" }
```

ローリングログを有効にする

AWS Glue ジョブの SparkUI とローリングログファイルを有効にすると、いくつかの利点があります。

- ローリングログファイル – ローリングログファイルを有効にすると、はジョブ実行の各ステップに個別のログファイル AWS Glue を生成し、特定のステージまたは変換に固有の問題の特定とトラブルシューティングを容易にします。
- ログ管理の向上 – ローリングログファイルは、ログファイルをより効率的に管理するのに役立ちます。潜在的に大きなログファイルを 1 つ持つ代わりに、ログはジョブの実行ステージに基づいて、より小さく管理しやすいファイルに分割されます。これにより、ログのアーカイブ、分析、トラブルシューティングを簡素化できます。
- 耐障害性の向上 – AWS Glue ジョブが失敗または中断された場合、ローリングログイベントファイルは、最後に成功したステージに関する貴重な情報を提供するため、最初から開始するのではなく、その時点からジョブを再開することが容易になります。
- コストの最適化 – ローリングログファイルを有効にすることで、ログファイルに関連するストレージコストを節約できます。潜在的に大きなログファイルを 1 つ保存する代わりに、より小さく、管理しやすいログファイルを保存します。これは、特に長時間実行されるジョブや複雑なジョブの場合に、よりコスト効率が高くなります。

新しい環境では、ユーザーは以下を通じてローリングログを明示的に有効にできます。

```
'-conf': 'spark.eventLog.rolling.enabled=true'
```

または

```
'-conf': 'spark.eventLog.rolling.enabled=true -conf  
spark.eventLog.rolling.maxFileSize=128m'
```

ローリングログがアクティブ化されると、ロールオーバーする前にイベントログファイルの最大サイズ `spark.eventLog.rolling.maxFileSize` を指定します。このオプションパラメータを指定しない場合のデフォルト値は 128 MB です。最小は 10 MB です。

生成されたすべてのロールされたログファイルの最大合計は 2 GB です。ローリングログがサポートされていない AWS Glue ジョブの場合、SparkUI でサポートされるログファイルの最大サイズは 0.5 GB です。

追加の設定を渡すことで、ストリーミングジョブのローリングログをオフにすることができます。ログファイルが非常に大きいと、管理にコストがかかる場合があることに注意してください。

ローリングログをオフにするには、次の設定を提供します。

```
'--spark-ui-event-logs-path': 'true',
```

```
'--conf': 'spark.eventLog.rolling.enabled=false'
```

Spark 履歴サーバーの起動

Spark 履歴サーバーを使用して、独自のインフラストラクチャ上で Spark ログを視覚化できます。(レガシーではなく) 標準形式で生成されたログで、AWS Glue 4.0 以降のバージョンで AWS Glue ジョブが実行された場合、AWS Glue コンソールで同じ視覚化を確認することができます。詳細については、「[the section called “Spark UI を使用したモニタリング”](#)」を参照してください。

EC2 インスタンスでサーバーをホストする AWS CloudFormation テンプレートを使用して Spark 履歴サーバーを起動するか、Docker を使用してローカルで起動できます。

トピック

- [AWS CloudFormation を使用した Spark 履歴サーバーの起動と Spark UI の表示](#)
- [Docker を使用した Spark 履歴サーバーの起動と Spark UI の表示](#)

AWS CloudFormation を使用した Spark 履歴サーバーの起動と Spark UI の表示

AWS CloudFormation テンプレートを使用して Apache Spark 履歴サーバーを起動し、Spark ウェブ UI を表示できます。これらのテンプレートは、要件に応じて変更する必要があるサンプルです。

AWS CloudFormation を使用して Spark 履歴サーバーを起動し、Spark UI を表示するには

1. 次の表に示す起動スタックボタンの 1 つを選択します。これにより、AWS CloudFormation コンソールでスタックが起動されます。

リージョン	起動する
米国東部 (オハイオ)	
米国東部 (バージニア北部)	
米国西部 (北カリフォルニア)	
米国西部 (オレゴン)	

リージョン	起動する
アフリカ (ケープタウン)	Launch Stack 
アジアパシフィック (香港)	Launch Stack 
アジアパシフィック (ムンバイ)	Launch Stack 
アジアパシフィック (大阪)	Launch Stack 
アジアパシフィック (ソウル)	Launch Stack 
アジアパシフィック (シンガポール)	Launch Stack 
アジアパシフィック (シドニー)	Launch Stack 
アジアパシフィック (東京)	Launch Stack 
カナダ (中部)	Launch Stack 
欧州 (フランクフルト)	Launch Stack 
欧州 (アイルランド)	Launch Stack 
欧州 (ロンドン)	Launch Stack 
欧州 (ミラノ)	Launch Stack 
ヨーロッパ (パリ)	Launch Stack 
ヨーロッパ (ストックホルム)	Launch Stack 

リージョン	起動する
中東 (バーレーン)	
南米 (サンパウロ)	

2. [Specify template (テンプレートの指定)] ページで、[次へ] を選択します。
3. [Specify stack details (スタック詳細の指定)] ページで、[Stack name (スタック名)] に入力します。パラメータの下に追加情報を入力します。

a. Spark UI の設定

以下の情報を指定します。

- IP アドレス範囲 – Spark UI の表示に使用できる IP アドレス範囲。特定の IP アドレス範囲からのアクセスを制限する場合は、カスタム値を使用する必要があります。
- 履歴サーバーポート – Spark UI のポート。デフォルト値を使用できます。
- イベントログディレクトリ – AWS Glue ジョブまたは開発エンドポイントから保存された Spark イベントログの場所を選択します。イベントログのパススキームとして `s3a://` を使用する必要があります。
- Spark パッケージの場所 – デフォルト値を使用できます。
- キーストアパス – HTTPS の SSL/TLS キーストアパス。カスタムキーストアファイルを使用する場合は、ここで S3 パス `s3://path_to_your_keystore_file` を指定できます。このパラメータを空のままにすると、自己署名証明書ベースのキーストアが生成されて使用されます。
- キーストアパスワード – HTTPS 用の SSL/TLS キーストアパスワードを入力します。

b. EC2 インスタンスの設定

以下の情報を指定します。

- インスタンスタイプ – Spark 履歴サーバーをホストする Amazon EC2 インスタンスのタイプ。このテンプレートによってアカウントで Amazon EC2 インスタンスが起動されるため、Amazon EC2 のコストがアカウントに別個に課金されます。
- 最新の AMI ID – Spark 履歴サーバーインスタンスの Amazon Linux 2 の AMI ID。デフォルト値を使用できます。

- VPC ID – Spark 履歴サーバーインスタンスの Virtual Private Cloud (VPC) の ID。アカウントで利用可能な VPC のいずれかを使用できます。デフォルトの VPC と [デフォルトのネットワーク ACL](#) を使用することは推奨されません。詳細については、Amazon VPC ユーザーガイドの「[デフォルト VPC とデフォルトサブネット](#)」および「[VPC を作成する](#)」を参照してください。
 - サブネット ID – Spark 履歴サーバーインスタンスの ID。VPC 内の任意のサブネットを使用できます。クライアントからサブネットのネットワークにアクセスできる必要があります。インターネット経由でアクセスする場合は、ルートテーブルにインターネットゲートウェイがあるパブリックサブネットを使用する必要があります。
- c. [Next] を選択します。
4. [Configure stack options] ページで、CloudFormation がスタック内のリソースを作成、変更、または削除する方法を決定するために現在のユーザー認証情報を使用する場合は、[Next] をクリックします。アクセス許可セクションで、現在のユーザー権限の代わりに使用するロールを指定し、[次へ] を選択することもできます。
 5. [確認] ページで、テンプレートを確認します。

[I acknowledge that AWS CloudFormation might create IAM resources] (AWS CloudFormation によって IAM リソースが作成される場合があることを承認します) を選択し、[Create stack] (スタックの作成) を選択します。

6. スタックが作成されるまで待ちます。
7. [Outputs (出力)] タブをクリックします。
 - a. パブリックサブネットを使用している場合は、SparkUiPublicUrl の URL をコピーします。
 - b. プライベートサブネットを使用している場合は、SparkUiPrivateUrl の URL をコピーします。
8. ウェブブラウザを開き、URL を貼り付けます。これにより、指定したポートで HTTPS を使用してサーバーにアクセスできます。ブラウザがサーバーの証明書を認識しない場合があります。この場合、保護を上書きして続行する必要があります。

Docker を使用した Spark 履歴サーバーの起動と Spark UI の表示

ローカルアクセスを希望する (Apache Spark 履歴サーバーの EC2 インスタンスを使用しない) 場合は、Docker を使用して Apache Spark 履歴サーバーを起動し、Spark UI をローカルに表示することもできます。この Dockerfile は、要件に応じて変更する必要があるサンプルです。

前提条件

ラップトップに Docker をインストールする方法については、[Docker Engine community](#) を参照してください。

Docker を使用して Spark 履歴サーバーを起動し、Spark UI をローカルに表示するには

1. GitHub からファイルをダウンロードします。

Dockerfileと pom.xml を [[AWS Glueコードサンプル](#)]からダウンロードします。

2. アクセス AWS にユーザー認証情報を使用するか、フェデレーテッドユーザー認証情報を使用するかを決定します。

- 現在のユーザー認証情報を使用してアクセスを AWS にするには、AWS_ACCESS_KEY_ID と AWS_SECRET_ACCESS_KEY で使用する値を docker run コマンドで取得します。詳細については、「IAM ユーザーガイド」の「[IAM ユーザーのアクセスキーの管理](#)」を参照してください。
- SAML 2.0 フェデレーテッドユーザーを使用して AWS にアクセスするには、AWS_ACCESS_KEY_ID、AWS_SECRET_ACCESS_KEY、AWS_SESSION_TOKEN で値を取得します。詳細については、「[一時的なセキュリティ認証情報のリクエスト](#)」を参照してください。

3. docker run コマンドで使用するイベントログディレクトリの場所を決定します。
4. ローカルディレクトリ内のファイルを使用して、名前 glue/sparkui とタグ latest を指定して、Docker イメージを構築します。

```
$ docker build -t glue/sparkui:latest .
```

5. Docker コンテナを作成してスタートします。

次のコマンドでは、ステップ 2 および 3 で事前に取得した値を使用します。

- a. ユーザー認証情報を使用して docker コンテナを作成するには、次のようなコマンドを使用します。

```
docker run -itd -e SPARK_HISTORY_OPTS="$SPARK_HISTORY_OPTS -
Dspark.history.fs.logDirectory=s3a://path_to_eventlog
-Dspark.hadoop.fs.s3a.access.key=AWS_ACCESS_KEY_ID -
Dspark.hadoop.fs.s3a.secret.key=AWS_SECRET_ACCESS_KEY"
-p 18080:18080 glue/sparkui:latest "/opt/spark/bin/spark-class
org.apache.spark.deploy.history.HistoryServer"
```

- b. 一時的な認証情報を使用して Docker コンテナを作成するには、`org.apache.hadoop.fs.s3a.TemporaryAWSCredentialsProvider` をプロバイダとして指定し、ステップ 2 で取得した認証情報値を指定します。詳細については、「[一時的な AWSCredentialsProvider でのセッション資格情報の使用](#)」の Hadoop: Amazon Web Services との統合」のドキュメントを参照してください。

```
docker run -itd -e SPARK_HISTORY_OPTS="$SPARK_HISTORY_OPTS -
Dspark.history.fs.logDirectory=s3a://path_to_eventlog
-Dspark.hadoop.fs.s3a.access.key=AWS_ACCESS_KEY_ID -
Dspark.hadoop.fs.s3a.secret.key=AWS_SECRET_ACCESS_KEY
-Dspark.hadoop.fs.s3a.session.token=AWS_SESSION_TOKEN
-
Dspark.hadoop.fs.s3a.aws.credentials.provider=org.apache.hadoop.fs.s3a.TemporaryAWSCred
-p 18080:18080 glue/sparkui:latest "/opt/spark/bin/spark-class
org.apache.spark.deploy.history.HistoryServer"
```

Note

これらの設定パラメータは、[Hadoop-AWS モジュール](#)に由来します。場合によっては、ユースケースに基づいて特定の設定を追加する必要があります。たとえば、隔離されたリージョンのユーザーは、`spark.hadoop.fs.s3a.endpoint` を設定する必要があります。

6. ブラウザで `http://localhost:18080` を開き、Spark UI をローカルに表示します。

AWS Glue ジョブ実行インサイトでのモニタリング

AWS Glue ジョブ実行インサイトは、AWS Glue AWS Glue ジョブのデバッグとジョブの最適化を簡素化する機能です。AWS Glue [Spark UI](#)、[CloudWatch ジョブを監視するためのログとメトリクスを提供します](#)。AWS Glue この機能により、AWS Glue ジョブの実行に関する次の情報を取得できます。

- エラーが発生した AWS Glue ジョブスクリプトの行番号。
- ジョブの失敗の直前に Spark クエリプランで最後に実行された Spark アクション。
- 時系列ログストリームに表示される障害に関連する Spark 例外イベント。
- 根本原因の分析と、問題を解決するための推奨アクション (スクリプトのチューニングなど)。

- 根本原因に対処する推奨アクションを含む共通の Spark イベント (Spark アクションに関連するログメッセージ)。

これらのインサイトはすべて、AWS Glueジョブのログにある 2 CloudWatch つの新しいログストリームを使用して取得できます。

要件

AWS Glue AWS Glue ジョブ実行インサイト機能はバージョン 2.0、3.0、4.0 で使用できます。既存のジョブの[移行ガイド](#)に従って、古い AWS Glue バージョンからアップグレードすることができます。

AWS Glue ETL ジョブのジョブ実行インサイトを有効にする

ジョブ実行インサイトは、AWS Glue Studio または CLI によって有効にできます。

AWS Glue Studio

AWS Glue Studio でジョブを作成する場合、[ジョブの詳細] タブでジョブ実行インサイトを有効または無効にすることができます。「ジョブインサイトの生成」ボックスが選択されていることを確認します。

Requested number of workers

The number of workers you want AWS Glue to allocate to this job.

The maximums are 299 for G.1X and 149 for G.2X, and the minimum is 2.

Generate job insights

AWS Glue will analyze your job runs and provide insights on how to optimize your jobs and the reasons for job failures.

コマンドライン

CLI を使用してジョブを作成する場合は、単一の新しい[ジョブパラメータ](#)でジョブの実行を開始できます: `--enable-job-insights = true`。

デフォルトでは、ジョブ実行インサイトログストリームは、[AWS Glue連続ログ記録](#) (つまり `/aws-glue/jobs/logs-v2/`) で使用されるのと同じデフォルトロググループの下に作成されます。連続ログ記録に同じ引数のセットを使用して、カスタムロググループ名、ログフィルター、およびロググループ構成を設定できます。詳細については、「[AWS Glue ジョブの連続ログ記録の有効化](#)」を参照してください。

ジョブ実行インサイトログにアクセスすると、ログがストリーミングされます。CloudWatch

ジョブ実行インサイト機能を有効にすると、ジョブ実行が失敗したときに2つのログストリームが作成されることがあります。ジョブが正常に終了すると、いずれのストリームも生成されません。

1. 例外分析ログストリーム:<job-run-id>-job-insights-rca-driver。このストリームは、以下を提供します:
 - エラーが発生した AWS Glue ジョブスクリプトの行番号。
 - Spark クエリプラン (DAG) で最後に実行された Spark アクション。
 - Spark ドライバおよび例外に関連するエグゼキュータからの簡潔な時系列イベント。完全なエラーメッセージ、失敗した Spark タスク、およびそのエグゼキューター ID などの詳細を調べて、必要に応じてさらに詳細な調査を行うために、特定のエグゼキューターのログストリームに集中できます。
2. ルールベースのインサイトストリーム:
 - エラーの修正方法に関する根本原因の分析とレコメンデーション (特定のジョブパラメータを使用してパフォーマンスを最適化するなど)。
 - 関連する Spark イベントは、根本原因の分析および推奨アクションの基礎として機能します。

Note

最初のストリームは、失敗したジョブ実行で例外の Spark イベントが使用できる場合にのみ存在し、2番目のストリームは、失敗したジョブ実行に対して利用可能なインサイトがある場合にのみ存在します。例えば、ジョブが正常に終了した場合、どちらのストリームも生成されません。ジョブが失敗しても、失敗シナリオと一致するサービス定義のルールがない場合は、最初のストリームだけが生成されます。

AWS Glue Studio からジョブが作成された場合、上記のストリームへのリンクは、[ジョブ実行の詳細] タブ (ジョブ実行インサイト) の下で「簡潔で統合されたエラーログ」および「エラー分析とガイド」としても使用できます。

Job run - jr_ [REDACTED]

Run details [Info](#)

⊗ An error occurred while calling o134.pyWriteDynamicFrame. No such file or directory 's3://[REDACTED]'

Job name [REDACTED]	Run status ✔ Success	Glue version 2.0	Recent attempt 2
Start time May 17, 2021 1:10 PM	End time May 17, 2021 1:10 PM	Start-up time 4 seconds	Execution time 1 minute
Trigger name -	Last modified on May 17, 2021 1:10 PM	Security configuration -	Timeout 2880 minutes
Allocated capacity 10	Max capacity 10	Number of workers 10	Worker type G.1X
Cloudwatch logs All logs Output logs Error logs	Job run insights Info Concise and consolidated error logs Error analysis and guidance		

AWS Glue ジョブ実行インサイトの例

このセクションでは、ジョブ実行インサイト機能が、失敗したジョブの問題を解決するためにどのように役立つかの例を紹介します。この例では、ユーザーは AWS Glue ジョブに必要なモジュール (tensorflow) をインポートし、データに基づいて機械学習モデルを分析および構築するのを忘れていました。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from pyspark.sql.types import *
from pyspark.sql.functions import udf,col

args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
```

```
job = Job(glueContext)
job.init(args['JOB_NAME'], args)

data_set_1 = [1, 2, 3, 4]
data_set_2 = [5, 6, 7, 8]

scoresDf = spark.createDataFrame(data_set_1, IntegerType())

def data_multiplier_func(factor, data_vector):
    import tensorflow as tf
    with tf.compat.v1.Session() as sess:
        x1 = tf.constant(factor)
        x2 = tf.constant(data_vector)
        result = tf.multiply(x1, x2)
        return sess.run(result).tolist()

data_multiplier_udf = udf(lambda x: data_multiplier_func(x, data_set_2),
    ArrayType(IntegerType(), False))
factoredDf = scoresDf.withColumn("final_value", data_multiplier_udf(col("value")))
print(factoredDf.collect())
```

ジョブ実行インサイト機能がないと、ジョブが失敗すると、Spark によってスローされる次のメッセージのみが表示されます。

```
An error occurred while calling o111.collectToPython. Traceback (most recent call last):
```

メッセージがあいまいなので、デバッグ体験が制限されます。この場合、この機能は次の 2 CloudWatch つのログストリームで追加のインサイトを提供します。

1. job-insights-rca-driver ログストリーム:

- 例外イベント: このログストリームは、Spark ドライバおよびさまざまな分散ワーカーから収集された障害に関連する Spark 例外イベントを提供します。これらのイベントは、障害のあるコードが Spark タスク、エグゼキューター、および AWS Glue ワーカーに分散されたステージ全体で実行されるときに、例外の時系列の伝播を理解するのに役立ちます。
- 行番号: このログストリームは、失敗の原因となった欠落している Python モジュールをインポートするための呼び出しを行った 21 行目を識別します。また、スクリプトで最後に実行された行として、Spark アクション `collect()` の呼び出しである 24 行目も識別します。

Timestamp	Message
	No older events at this moment. Retry
2022-01-31T06:07:04.750-08:00	22/01/31 14:07:04 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Event: GlueExceptionAnalysisTaskFailed Failure Reason: Traceb...
2022-01-31T06:07:04.870-08:00	22/01/31 14:07:04 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Event: GlueExceptionAnalysisTaskFailed Stage ID: 0, Task ID: ...
2022-01-31T06:07:04.888-08:00	22/01/31 14:07:04 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Event: GlueExceptionAnalysisTaskFailed Stage ID: 0, Task ID: ...
2022-01-31T06:07:04.940-08:00	22/01/31 14:07:04 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Event: GlueExceptionAnalysisTaskFailed Stage ID: 0, Task ID: ...
2022-01-31T06:07:04.998-08:00	22/01/31 14:07:04 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Event: GlueExceptionAnalysisStageFailed Failure Reason: Job a...
2022-01-31T06:07:05.044-08:00	22/01/31 14:07:05 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Event: GlueExceptionAnalysisJobFailed Failure Reason: JobFail...
2022-01-31T06:07:05.105-08:00	22/01/31 14:07:05 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Root Cause Analysis Result: line 24 in script jobInsightsDemo... 22/01/31 14:07:05 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Root Cause Analysis Result: line 24 in script jobInsightsDemo.py.
	Copy
2022-01-31T06:07:05.427-08:00	22/01/31 14:07:05 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Event: GlueETLJobExceptionEvent Failure Reason: Traceback (mo...
2022-01-31T06:07:05.430-08:00	22/01/31 14:07:05 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Last Executed Line number from script jobInsightsDemo.py: 33 22/01/31 14:07:05 ERROR GlueExceptionAnalysisListener: [Glue Exception Analysis] Last Executed Line number from script jobInsightsDemo.py: 33
	Copy

2. job-insights-rule-driver ログストリーム:

- **根本原因とレコメンデーション** : スクリプトの障害の行番号と最後に実行された行番号に加えて、このログストリームには、根本原因の分析と、AWS Glue ドキュメントに従って、AWS Glue ジョブの追加の Python モジュールを使用するために必要なジョブパラメーターを設定するためのレコメンデーションが表示されます。
- **基本イベント** : このログストリームには、根本原因を推測してレコメンデーションを提供するために、サービス定義ルールで評価された Spark 例外イベントも表示されます。

2022-01-31T06:07:05.499-08:00	22/01/31 14:07:05 ERROR Analyzer: 2022-01-31 14:07:05,499 ERROR [pool-2-thread-1] app.GlueJobAnalyzerApp\$ (Logging.scala:logError(9)) - [Glue ... 22/01/31 14:07:05 ERROR Analyzer: 2022-01-31 14:07:05,499 ERROR [pool-2-thread-1] app.GlueJobAnalyzerApp\$ (Logging.scala:logError(9)) - [Glue Insights]
	<pre>{ "details": { "time": 1643638025489, "rootCauseAnalysis": "Module that is referenced in Glue job was not found.", "action": "Include all modules used in Glue job, refer documentation on how to include external modules, https://aws.amazon.com/premiumsupport/knowledge-center/glue-version2-external-python-libraries/" }, "cause": { "module": "data_multiplier_func", "issue": "ModuleNotFoundError: No module named 'tensorflow'", "fileName": "jobInsightsDemo.py", "lineOfCode": 24 }, "basis": [{ "event": { "timestamp": 1643638024940, "failureReason": "Traceback (most recent call last):\n File \"/opt/amazon/spark/python/lib/pyspark.zip/pyspark/worker.py", line 377, in main\n process()\n File \"/opt/amazon/spark/python/lib/pyspark.zip/pyspark/worker.py", line 372, in process\n serializer.dump_stream(func(split_index, iterator),\n outfile)\n File \"/opt/amazon/spark/python/lib/pyspark.zip/pyspark/serializers.py", line 345, in dump_stream\n self.serializer.dump_stream(self._batched(iterator), stream)\n File \"/opt/amazon/spark/python/lib/pyspark.zip/pyspark/serializers.py", line 141, in dump_stream\n for obj in iterator:\n File \"/opt/amazon/spark/python/lib/pyspark.zip/pyspark/serializers.py", line 334, in _batched\n for item in iterator:\n File\n \<string>", line 1, in <lambda>\n File \"/opt/amazon/spark/python/lib/pyspark.zip/pyspark/worker.py", line 85, in <lambda>\n return lambda *a: f(*a)\n File\n \"/opt/amazon/spark/python/lib/pyspark.zip/pyspark/util.py", line 99, in wrapper\n return f(*args, **kwargs)\n File \"/tmp/jobInsightsDemo.py", line 31, in\n <lambda>\n File \"/tmp/jobInsightsDemo.py", line 24, in data_multiplier_func\nModuleNotFoundError: No module named 'tensorflow'\n", "stackTrace": [{ "declaringClass": "data_multiplier_func", "methodName": "ModuleNotFoundError: No module named 'tensorflow'", "fileName": "/tmp/jobInsightsDemo.py", "lineNumber": 24 }] }] } }</pre>
	Copy

Amazon CloudWatch によるモニタリング

Amazon CloudWatch を使用して AWS Glue をモニタリングすることで、AWS Glue から raw データを収集し、ほぼリアルタイムの読み取り可能なメトリクスに加工することができます。これらの統計は 2 週間記録されるため、履歴情報にアクセスしてウェブアプリケーションまたはサービスの動作をよりの確に把握することができます。デフォルトでは、AWS Glue メトリクスデータは

CloudWatch に自動的に送信されます。詳細については、Amazon CloudWatch ユーザーガイドの「[Amazon CloudWatch とは](#)」、および「[AWS Glue のメトリクス](#)」を参照してください。

継続的なログ記録

AWS Glue は、AWS Glue ジョブのリアルタイムの連続ログ記録もサポートします。ジョブの連続ログ記録が有効になっている場合は、AWS Glue コンソールまたは CloudWatch コンソールダッシュボードでリアルタイムのログを表示できます。詳細については、「[AWS Glue ジョブの連続ログ記録](#)」を参照してください。

オブザーバビリティメトリクス

[ジョブのオブザーバビリティメトリクス] が有効になっている場合、ジョブの実行時に追加の Amazon CloudWatch メトリクスが生成されます。AWS Glue オブザーバビリティメトリクスを使用して、AWS Glue の内部で何が起きているかに関するインサイトを生成し、問題の優先順位付けと分析を改善できます。

トピック

- [Amazon CloudWatch メトリクスを使用した AWS Glue のモニタリング](#)
- [AWS Glue ジョブプロファイルでの Amazon CloudWatch アラームの設定](#)
- [AWS Glue ジョブの連続ログ記録](#)
- [AWS Glue オブザーバビリティメトリクスを使用したモニタリング](#)

Amazon CloudWatch メトリクスを使用した AWS Glue のモニタリング

AWS Glue ジョブプロファイラーを使用して AWS Glue オペレーションをプロファイルおよびモニタリングできます。AWS Glue ジョブから raw データが収集され、ほぼリアルタイムの読み取り可能なメトリクスに加工されて、Amazon CloudWatch に保存されます。これらの統計は CloudWatch に保持されて集計されるため、履歴情報にアクセスしてアプリケーションの動作をよりの確に把握できます。

Note

ジョブメトリクスを有効にし、CloudWatch カスタムメトリクスが作成されると、追加料金が発生することがあります。詳細については、「[Amazon CloudWatch の料金](#)」を参照してください。

AWS Glue メトリクスの概要

AWS Glue を操作するときに、CloudWatch にメトリクスが送信されます。これらのメトリクスは、AWS Glue コンソール (推奨される方法)、CloudWatch コンソールダッシュボード、または AWS Command Line Interface (AWS CLI) で表示できます。

AWS Glue コンソールダッシュボードを使ってメトリクスを表示するには

ジョブのメトリクスの概要または詳細なグラフを表示したり、ジョブの実行の詳細なグラフを表示したりできます。

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインで、[ジョブ実行モニタリング] を選択します。
3. [Job の実行] で [アクション] を選択して、現在実行中のジョブの停止、ジョブの表示、またはジョブのブックマークの巻き戻しを実行します。
4. ジョブを選択し、[実行の詳細を表示する] を選択すると、そのジョブ実行に関する追加情報が表示されます。

CloudWatch コンソールダッシュボードを使用してメトリクスを表示するには

メトリクスはまずサービスの名前空間ごとにグループ化され、次に各名前空間内のさまざまなディメンションの組み合わせごとにグループ化されます。

1. CloudWatch コンソール (<https://console.aws.amazon.com/cloudwatch/>) を開きます。
2. ナビゲーションペインで [Metrics (メトリクス)] を選択します。
3. [Glue] 名前空間を選択します。

AWS CLI を使ってメトリクスを表示するには

- コマンドプロンプトで、次のコマンドを使用します。

```
aws cloudwatch list-metrics --namespace Glue
```

メトリクスは 30 秒ごとに AWS Glue から CloudWatch に報告され、CloudWatch メトリクスダッシュボードは 1 分ごとにそれらを表示するように設定されています。AWS Glue メトリクスは、以

前に報告された値からデルタ値を表示します。必要に応じて、メトリクスのダッシュボードにより 30 秒の値が集計 (合計) され、直近 1 分間分の値を取得できます。

Spark ジョブの AWS Glue メトリクス動作

AWS Glue メトリクスは、スクリプトでの `GlueContext` の初期化時に有効になり、基本的には Apache Spark タスクの終了時にのみ更新されます。また、これまでに完了した Spark タスク全体の集計値を表します。

ただし、AWS Glue が CloudWatch に渡す Spark メトリクスは、一般的に報告された時点で現在の状態を表す絶対値です。それらは 30 秒ごとに AWS Glue から CloudWatch に報告され、一般的にメトリクスダッシュボードには直近 1 分間に受け取ったデータポイントの平均が表示されます。

すべての AWS Glue メトリクス名の前には、必ず次のいずれかの種類のプレフィックスが付きます。

- `glue.driver.` – 名前がこのプレフィックスで始まるメトリクスは、Spark ドライバーのすべてのエグゼキュターから集計された AWS Glue メトリクスか、Spark ドライバーに対応する Spark メトリクスを表します。
- `.executorIdglue.` – `executorId` は、特定の Spark エグゼキュターの番号です。ログに表示されているエグゼキュターに対応します。
- `glue.ALL.` – 名前がこのプレフィックスで始まるメトリクスは、すべての Spark エグゼキュターからの値を集計します。

AWS Glue のメトリクス

AWS Glue によって、プロファイルが作成され、次のメトリクスが 30 秒ごとに CloudWatch に送信され、AWS Glue メトリクスダッシュボードに 1 分に 1 回表示されます。

メトリクス	説明
<code>glue.driver.aggregate.bytesRead</code>	<p>すべてのエグゼキュターで実行され完了したすべての Spark タスクによるすべてのデータソースから読み取られたバイト数。</p> <p>有効なディメンション: <code>JobName</code> (AWS Glue Job の名前)、<code>JobRunId</code> (JobRun ID. または ALL)、<code>Type</code> (カウント)。</p>

メトリクス	説明
	<p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。</p> <p>単位: バイト</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none">読み込みバイト。ジョブ進行状況。JDBC データソース。ジョブブックマークの問題。ジョブ実行間の差異。 <p>このメトリクスは、<code>glue.ALL.s3.filesystem.read_bytes</code> メトリクスと同様に使用できます。ただし、このメトリクスは Spark タスクの最後に更新され、S3 以外のデータソースも取得するという点が異なります。</p>

メトリクス	説明
<code>glue.driver.aggregate.elapsedTime</code>	<p>ETL の経過時間 (ミリ秒単位) (ジョブのブートストラップ時間は含まれません)。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)。</p> <p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。</p> <p>単位: ミリ秒</p> <p>ジョブの実行に平均どのくらいの時間がかかるかを調べるために利用できます。</p> <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none">• ストラグラーにアラームを設定する。• ジョブ実行間の差異を測定する。

メトリクス	説明
<code>glue.driver.aggregate.numCompletedStages</code>	<p>ジョブの完了したステージの数</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)。</p> <p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。</p> <p>単位: 個</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none">• ジョブ進行状況。• 他のメトリクスとの関連におけるジョブ実行のステージごとのタイムライン。 <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none">• ジョブの実行における要求の厳しい段階を特定する。• ジョブ実行全体の中で関連するスパイク (要求の厳しいステージ) にアラームを設定する。

メトリクス	説明
<code>glue.driver.aggregate.numCompletedTasks</code>	<p>ジョブで完了したタスクの数。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)。</p> <p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。</p> <p>単位: 個</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none">• ジョブ進行状況。• ステージ内の並列性。

メトリクス	説明
<code>glue.driver.aggregate.numFailedTasks</code>	<p>失敗したタスクの数。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)。</p> <p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。</p> <p>単位: 個</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none">• ジョブのタスクが失敗する原因となるデータ異常。• ジョブのタスクが失敗する原因となるクラスター異常。• ジョブのタスクが失敗する原因となるスクリプト異常。 <p>データを使用して、データ、クラスター、スクリプトの異常を示唆する可能性のある障害の増加に対するアラームを設定できます。</p>

メトリクス	説明
<code>glue.driver.aggregate.numKilledTasks</code>	<p>強制終了したタスクの数。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)。</p> <p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。</p> <p>単位: 個</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none">タスクが強制終了される例外 (OOM) が発生するデータスキュー異常。タスクが強制終了される例外 (OOM) が発生するスクリプト異常。 <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none">データ異常を示す障害が増加した場合のアラームを設定する。クラスター異常を示す障害が増加した場合のアラームを設定する。スクリプト異常を示す障害が増加した場合のアラームを設定する。

メトリクス	説明
<code>glue.driver.aggregate.recordsRead</code>	<p>すべてのエグゼキューターで実行され完了したすべての Spark タスクによるすべてのデータソースから読み取られたレコードの数。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)。</p> <p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。</p> <p>単位: 個</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none">読み込まれたレコード。ジョブ進行状況。JDBC データソース。ジョブブックマークの問題。ジョブ実行の日々のスキュー。 <p>このメトリクスは、<code>glue.ALL.s3.filesystem.read_bytes</code> メトリクスと同様に使用できますが、このメトリクスは Spark タスクの最後に更新される点が異なります。</p>

メトリクス	説明
<code>glue.driver.aggregate.shuffleBytesWritten</code>	<p>前のレポート以降にデータをシャッフルするためにそれらのエグゼキューターにより書き込まれたバイト数 (直前の 1 分間にこの目的のために書き込まれたバイト数として AWS Glue メトリクスダッシュボードにより集計)。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)。</p> <p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。</p> <p>単位: バイト</p> <p>ジョブのデータシャッフル (大規模な結合、グループ化、再分割、合体) をモニタリングするために使用できます。</p> <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none">• さらに処理する前に、大きな入力ファイルを再分割または解凍します。• ホットキーを避けるために、データをより均一に再分割します。• 結合またはグループ化オペレーションの前にデータを事前にフィルタリングします。

メトリクス	説明
<code>glue.driver.aggregate.shuffleLocalBytesRead</code>	<p>前のレポート以降にデータをシャッフルするためにそれらのエグゼキューターにより読み取られたバイト数 (直前の 1 分間にこの目的のために読み取られたバイト数として AWS Glue メトリクスダッシュボードにより集計)。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)。</p> <p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。</p> <p>単位: バイト</p> <p>ジョブのデータシャッフル (大規模な結合、グループ化、再分割、合体) をモニタリングするために使用できます。</p> <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none">• さらに処理する前に、大きな入力ファイルを再分割または解凍します。• ホットキーを使用して、データをより均一に再分割します。• 結合またはグループ化オペレーションの前にデータを事前にフィルタリングします。

メトリクス	説明
<code>glue.driver.BlockManager.disk.diskSpaceUsed_MB</code>	<p>すべてのエグゼキュターで使用されたディスク領域のメガバイト数。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)。</p> <p>有効な統計: Average これは Spark メトリクスで、絶対値としてレポートされます。</p> <p>単位: メガバイト</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none">• キャッシュされた RDD パーティションを表すブロックに使用されるディスク領域。• 中間シャッフル出力を表すブロックに使用されるディスク領域。• ブロードキャストを表すブロックに使用されるディスク領域。 <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none">• ディスク使用量の増加によるジョブの障害を特定する。• スピルやシャッフルを引き起こす大きなパーティションを特定する。• これらの問題を解決するためにプロビジョニングされた DPU 容量を増やす。

メトリクス	説明
<pre>glue.driver.ExecutorAllocationManager.executors.numberAllExecutors</pre>	<p>アクティブに実行されているジョブエグゼキューターの数。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)。</p> <p>有効な統計: Average これは Spark メトリクスで、絶対値としてレポートされます。</p> <p>単位: 個</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none"> • ジョブアクティビティ。 • 高負荷エグゼキューター (少数のエグゼキューターのみが実行されている) • 現在のエグゼキューターレベルの並列処理。 <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none"> • クラスターの使用率が低い場合に、事前に大きな入力ファイルを再分割または解凍する。 • 一部が高負荷であるために起こるステージまたはジョブ実行の遅延を特定する。 • numberMaxNeededExecutors と比較して、より多くの DPU をプロビジョニングするためにバックログを理解する。

メトリクス	説明
<pre>glue.driver.ExecutorAllocationManager.executors.numberMaxNeededExecutors</pre>	<p>現在の負荷を満たすために必要な (アクティブに実行中および保留中の) ジョブエグゼキュターの最大数。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)。</p> <p>有効な統計: Maximum。これは Spark メトリクスで、絶対値としてレポートされます。</p> <p>単位: 個</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none"> • ジョブアクティビティ。 • DPU の容量または強制終了/失敗したエグゼキュターのために使用できないエグゼキュターが原因で、まだスケジューリングされていない現在のエグゼキュターレベルの並列処理と保留タスクのバックログ。 <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none"> • スケジューリングキューの保留中/バックログを特定する。 • 一部が高負荷であるために起こるステージまたはジョブ実行の遅延を特定する。 • numberAllExecutors と比較して、より多くの DPU をプロビジョニングするためにバックログを理解する。 • プロビジョニングされた DPU 容量を増やして、保留中のエグゼキュターのバックログを修正する。

メトリクス	説明
glue.driver.jvm.heap.usage	<p>ドライバー、executorId により識別されるエグゼキューター、またはすべてのエグゼキューターに対する、このドライバー (スケール: 0~1) の JVM ヒープにより使用されるメモリの割合。</p>
glue.executorId.jvm.heap.usage	<p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)。</p>
glue.ALL.jvm.heap.usage	<p>有効な統計: Average これは Spark メトリクスで、絶対値としてレポートされます。</p> <p>単位: パーセント</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none"> glue.driver.jvm.heap.usage 使用時のドライバーのメモリ不足状態 (OOM)。 glue.ALL.jvm.heap.usage 使用時のエグゼキューターのメモリ不足状態 (OOM)。 <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none"> メモリを消費するエグゼキューター ID とステージを特定する。 高負荷のエグゼキューター ID とステージを特定する。 ドライバーのメモリ不足状態 (OOM) を特定する。 エグゼキューターのメモリ不足状態 (OOM) を特定し、対応するエグゼキューター ID を取得して、エグゼキューターログからスタックトレースを取得できるようにする。 高負荷またはメモリ不足状態 (OOM) の原因となるデータのスキューがある可能性があるファイルまたはパーティションを特定します。

メトリクス	説明
glue.driver.jvm.heap.used	<p>ドライバー、executorId によって特定されたエグゼキュター、またはすべてのエグゼキュターのために JVM ヒープが使用したメモリバイト数。</p>
glue.executorId.jvm.heap.used	<p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)。</p>
glue.ALL.jvm.heap.used	<p>有効な統計: Average これは Spark メトリクスで、絶対値としてレポートされます。</p> <p>単位: バイト</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none"> • ドライバーのメモリ不足状態 (OOM)。 • エグゼキュターのメモリ不足状態 (OOM)。 <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none"> • メモリを消費するエグゼキュター ID とステージを特定する。 • 高負荷のエグゼキュター ID とステージを特定する。 • ドライバーのメモリ不足状態 (OOM) を特定する。 • エグゼキュターのメモリ不足状態 (OOM) を特定し、対応するエグゼキュター ID を取得して、エグゼキュターログからスタックトレースを取得できるようにする。 • 高負荷またはメモリ不足状態 (OOM) の原因となるデータのスキューがある可能性があるファイルまたはパーティションを特定します。

メトリクス	説明
<pre>glue.driver.s3.filesystem.read_bytes</pre>	<p>前のレポート以降に、ドライバー、executorId によって特定されるエグゼキューター、またはすべてのエグゼキューターが Amazon S3 から読み取ったバイト数 (直前 1 分間に読み取ったバイト数として AWS Glue メトリクスダッシュボードに集計されます)。</p>
<pre>glue.executorId.s3.filesystem.read_bytes</pre>	<p>有効なディメンション: JobName、JobRunId、および Type (ゲージ)。</p>
<pre>glue.ALL.s3.filesystem.read_bytes</pre>	<p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。AWS Glue メトリクスダッシュボードのカーブの下にある領域は、2 つの異なるジョブ実行によって読み取られたバイト数を視覚的に比較するために使用できます。</p> <p>単位: バイト</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none"> • ETL データ移動。 • ジョブ進行状況。 • ジョブのブックマークの問題 (処理されたデータ、再処理されたデータ、およびスキップされたデータ) • 読み取りと外部データソースからの取り込みとの比率の比較。 • ジョブ実行間の差異。 <p>結果データは、次の目的で使用できます。</p> <ul style="list-style-type: none"> • DPU 容量の計画。

メトリクス	説明
	<ul style="list-style-type: none">ジョブ実行およびジョブステージのデータ読み取りの急激な増加または落ち込みに対するアラームの設定

メトリクス	説明
<pre>glue.driver.s3.filesystem.write_bytes</pre> <pre>glue.executorId.s3.filesystem.write_bytes</pre> <pre>glue.ALL.s3.filesystem.write_bytes</pre>	<p>前のレポート以降に、ドライバー、executorId によって特定されるエグゼキューター、またはすべてのエグゼキューターが Amazon S3 に書き込んだバイト数 (直前 1 分間に書き込んだバイト数として AWS Glue メトリクスダッシュボードに集計されます)。</p> <p>有効なディメンション: JobName、JobRunId、および Type (ゲージ)。</p> <p>有効な統計: SUM。このメトリクスは、最後に報告された値からのデルタ値であるため、AWS Glue メトリクスダッシュボードでは、集計に SUM 統計が使用されます。AWS Glue メトリクスダッシュボードのカーブの下にある領域は、2 つの異なるジョブ実行によって書き込まれたバイト数を視覚的に比較するために使用できます。</p> <p>単位: バイト</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none"> • ETL データ移動。 • ジョブ進行状況。 • ジョブのブックマークの問題 (処理されたデータ、再処理されたデータ、およびスキップされたデータ) • 読み取りと外部データソースからの取り込みとの比率の比較。 • ジョブ実行間の差異。 <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none"> • DPU 容量の計画。

メトリクス	説明
	<ul style="list-style-type: none">ジョブ実行およびジョブステージのデータ読み取りの急激な増加または落ち込みに対するアラームの設定
<code>glue.driver.streaming.numRecords</code>	<p>マイクロバッチで受信されたレコードの数。このメトリクスは AWS Glue バージョン 2.0 以降の AWS Glue ストリーミングジョブにのみ使用できます。です。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)。</p> <p>有効な統計: Sum, Maximum, Minimum, Average, Percentile</p> <p>単位: 個</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none">読み込まれたレコード。ジョブ進行状況。

メトリクス	説明
<code>glue.driver.streaming.batchProcessingTimeInMs</code>	<p>バッチの処理に要する時間 (ミリ秒単位)。このメトリクスは AWS Glue バージョン 2.0 以降の AWS Glue ストリーミングジョブにのみ使用できます。です。</p> <p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)。</p> <p>有効な統計: Sum, Maximum, Minimum, Average, Percentile</p> <p>単位: 個</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none">• ジョブ進行状況。• スクリプトのパフォーマンス。

メトリクス	説明
glue.driver.system.cpuSystemLoad	<p>ドライバー、executorId により識別されるエグゼキューター、またはすべてのエグゼキューターに対する使用された CPU システムロードの割合 (スケール: 0 ~ 1)。</p>
glue.executorId.system.cpuSystemLoad	<p>有効なディメンション: JobName (AWS Glue Job の名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)。</p>
glue.ALL.system.cpuSystemLoad	<p>有効な統計: Average このメトリクスは絶対値としてレポートされます。</p> <p>単位: パーセント</p> <p>次のモニタリングに使用できます。</p> <ul style="list-style-type: none"> • ドライバー CPU 負荷。 • エグゼキューター CPU 負荷。 • ジョブ内の CPU 制約または IO 制約のエグゼキューターまたはステージの検出。 <p>データの使用方法をいくつか次に示します。</p> <ul style="list-style-type: none"> • IO メトリクス (読み取りバイト数/シャッフルバイト数、タスクの並列処理) および必要な最大エグゼキューター数のメトリクスとともに、DPU 容量を計画する。 • CPU/IO 制約比を確認する。これにより、CPU 使用率が低い分割可能なデータセットを使用して、長時間実行されるジョブに対して、再分割とプロビジョニングされた容量の増加が可能になります。

AWS Glue メトリクスのディメンション

AWS Glue メトリクスは AWS Glue 名前空間を使用し、以下のディメンションのメトリクスを提供しています。

ディメンション	説明
JobName	このディメンションで、特定の AWS Glue ジョブのすべてのジョブ実行のメトリクスがフィルタリングされます。
JobRunId	このディメンションで、JobRun ID によって特定の AWS Glue ジョブ実行のメトリクスがフィルタリングされるか、ALL となります。
Type	このディメンションで、count (集計した数) または gauge (ある時点での値) によってメトリクスがフィルタリングされます。

詳細については、『[Amazon CloudWatch ユーザーガイド](#)』を参照してください。

AWS Glue ジョブプロファイルでの Amazon CloudWatch アラームの設定

AWS Glue メトリクスは、Amazon CloudWatch でも利用できます。スケジュールされたジョブのどの AWS Glue メトリクスでもアラームを設定できます。

アラームを設定するいくつかの一般的なシナリオは次のとおりです。

- ジョブのメモリ不足 (OOM): メモリ使用率が AWS Glue ジョブのドライバーまたはエグゼキューターの通常平均を超えた場合にアラームを設定します。
- 散在したエグゼキューター: AWS Glue ジョブでエグゼキューターの数長時間一定のしきい値を下回ったときにアラームを設定します。
- データバックログまたは再処理: CloudWatch の数式を使用して、ワークフローの個別のジョブから得られたメトリクスを比較します。その後、生成された式の値 (あるジョブにより書き込まれたバイトと後続のジョブにより読み込まれたバイトの割合など) でアラームをトリガーできます。

アラームの設定に関する詳細については、[Amazon CloudWatch Events ユーザーガイド](#)の「[Create or Edit a CloudWatch Alarm](#)」を参照してください。

CloudWatch を使用したモニタリングとデバッグのシナリオについては、「[ジョブのモニタリングとデバッグ](#)」を参照してください。

AWS Glue ジョブの連続ログ記録

AWS Glue は、AWS Glue ジョブのリアルタイム連続ログ記録を提供します。ドライバーログ、エグゼキュターログ、および Apache Spark ジョブの進行状況バーを含む、Amazon CloudWatch のリアルタイムの Apache Spark ジョブログを表示できます。リアルタイムのログを表示すると、実行中のジョブについてよりの確に把握することができます。

AWS Glue ジョブを開始すると、Spark アプリケーションの実行開始後、(各エグゼキュターが終了するまで 5 秒ごとに) リアルタイムのログ記録情報が CloudWatch に送信されます。このログは、AWS Glue コンソールまたは CloudWatch コンソールダッシュボードで表示できます。

この連続ログ記録機能には、次の機能が含まれます。

- 連続ログ記録
- アプリケーション固有のメッセージを記録するカスタムスクリプトロガー
- コンソールの進行状況バーにより、現在の AWS Glue ジョブの実行ステータスを追跡する

連続ログ記録が AWS Glue バージョン 2.0 でサポートされる方法の詳細については、「[Running Spark ETL Jobs with Reduced Startup Times](#)」を参照してください。

IAM ロールがログを読み取るために CloudWatch ロググループまたはストリームへのアクセスを制限できます。アクセスの制限の詳細については、CloudWatch のドキュメントを「[CloudWatch Logs でのアイデンティティベースのポリシー \(IAM ポリシー\) の使用](#)」を参照してください。

Note

連続ログを有効にし、追加の CloudWatch ログイベントが作成されると、追加料金が発生することがあります。詳細については、「[Amazon CloudWatch の料金](#)」を参照してください。

トピック

- [AWS Glue ジョブの連続ログ記録の有効化](#)
- [AWS Glue ジョブの連続ログ記録の表示](#)

AWS Glue ジョブの連続ログ記録の有効化

AWS Glue コンソールを使用するか、AWS Command Line Interface (AWS CLI) を通じて、連続ログ記録を有効にできます。

新しいジョブを作成するか、既存のジョブを編集または AWS CLI を介して有効にするときに、連続ログ記録を有効にできます。

また、Amazon CloudWatch ロググループ名、AWS Glue ジョブ実行 ID ドライバー/エグゼキューター ID の前の CloudWatch ログストリームプレフィックス、ログメッセージのログ変換パターンなどのカスタム設定オプションを指定することもできます。これらの設定は、異なる有効期限ポリシーを持つカスタム CloudWatch ロググループに集約ログを設定し、カスタムログストリームプレフィックスと変換パターンを使用してさらに分析するのに役立ちます。

トピック

- [AWS Management Console の使用](#)
- [カスタムスクリプトロガーを使用したアプリケーション固有のメッセージのログ記録](#)
- [進行状況バーでのジョブ進行状況の表示](#)
- [連続ログ記録によるセキュリティ設定。](#)

AWS Management Console の使用

次のステップに従ってコンソールを使用して、AWS Glue ジョブを作成または編集するときに連続ログ記録を有効にします。

連続ログ記録を有効にして新しい AWS Glue ジョブを作成するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインで [ETL ジョブ] を選択します。
3. [ビジュアル ETL] を選択します。
4. [Job 詳細] タブで、[詳細プロパティ] セクションを展開します。
5. [継続的なログ記録] で [CloudWatch でログの有効化] を選択します。

既存の AWS Glue ジョブに対して連続ログ記録を有効にするには

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。

2. ナビゲーションペインで [Jobs] (ジョブ) を選択します。
3. [Jobs (ジョブ)] リストから既存のジョブを選択します。
4. [Action (アクション)], [Edit job (ジョブの編集)] の順に選択します。
5. [Job 詳細] タブで、[詳細プロパティ] セクションを展開します。
6. [継続的なログ記録] で [CloudWatch でログの有効化] を選択します。

AWS CLI の使用

連続ログ記録を有効にするには、AWS Glue ジョブにジョブパラメータを渡します。他の AWS Glue ジョブパラメータに類似した次の特殊なジョブパラメータを渡します。詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。

```
'--enable-continuous-cloudwatch-log': 'true'
```

カスタム Amazon CloudWatch ロググループ名を指定できます。指定しない場合、デフォルトのロググループ名は、/aws-glue/jobs/logs-v2/ になります。

```
'--continuous-log-logGroup': 'custom_log_group_name'
```

カスタム Amazon CloudWatch ログストリームプレフィックスを指定できます。指定しない場合、デフォルトのログストリームプレフィックスはジョブ実行 ID になります。

```
'--continuous-log-logStreamPrefix': 'custom_log_stream_prefix'
```

カスタムの連続ロギング変換パターンを指定できます。指定しない場合、デフォルトの変換パターンは %d{yy/MM/dd HH:mm:ss} %p %c{1}: %m%n になります。変換パターンは、ドライバーログとエグゼキューターログにのみ適用されます。AWS Glue の進行状況バーには影響しません。

```
'--continuous-log-conversionPattern': 'custom_log_conversion_pattern'
```

カスタムスクリプトロガーを使用したアプリケーション固有のメッセージのログ記録

AWS Glue ロガーを使用すると、ドライバーログストリームにリアルタイムで送信される、スクリプトのアプリケーション固有のメッセージを記録できます。

Python スクリプトの例を以下に示します。

```
from awsglue.context import GlueContext
```

```
from pyspark.context import SparkContext

sc = SparkContext()
glueContext = GlueContext(sc)
logger = glueContext.get_logger()
logger.info("info message")
logger.warn("warn message")
logger.error("error message")
```

Scala スクリプトの例を以下に示します。

```
import com.amazonaws.services.glue.log.GlueLogger

object GlueApp {
  def main(sysArgs: Array[String]) {
    val logger = new GlueLogger
    logger.info("info message")
    logger.warn("warn message")
    logger.error("error message")
  }
}
```

進行状況バーでのジョブ進行状況の表示

AWS Glue は、JOB_RUN_ID-progress-bar ログストリームでリアルタイムの進行状況バーを提供して AWS Glue ジョブ実行のステータスを確認します。現時点では、glueContext を初期化するジョブのみがサポートされています。glueContext を初期化せずに純粋な Spark ジョブを実行すると、AWS Glue 進行状況バーは表示されません。

進行状況バーでは、以下の進行状況が 5 秒ごとに更新されます。

```
Stage Number (Stage Name): > (numCompletedTasks + numActiveTasks) /
totalNumOfTasksInThisStage]
```

連続ログ記録によるセキュリティ設定。

CloudWatch ログのセキュリティ設定が有効になっている場合、AWS Glue では、連続ログに対して次のような名前のロググループが作成されます。

```
<Log-Group-Name>-<Security-Configuration-Name>
```

デフォルトのロググループおよびカスタムロググループは次のようになります。

- デフォルトの連続ロググループは `/aws-glue/jobs/logs-v2-<Security-Configuration-Name>` となります。
- カスタム連続ロググループは `<custom-log-group-name>-<Security-Configuration-Name>` となります。

CloudWatch Logs でセキュリティ設定を有効にする場合、`logs:AssociateKmsKey` を IAM ロールのアクセス許可に追加する必要があります。そのアクセス許可が含まれていない場合、連続ログ記録は無効になります。また、CloudWatch Logs 暗号化を設定するには、Amazon CloudWatch Logs ユーザーガイドの「[AWS Key Management Service を使用して CloudWatch Logs のログデータを暗号化する](#)」を参照してください。

セキュリティ設定作成の詳細については、「[AWS Glue コンソールでのセキュリティ設定の使用](#)」を参照してください。

AWS Glue ジョブの連続ログ記録の表示

リアルタイムのログは、AWS Glue コンソールまたは Amazon CloudWatch コンソールを使用して表示できます。

AWS Glue コンソールダッシュボードを使用してリアルタイムのログを表示するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインで ジョブを選択します。
3. 既存のジョブを追加または開始します。[アクション]、[ジョブの実行] の順に選択します。

ジョブの実行を開始するときは、実行中のジョブに関する情報が含まれているページに移動します。

- [ログ] タブには、集約された古いアプリケーションログが表示されます。
 - [Continuous logging (連続ログ記録)] タブには、`glueContext` が初期化されたジョブが実行されているときにリアルタイムの進行状況バーが表示されます。
 - [Continuous logging (連続ログ記録)] タブには、[Driver logs (ドライバーログ)] も含まれています。このログには、リアルタイムの Apache Spark ドライバーログと、ジョブが実行中に AWS Glue アプリケーションを使用して記録されたスクリプトからのアプリケーションログがキャプチャされます。
4. 古いジョブでは、リアルタイムログを [Job History] (ジョブ履歴) ビューで [Logs] (ログ) を選択して表示することもできます。このアクションを実行すると、そのジョブ実行のすべての Spark

ドライバー、エグゼキュター、および進捗状況バーのログストリームを表示する CloudWatch コンソールに移動します。

CloudWatch コンソールダッシュボードを使用してリアルタイムのログを表示するには

1. CloudWatch コンソール (<https://console.aws.amazon.com/cloudwatch/>) を開きます。
2. ナビゲーションペインで [ログ] を選択します。
3. `/aws-glue/jobs/logs-v2/` ロググループを選択します。
4. [フィルタ] ボックスに、ジョブ実行 ID を貼り付けます。

ドライバーログ、エグゼキュターログ、および進行状況バーを表示できます ([標準フィルタ] を使用している場合)。

AWS Glue オブザーバビリティメトリクスを使用したモニタリング

Note

AWS Glue オブザーバビリティメトリクスは AWS Glue 4.0 以降のバージョンで使用できません。

AWS Glue オブザーバビリティメトリクスを使用して、AWS Glue for Apache Spark ジョブの内部で何が起きているかに関するインサイトを生成し、問題の優先順位付けと分析を改善できます。オブザーバビリティメトリクスは Amazon CloudWatch ダッシュボードを通じて視覚化され、エラーの根本原因分析の実行に役立てたり、パフォーマンスのボトルネックを診断したりするために使用できます。問題のデバッグにかかる時間を大幅に削減できるため、より迅速かつ効果的に問題を解決することに注力できます。

AWS Glue オブザーバビリティは、次の 4 つのグループに分類された Amazon CloudWatch メトリクスを提供します。

- [信頼性] (エラークラス) – 対処する特定の時間範囲における最も一般的な障害の理由を簡単に特定します。
- [パフォーマンス] (歪み) – パフォーマンスのボトルネックを特定し、チューニング手法を適用します。例えば、ジョブの歪みを理由としてパフォーマンスの低下が発生した場合、Spark Adaptive Query Execution を有効にして、歪み結合のしきい値を微調整することをお勧めします。

- [スループット] (ソース/シンクあたりのスループット) – データの読み取りと書き込みの傾向をモニタリングします。異常についての Amazon CloudWatch アラームを設定することもできます。
- [リソースの使用率] (ワーカー、メモリ、ディスクの使用率) – キャパシティの使用率が低いジョブを効率的に見つけます。これらのジョブについて AWS Glue 自動スケーリングを有効にすることをお勧めします。

AWS Glue オブザーバビリティメトリクスの開始方法

Note

新しいメトリクスは、AWS Glue Studio コンソールでデフォルトで有効化されています。

AWS Glue Studio でオブザーバビリティメトリクスを設定するには:

1. AWS Glue コンソールにログインし、コンソールメニューから [ETL ジョブ] を選択します。
2. [自分のジョブ] セクションで、ジョブ名をクリックしてジョブを選択します。
3. [Job details] (ジョブの詳細) タブを選択します。
4. 一番下までスクロールして、[詳細プロパティ] を選択し、[ジョブのオブザーバビリティメトリクス] を選択します。

The screenshot shows the AWS Glue console interface for a job named "obs-test". The "Job details" tab is selected. Under the "Advanced properties" section, the "Job observability metrics" option is checked and highlighted with a red box. Other options like "Job metrics", "Continuous logging", "Spark UI", and "Serverless Spark UI" are also checked. The "Script filename" is "obs-test.py" and the "Script path" is "s3://aws-glue-assets-590186200215-us-east-1/scripts/".

AWS CLI を使用して AWS Glue オブザーバビリティメトリクスを有効にするには

- 入力 JSON ファイル内の次の key-value を `--default-arguments` マップに追加します。

```
--enable-observability-metrics, true
```

AWS Glue オブザーバビリティの使用

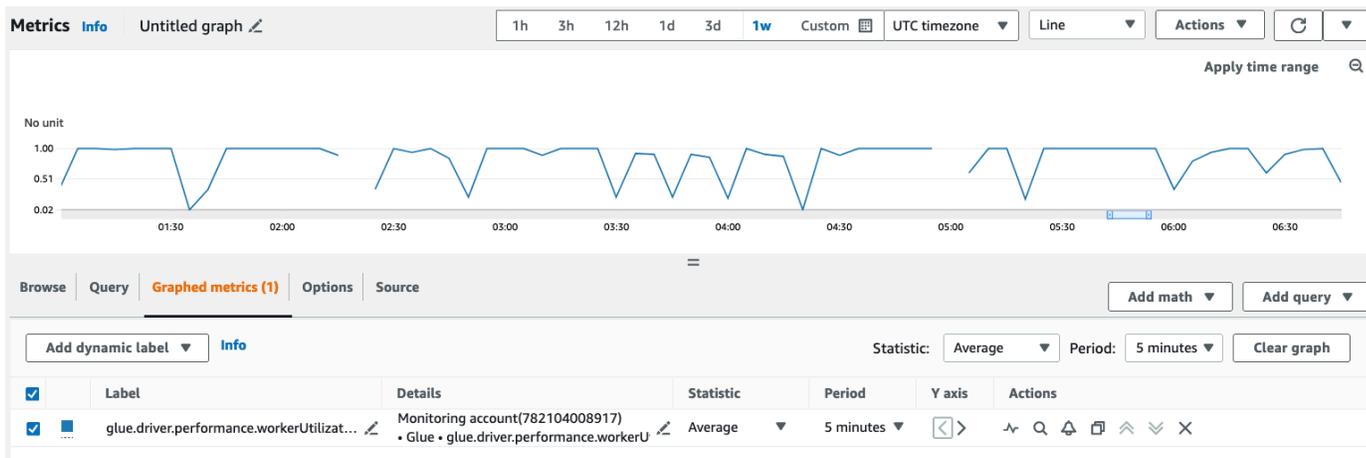
AWS Glue オブザーバビリティメトリクスは Amazon CloudWatch を通じて提供されるため、Amazon CloudWatch コンソール、AWS CLI、SDK、または API を使用してオブザーバビリティ

メトリクスデータポイントをクエリできます。AWS Glue オブザーバビリティメトリクスを使用する場合のユースケースの例については、「[Using Glue Observability for monitoring resource utilization to reduce cost](#)」を参照してください。

Amazon CloudWatch コンソールでの AWS Glue オブザーバビリティの使用

Amazon CloudWatch コンソールでメトリクスをクエリおよび視覚化するには:

1. Amazon CloudWatch コンソールを開き、[すべてのメトリクス] を選択します。
2. [カスタム名前空間] の下で、AWS Glue を選択します。
3. [ジョブのオブザーバビリティメトリクス]、[ソースごとのオブザーバビリティメトリクス]、または [シンクごとのオブザーバビリティメトリクス] を選択します。
4. 特定のメトリクス名、ジョブ名、ジョブ実行 ID を検索し、選択します。
5. [グラフ化されたメトリクス] タブで、任意の統計、期間、その他のオプションを設定します。



AWS CLI を使用してオブザーバビリティメトリクスをクエリするには:

1. メトリクス定義 JSON ファイルを作成し、`your-Glue-job-name` と `your-Glue-job-run-id` を実際のものに置き換えます。

```
$ cat multiplequeries.json
[
  {
    "Id": "avgWorkerUtil_0",
    "MetricStat": {
      "Metric": {
        "Namespace": "Glue",
        "MetricName": "glue.driver.workerUtilization",
```

```

        "Dimensions": [
            {
                "Name": "JobName",
                "Value": "<your-Glue-job-name-A>"
            },
            {
                "Name": "JobRunId",
                "Value": "<your-Glue-job-run-id-A>"
            },
            {
                "Name": "Type",
                "Value": "gauge"
            },
            {
                "Name": "ObservabilityGroup",
                "Value": "resource_utilization"
            }
        ]
    },
    "Period": 1800,
    "Stat": "Minimum",
    "Unit": "None"
}
},
{
    "Id": "avgWorkerUtil_1",
    "MetricStat": {
        "Metric": {
            "Namespace": "Glue",
            "MetricName": "glue.driver.workerUtilization",
            "Dimensions": [
                {
                    "Name": "JobName",
                    "Value": "<your-Glue-job-name-B>"
                },
                {
                    "Name": "JobRunId",
                    "Value": "<your-Glue-job-run-id-B>"
                },
                {
                    "Name": "Type",
                    "Value": "gauge"
                }
            ]
        }
    }
}

```

```

        "Name": "ObservabilityGroup",
        "Value": "resource_utilization"
      }
    ]
  },
  "Period": 1800,
  "Stat": "Minimum",
  "Unit": "None"
}
]

```

2. get-metric-data コマンドを実行します。

```

$ aws cloudwatch get-metric-data --metric-data-queries file: //multiplequeries.json \
\
  --start-time '2023-10-28T18: 20' \
  --end-time '2023-10-28T19: 10' \
  --region us-east-1
{
  "MetricDataResults": [
    {
      "Id": "avgWorkerUtil_0",
      "Label": "<your-label-for-A>",
      "Timestamps": [
        "2023-10-28T18:20:00+00:00"
      ],
      "Values": [
        0.06718750000000001
      ],
      "StatusCode": "Complete"
    },
    {
      "Id": "avgWorkerUtil_1",
      "Label": "<your-label-for-B>",
      "Timestamps": [
        "2023-10-28T18:50:00+00:00"
      ],
      "Values": [
        0.5959183673469387
      ],
      "StatusCode": "Complete"
    }
  ]
}

```

```

    }
  ],
  "Messages": []
}

```

オブザーバビリティメトリクス

AWS Glue オブザーバビリティは以下のメトリクスをプロファイリングし、30 秒ごとに Amazon CloudWatch に送信します。これらのメトリクスの一部は、AWS Glue Studio の「ジョブ実行モニタリング」ページで確認できます。

メトリクス	説明	カテゴリ
glue.driver.skewness.stage	<p>メトリクスカテゴリ: job_performance</p> <p>Spark ステージの実行歪度: このメトリクスは、入力データの歪みや変換 (スキュー結合など) によって生じる可能性のある実行の歪みをキャプチャします。このメトリクスの値は [0, infinity] の範囲に入ります。ここで、0 は、ステージ内のすべてのタスクのうち、タスクの実行時間の最大値と中央値の比率が特定のステージ歪度係数未満であることを意味します。デフォルトのステージ歪度係数は `5` で、spark conf: spark.metrics.conf.driver.source.glue.jobPerformance.skewnessFactor で上書きされます</p>	job_performance

メトリクス	説明	カテゴリ
	<p>ステージ歪度の値が 1 の場合、比率はステージ歪度係数の 2 倍になります。</p> <p>ステージ歪度の値は、現在の歪度を反映して 30 秒ごとに更新されます。ステージ終了時の値は、最終ステージの歪度を反映しています。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (job_performance)</p> <p>有効な統計: Average、Maximum、Minimum、Percentile</p> <p>単位: 個</p>	

メトリクス	説明	カテゴリ
glue.driver.skewness.job	<p>メトリクスカテゴリ: job_performance</p> <p>ジョブの歪度は、ジョブステージの歪度の加重平均です。加重平均では、実行に時間がかかるステージの重みが高くなります。これは、非常に歪んだステージが、実際には他のステージに比べて実行時間が非常に短い(したがって、その歪度はジョブパフォーマンス全体にとって重要ではなく、歪みの解消に努力する価値がない)というコーナーケースを回避するためです。</p> <p>このメトリクスは各ステージの完了時に更新されるため、最後の値には実際の全体的なジョブの歪みが反映されません。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (job_performance)</p> <p>有効な統計: Average、Maximum、Minimum、Percentile</p> <p>単位: 個</p>	job_performance

メトリクス	説明	カテゴリ
glue.succeed.ALL	<p>メトリクスカテゴリ: error</p> <p>障害カテゴリの全体像を把握するための、正常なジョブ実行の総数</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)、ObservabilityGroup (エラー)</p> <p>有効な統計: SUM</p> <p>単位: 個</p>	エラー
glue.error.ALL	<p>メトリクスカテゴリ: error</p> <p>障害カテゴリの全体像を把握するための、ジョブ実行エラーの総数</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)、ObservabilityGroup (エラー)</p> <p>有効な統計: SUM</p> <p>単位: 個</p>	エラー

メトリクス	説明	カテゴリ
glue.error.[error category]	<p>メトリクスカテゴリ: error</p> <p>これは実際には一連のメトリクスで、ジョブの実行が失敗した場合にのみ更新されません。エラー分類はトリアージとデバッグに役立ちます。ジョブ実行に失敗すると、失敗の原因となったエラーが分類され、対応するエラーカテゴリメトリックが 1 に設定されます。これにより、時間の経過に伴う障害分析だけでなく、すべてのジョブエラー分析を実行して、最も一般的な障害カテゴリを特定して対処を開始できます。AWS Glue には、OUT_OF_MEMORY (ドライバとエグゼキューター)、PERMISSION、SYNTAX、THROTTLING など、28 のエラーカテゴリがあります。また、COMPILATION、LAUNCH、TIMEOUT といったエラーカテゴリもあります。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (カウント)、ObservabilityGroup (エラー)</p> <p>有効な統計: SUM</p>	エラー

メトリクス	説明	カテゴリ
	単位: 個	
glue.driver.workerUtilization	<p>メトリクスカテゴリ: resource_utilization</p> <p>割り当てられたワーカーのうち、実際に使用されているワーカーの割合。うまく機能しない場合は、自動スケールリングが役に立ちます。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average、Maximum、Minimum、Percentile</p> <p>単位: パーセント</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.driver.memory.heap.[available used]	<p>メトリクスカテゴリ: resource_utilization</p> <p>ジョブ実行中のドライバの使用可能/使用済みヒープメモリ。これは、特に時間の経過に伴うメモリ使用量の傾向を把握するのに役立ち、メモリ関連の障害をデバッグするだけでなく、潜在的な障害を回避するのにも役立ちます。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: バイト</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.driver.memory.heap.used.percentage	<p>メトリクスカテゴリ: resource_utilization</p> <p>ジョブ実行中にドライバが (%) 使用したヒープメモリ。これは、特に時間の経過に伴うメモリ使用量の傾向を把握するのに役立ち、メモリ関連の障害をデバッグするだけでなく、潜在的な障害を回避するのに役立ちます。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: パーセント</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.driver.memory.non-heap. [available used]	<p>メトリクスカテゴリ: resource_utilization</p> <p>ジョブ実行中のドライバの使用可能/使用済み非ヒープメモリ。これは、特に時間の経過に伴うメモリ使用量の傾向を把握するのに役立ち、メモリ関連の障害をデバッグするだけでなく、潜在的な障害を回避するのにも役立ちます。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: バイト</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.driver.memory.non-heap.used.percentage	<p>メトリクスカテゴリ: resource_utilization</p> <p>ジョブ実行中にドライバが (%) 使用した非ヒープメモリ。これは、特に時間の経過に伴うメモリ使用量の傾向を把握するのに役立ち、メモリ関連の障害をデバッグするだけでなく、潜在的な障害を回避するのに役立ちます。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: パーセント</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.driver.memory.total.[available used]	<p>メトリクスカテゴリ: resource_utilization</p> <p>ジョブ実行中のドライバの使用可能/使用済み合計メモリ。これは、特に時間の経過に伴うメモリ使用量の傾向を把握するのに役立ち、メモリ関連の障害をデバッグするだけでなく、潜在的な障害を回避するのに役立ちます。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: バイト</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.driver.memory.total.used.percentage	<p>メトリクスカテゴリ: resource_utilization</p> <p>ジョブ実行中にドライバが (%) 使用した合計メモリ。これは、特に時間の経過に伴うメモリ使用量の傾向を把握するのに役立ち、メモリ関連の障害をデバッグするだけでなく、潜在的な障害を回避するのにも役立ちます。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: パーセント</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.ALL.memory.heap.[available used]	<p>メトリクスカテゴリ: resource_utilization</p> <p>エグゼキュターの使用可能/使用済みヒープメモリ。ALL はすべてのエグゼキュターを意味します。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: バイト</p>	resource_utilization
glue.ALL.memory.heap.used.percentage	<p>メトリクスカテゴリ: resource_utilization</p> <p>エグゼキュターの (%) 使用済みヒープメモリ。ALL はすべてのエグゼキュターを意味します。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: パーセント</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.ALL.memory.non-heap.[available used]	<p>メトリクスカテゴリ: resource_utilization</p> <p>エグゼキュターの使用可能/使用済み非ヒープメモリ。ALL はすべてのエグゼキュターを意味します。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: バイト</p>	resource_utilization
glue.ALL.memory.non-heap.used.percentage	<p>メトリクスカテゴリ: resource_utilization</p> <p>エグゼキュターの (%) 使用済み非ヒープメモリ。ALL はすべてのエグゼキュターを意味します。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: パーセント</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.ALL.memory.total.[available used]	<p>メトリクスカテゴリ: resource_utilization</p> <p>エグゼキュターの使用可能/使用済み合計メモリ。ALL はすべてのエグゼキュターを意味します。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: バイト</p>	resource_utilization
glue.ALL.memory.total.used.percentage	<p>メトリクスカテゴリ: resource_utilization</p> <p>エグゼキュターの (%) 使用済み合計メモリ。ALL はすべてのエグゼキュターを意味します。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: パーセント</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.driver.disk.[available_GB used_GB]	<p>メトリクスカテゴリ: resource_utilization</p> <p>ジョブ実行中のドライバの使用可能/使用済みディスク容量。これは、特に時間の経過に伴うディスク使用量の傾向を把握するのに役立ち、潜在的な障害を回避できるだけでなく、ディスク容量不足に関連する障害をデバッグするのにも役立ちます。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: ギガバイト</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.driver.disk.used.percentage]	<p>メトリクスカテゴリ: resource_utilization</p> <p>ジョブ実行中のドライバの使用可能/使用済みディスク容量。これは、特に時間の経過に伴うディスク使用量の傾向を把握するのに役立ち、潜在的な障害を回避できるだけでなく、ディスク容量不足に関連する障害をデバッグするのにも役立ちます。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: パーセント</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.ALL.disk.[available_GB used_GB]	<p>メトリクスカテゴリ: resource_utilization</p> <p>エグゼキュターの使用可能/使用済みディスク容量。ALL はすべてのエグゼキュターを意味します。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: ギガバイト</p>	resource_utilization
glue.ALL.disk.used.percentage	<p>メトリクスカテゴリ: resource_utilization</p> <p>エグゼキュターの使用可能/使用済み/(%) 使用済みディスク容量。ALL はすべてのエグゼキュターを意味します。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)</p> <p>有効な統計: Average</p> <p>単位: パーセント</p>	resource_utilization

メトリクス	説明	カテゴリ
glue.driver.bytesRead	<p data-bbox="591 226 889 310">メトリクスカテゴリ: throughput</p> <p data-bbox="591 352 1008 720">このジョブ実行で入力ソースごとに読み取られたバイト数、およびすべてのソースで読み取られたバイト数。これにより、データ量とその経時的な変化を把握でき、データの歪みなどの問題に対処しやすくなります。</p> <p data-bbox="591 762 1019 1087">有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)、Source (ソースデータの場所)</p> <p data-bbox="591 1129 889 1171">有効な統計: Average</p> <p data-bbox="591 1213 766 1255">単位: バイト</p>	スループット

メトリクス	説明	カテゴリ
glue.driver.[recordsRead filesRead]	<p>メトリクスカテゴリ: throughput</p> <p>このジョブ実行で入力ソースごとに読み取られたレコード/ファイルの数、およびすべてのソースで読み取られたレコード/ファイルの数。これにより、データ量とその経時的な変化を把握でき、データの歪みなどの問題に対処しやすくなります。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)、Source (ソースデータの場所)</p> <p>有効な統計: Average</p> <p>単位: 個</p>	スループット

メトリクス	説明	カテゴリ
glue.driver.partitionsRead	<p>メトリクスカテゴリ: throughput</p> <p>このジョブ実行で Amazon S3 入力ソースごとに読み取られたパーティションの数、およびすべてのソースで読み取られたのパーティション数。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)、Source (ソースデータの場所)</p> <p>有効な統計: Average</p> <p>単位: 個</p>	スループット

メトリクス	説明	カテゴリ
glue.driver.bytesWritten	<p>メトリクスカテゴリ: throughput</p> <p>このジョブ実行で出力シンク1つあたりに書き込まれたバイト数、およびすべてのシンクで書き込まれたバイト数。これにより、データ量とその経時的変化を把握でき、処理の歪みなどの問題に対処しやすくなります。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)、Sink (シンクデータの場所)</p> <p>有効な統計: Average</p> <p>単位: バイト</p>	スループット

メトリクス	説明	カテゴリ
glue.driver.[recordsWritten filesWritten]	<p>メトリクスカテゴリ: throughput</p> <p>このジョブ実行で出力シンクごとに書き込まれるレコード/ファイルの数、およびすべてのシンクで書き込まれるレコード/ファイルの数。これにより、データ量とその経時的変化を把握でき、処理の歪みなどの問題に対処しやすくなります。</p> <p>有効なディメンション: JobName (AWS Glue ジョブの名前)、JobRunId (JobRun ID. または ALL)、Type (ゲージ)、ObservabilityGroup (resource_utilization)、Sink (シンクデータの場所)</p> <p>有効な統計: Average</p> <p>単位: 個</p>	スループット

エラーカテゴリ

エラーカテゴリ	説明
COMPILATION_ERROR	Scala コードのコンパイル中に発生するエラー。
CONNECTION_ERROR	サービス/リモート、ホスト/データベースのサービスなどへの接続中に発生するエラー。

エラーカテゴリ	説明
DISK_NO_SPACE_ERROR	ドライバ/エグゼキュターのディスクに空きがないと発生するエラー。
OUT_OF_MEMORY_ERROR	ドライバ/エグゼキュターのメモリに空きがないと発生するエラー。
IMPORT_ERROR	依存関係をインポートすると発生するエラー。
INVALID_ARGUMENT_ERROR	入力引数が無効または不正なときにエラーが発生します。
PERMISSION_ERROR	サービス、データなどに対する権限がないと発生するエラー。
RESOURCE_NOT_FOUND_ERROR	データ、位置情報などが存在しないと発生するエラー。
QUERY_ERROR	Spark SQL クエリの実行によって発生するエラー。
SYNTAX_ERROR	スクリプトに構文エラーがあると発生するエラー。
THROTTLING_ERROR	サービスの同時実行数の制限に達したり、サービスクォータの制限を超えたりすると発生するエラー。
DATA_LAKE_FRAMEWORK_ERROR	AWS Glue がネイティブでサポートするデータレイクフレームワーク (Hudi、Iceberg など) から発生するエラー。
UNSUPPORTED_OPERATION_ERROR	サポートされていない操作を行うと発生するエラー。
RESOURCES_ALREADY_EXISTS_ERROR	作成または追加するリソースがすでに存在していると発生するエラー。

エラーカテゴリ	説明
GLUE_INTERNAL_SERVICE_ERROR	AWS Glue の内部サービスに問題があると発生するエラー。
GLUE_OPERATION_TIMEOUT_ERROR	AWS Glue 操作がタイムアウトになると発生するエラー。
GLUE_VALIDATION_ERROR	AWS Glue ジョブに必要な値を検証できなかった場合に発生するエラー。
GLUE_JOB_BOOKMARK_VERSION_MISMATCH_ERROR	同じソースバケットで同じジョブを実行し、同じ/異なる送信先に同時に書き込んだ場合 (同時実行数 >1) に発生するエラー。
LAUNCH_ERROR	AWS Glue ジョブの起動段階で発生するエラー。
DYNAMODB_ERROR	Amazon DynamoDB サービスから発生する一般的なエラー。
GLUE_ERROR	AWS Glue サービスから発生する一般的なエラー。
LAKEFORMATION_ERROR	AWS Lake Formation サービスから発生する一般的なエラー。
REDSHIFT_ERROR	Amazon Redshift サービスから発生する一般的なエラー。
S3_ERROR	Amazon S3 サービスから発生する一般的なエラー。
SYSTEM_EXIT_ERROR	一般的なシステム終了エラー。
TIMEOUT_ERROR	操作がタイムアウトしてジョブが失敗すると発生する一般的なエラー。
UNCLASSIFIED_SPARK_ERROR	Spark から発生する一般的なエラー。

エラーカテゴリ	説明
UNCLASSIFIED_ERROR	デフォルトのエラーカテゴリ。

制限事項

Note

メトリクスを公開するには、`glueContext` を初期化する必要があります。

ソースディメンションの値は、ソースタイプに応じて Amazon S3 パスまたはテーブル名のいずれかになります。さらに、ソースが JDBC でクエリオプションが使用されている場合、クエリ文字列はソースディメンションに設定されます。値が 500 文字を超える場合は、500 文字以内に切り捨てられます。値には次の制限があります。

- ASCII 以外の文字は削除されます。
- ソース名に ASCII 文字が含まれていない場合は、`<non-ASCII input>` に変換されます。

スループットメトリクスの制約事項と考慮事項

- `DataFrame` と `DataFrame` ベースの `DynamicFrame` (例:JDBC、Amazon S3 の `parquet` からの読み取り) はサポートされていますが、`RDD` ベースの `DynamicFrame` (Amazon S3 での `csv`、`json` の読み取りなど) はサポートされていません。技術的には、Spark UI に表示されるすべての読み取りと書き込みがサポートされています。
- データソースがカタログテーブルで、形式が JSON、CSV、テキスト、または Iceberg の場合、`recordsRead` メトリクスが出力されます。
- `glue.driver.throughput.recordsWritten`、`glue.driver.throughput.bytesWritten`、`glue.driver.throughput.recordsRead` メトリクスは JDBC テーブルと Iceberg テーブルでは使用できません。
- メトリクスは遅延する可能性があります。ジョブが約 1 分後に終了する場合、Amazon CloudWatch メトリクスにはスループットメトリクスがない可能性があります。

ジョブのモニタリングとデバッグ

AWS Glue ジョブに関するメトリクスを収集して AWS Glue および Amazon CloudWatch コンソールで可視化し、問題を特定して修正できます。AWS Glue ジョブのプロファイリングには、以下のステップが必要です。

1. メトリクスを有効にする:
 - a. ジョブ定義で [Job metrics (ジョブメトリクス)] オプションを有効にします。AWS Glue コンソールでプロファイリングを有効にするか、ジョブに対するパラメータとして有効にできます。詳細については、「[Spark ジョブのジョブプロパティの定義](#)」または「[AWS Glue ジョブのパラメータ](#)」を参照してください。
 - b. ジョブ定義で [AWS Glue オブザーバビリティメトリクス] オプションを有効にします。AWS Glue コンソールでオブザーバビリティを有効にするか、ジョブに対するパラメータとして有効にできます。詳細については、「[AWS Glue オブザーバビリティメトリクスを使用したモニタリング](#)」を参照してください。
2. ジョブスクリプトが GlueContext を初期化することを確認します。たとえば、次のスクリプトスニペットは GlueContext を初期化し、プロファイリングされたコードが配置されるスクリプト内の場所を示しています。この一般的な形式は、以下のデバッグシナリオで使用されません。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
import time

## @params: [JOB_NAME]
args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'], args)

...
...
```

```
code-to-profile
```

```
...  
...
```

```
job.commit()
```

3. ジョブを実行します。
4. メトリクスを可視化する:
 - a. AWS Glue コンソールでジョブメトリクスを可視化し、ドライバまたはエグゼキュターの異常なメトリクスを識別します。
 - b. オブザーバビリティメトリクスは、「ジョブ実行モニタリング」ページ、「ジョブ実行詳細」ページ、または Amazon CloudWatch で確認できます。詳細については、「[AWS Glue オブザーバビリティメトリクスを使用したモニタリング](#)」を参照してください。
5. 特定されたメトリクスを使用して根本原因を絞り込みます。
6. 必要に応じて、識別されたドライバーまたはジョブエグゼキュターのログストリームを使用して根本原因を確認します。

AWS Glue オブザーバビリティメトリクスのユースケース

- [OOM 例外とジョブの異常のデバッグ](#)
- [要求の厳しいステージとストラグラータスクのデバッグ](#)
- [複数のジョブの進行状況のモニタリング](#)
- [DPU の容量計画のモニタリング](#)
- [AWS Glue オブザーバビリティを使用してリソースの使用状況を監視し、コストを削減します。](#)

OOM 例外とジョブの異常のデバッグ

AWS Glue でメモリ不足 (OOM) 例外とジョブの異常をデバッグできます。以下のセクションでは、Apache Spark ドライバーや Spark エグゼキュターのメモリ不足例外をデバッグするためのシナリオについて説明します。

- [ドライバー OOM 例外のデバッグ](#)
- [エグゼキュター OOM 例外のデバッグ](#)

ドライバー OOM 例外のデバッグ

このシナリオでは、Spark ジョブで多数の小さいファイルが Amazon Simple Storage Service (Amazon S3) から読み込まれます。ファイルが Apache Parquet 形式に変換された後、Amazon S3 に書き込まれます。Spark ドライバーのメモリが不足しています。入力 Amazon S3 データの複数の Amazon S3 パーティションに 100 万を超えるファイルがあります。

プロファイルされたコードは次のとおりです。

```
data = spark.read.format("json").option("inferSchema", False).load("s3://input_path")
data.write.format("parquet").save(output_path)
```

AWS Glue コンソールでプロファイルされたメトリクスを可視化する

以下のグラフでは、ドライバーとエグゼキューターのメモリ使用率がパーセントで示されています。この使用率は、直近の 1 分間に報告された値を平均した 1 つのデータポイントとしてプロットされます。ジョブのメモリプロファイルでは、[ドライバーメモリ](#)が安全しきい値である使用率の 50% をすぐに超えることがわかります。一方、すべてのエグゼキューターにおける[平均メモリ使用率](#)は、まだ 4% 未満です。これは、この Spark ジョブにおけるドライバー実行に異常があることを明確に示しています。



ジョブの実行はすぐに失敗し、AWS Glue コンソールの [History] (履歴) タブにエラー「Command Failed with Exit Code 1」が表示されます。このエラー文字列は、システムエラーのためにジョブが失敗したことを意味します。この場合はドライバーのメモリ不足です。

Run ID	Retry attempt	Run status	Error	Logs	Error logs	Execution time	Timeout	Delay	Triggered by	Start time
jr_651bfc34...	-	Failed	!	Logs	Error logs	2 mins	2880 mins			7 June 2018 7:37 PM UTC-7
jr_5731b225...	-	Failed	Command failed with exit code 1			ins	2880 mins			7 June 2018 7:37 PM UTC-7

コンソールの [History] (履歴) タブにある [Error logs] (エラーログ) リンクを選択し、CloudWatch Logs からのドライバー OOM に関する詳細情報を確認します。ジョブのエラーログで「**Error**」を検索し、それが本当にジョブの失敗の原因となった OOM 例外であることを確認します。

```
# java.lang.OutOfMemoryError: Java heap space
# -XX:OnOutOfMemoryError="kill -9 %p"
# Executing /bin/sh -c "kill -9 12039"...
```

ジョブの [履歴] タブで、[ログ] を選択します。ジョブの開始時における CloudWatch Logs 内のドライバー実行の以下のトレースが表示されます。Spark ドライバーは、すべてのディレクトリのすべてのファイルのリストの取得を試み、InMemoryFileIndex を構築してファイルごとに 1 つのタスクを起動します。この結果、Spark ドライバーは、すべてのタスクを追跡するため、大量の状態をメモリに保持する必要が生じます。また、メモリ内インデックスの多数のファイルの完全なリストを取得するため、ドライバー OOM となります。

グループ化を使用した複数のファイルの処理を修正する

でグループ化AWS Glue機能を使用して、複数のファイルの処理を修正できます。グループ化は、動的なフレームを使用しているときと、入力データセットに多数のファイル (50,000 超) があるときに自動的に有効になります。グループ化により、複数のファイルを 1 つのグループにまとめることができ、タスクは単一のファイルではなくグループ全体を処理できるようになります。その結果、Spark ドライバーがメモリに保存する状態がかなり少なくなり、追跡するタスクが減少します。

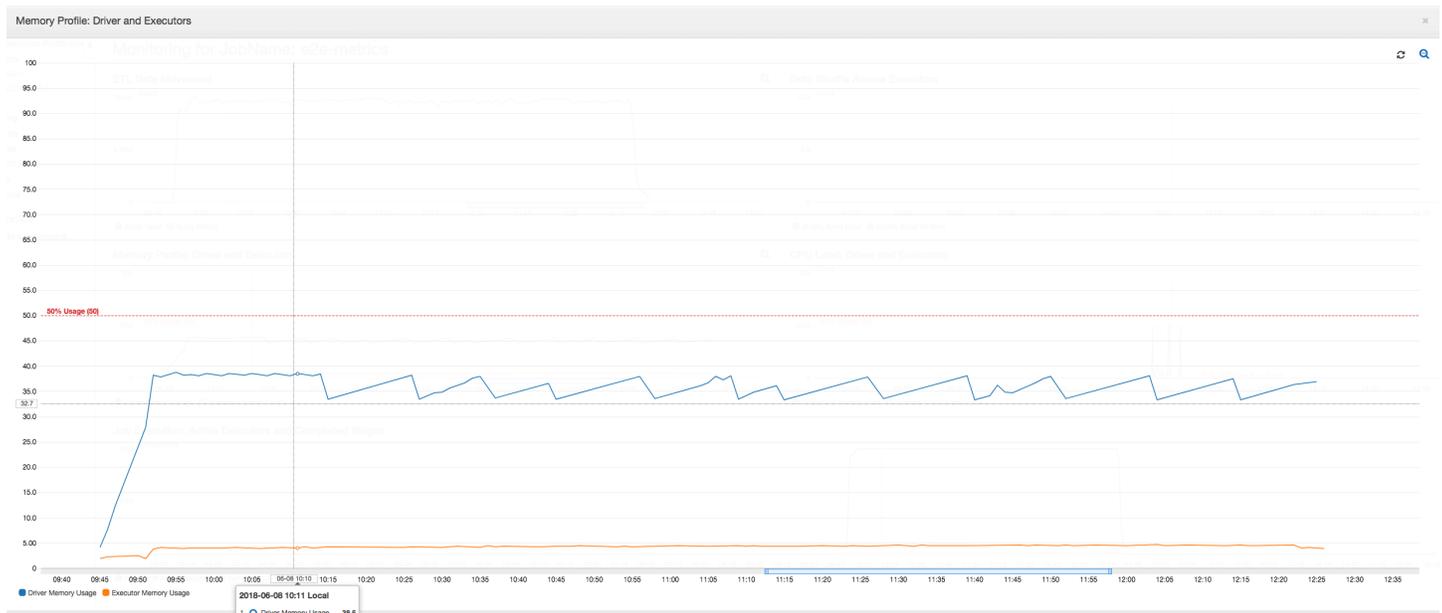
データセットのグループ化を手動で有効にする方法の詳細については、「[大きなグループの入カファイルの読み取り](#)」を参照してください。

AWS Glue ジョブのメモリプロファイルを確認するには、グループ化を有効にして次のコードをプロファイルします。

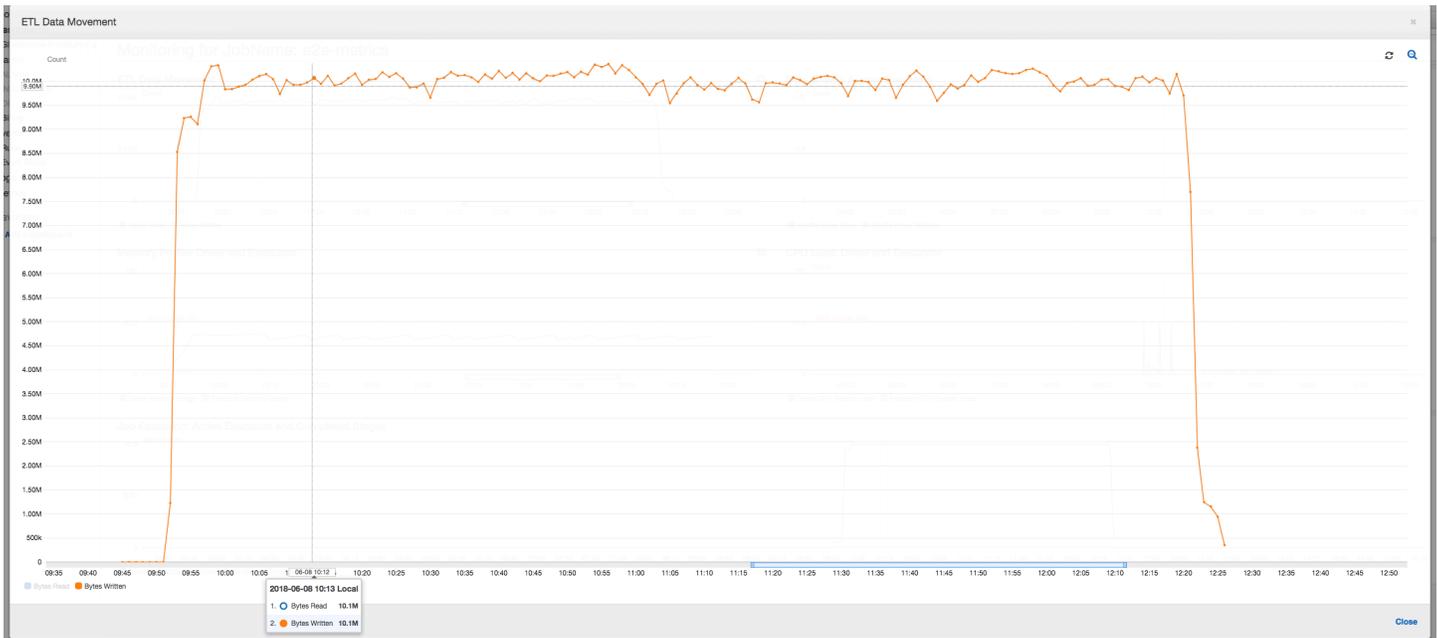
```
df = glueContext.create_dynamic_frame_from_options("s3", {'paths': ["s3://input_path"],
  "recurse":True, 'groupFiles': 'inPartition'}, format="json")
datasink = glueContext.write_dynamic_frame.from_options(frame = df, connection_type
  = "s3", connection_options = {"path": output_path}, format = "parquet",
  transformation_ctx = "datasink")
```

AWS Glue ジョブプロファイルでは、メモリプロファイルと ETL データ移動を監視できます。

ドライバーは、AWS Glue ジョブの持続期間全体にわたって、メモリ使用率のしきい値である 50% 未満で実行されます。エグゼキューターは、Amazon S3 からデータをストリーミングして処理し、Amazon S3 に書き出します。その結果、消費されるメモリは常に 5% 未満となります。



以下のデータ移動プロファイルは、ジョブが進行するにつれてすべてのエグゼキューターにより直近 1 分間に読み取りおよび書き込みされた Amazon S3 バイトの合計数を示しています。データはすべてのエグゼキューターでストリーミングされるため、どちらも同様のパターンに従っています。ジョブは、100 万ファイルすべての処理を 3 時間未満で完了します。



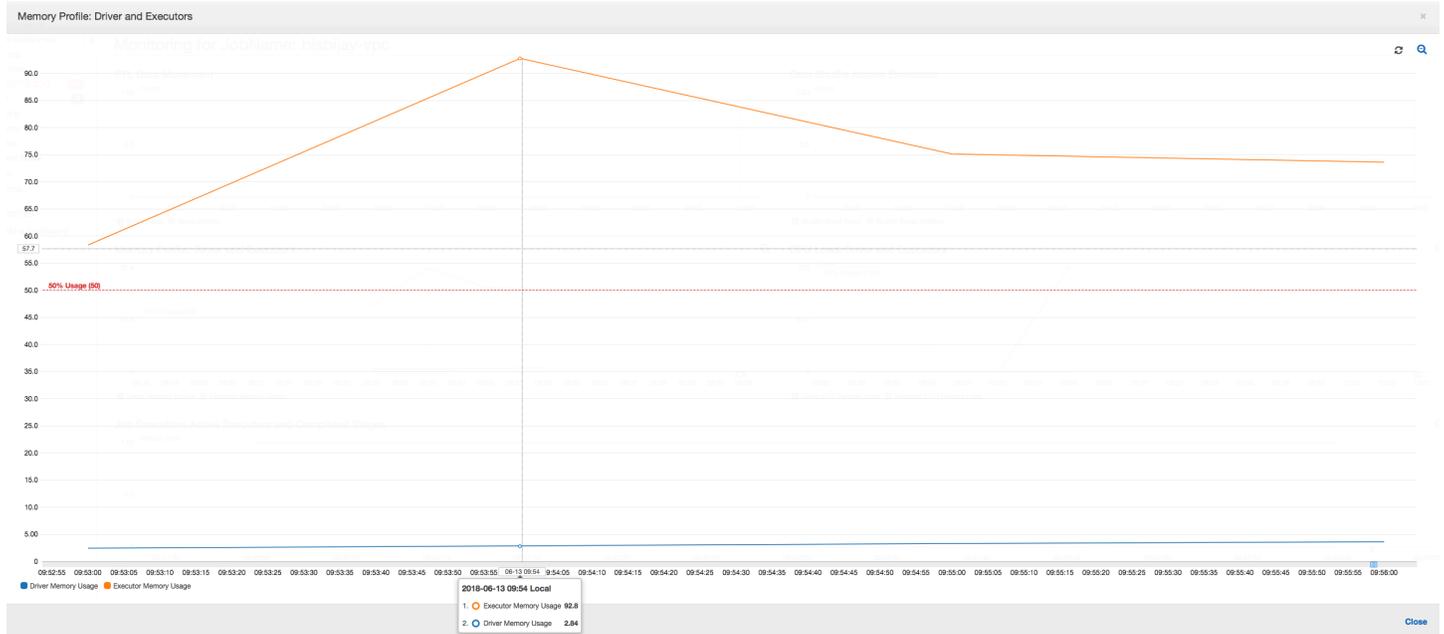
エグゼキュター OOM 例外のデバッグ

このシナリオでは、Apache Spark エグゼキュターで発生する可能性がある OOM 例外をデバッグする方法について説明します。次のコードでは、Spark MySQL リーダーを使用して、約 34 万行の大きなテーブルを Spark データフレームに取り込みます。次に、Parquet 形式で Amazon S3 に書き出します。接続プロパティを用意し、デフォルト Spark 設定を使用してテーブルを読み取ることができます。

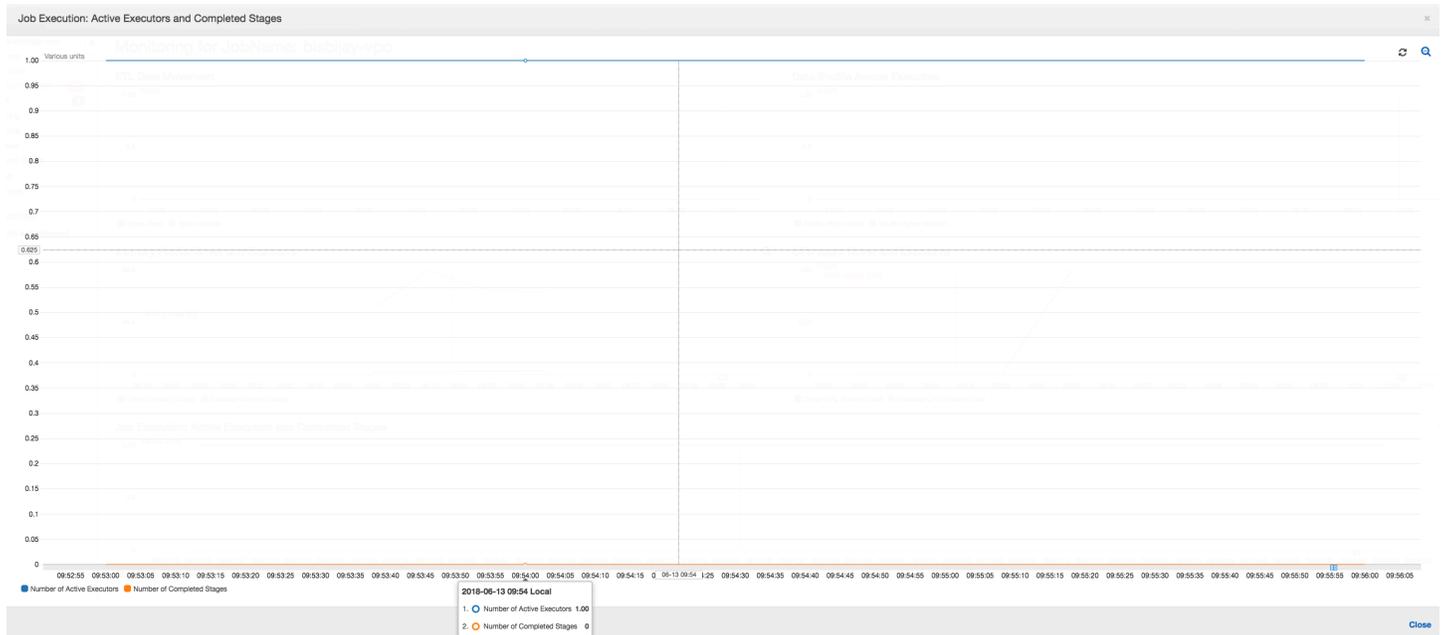
```
val connectionProperties = new Properties()
connectionProperties.put("user", user)
connectionProperties.put("password", password)
connectionProperties.put("Driver", "com.mysql.jdbc.Driver")
val sparkSession = glueContext.sparkSession
val dfSpark = sparkSession.read.jdbc(url, tableName, connectionProperties)
dfSpark.write.format("parquet").save(output_path)
```

AWS Glue コンソールでプロファイルされたメトリクスを可視化する

メモリ使用量グラフの傾きが正で 50% を超えた場合、および次のメトリクスが出力される前にジョブが失敗した場合は、メモリ枯渇が原因の候補となります。以下のグラフは、実行してから 1 分以内に、すべてのエグゼキュターにおける [平均メモリ使用率](#) がすぐに 50% を超えることを示しています。使用率は最大 92% に上昇し、エグゼキュターを実行するコンテナは Apache Hadoop YARN により停止されます。

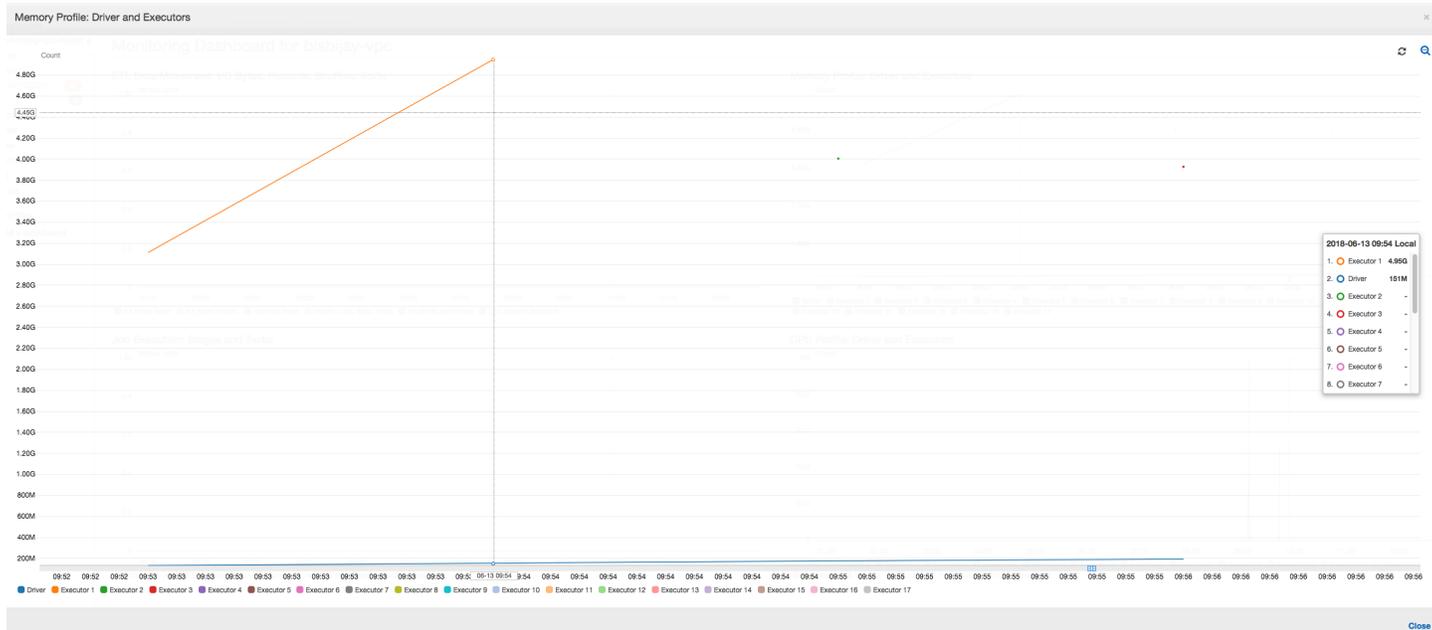


次のグラフが示すよう、ジョブが失敗するまで常に単一のエグゼキューターが実行されています。これは、新しいエグゼキューターが起動され、停止されたエグゼキューターが置き換えられるためです。JDBC データソースの読み取りは、列上のテーブルをパーティション化して複数の接続を開く必要があるため、デフォルトでは並列化されません。その結果、1つのエグゼキューターだけがテーブル全体を順番に読み込みます。



次のグラフが示すように、Spark はジョブが失敗するまでに 4 回新しいタスクを起動しようとして、3つのエグゼキューターのメモリプロファイルを確認できます。各エグゼキューターはすべてのメモ

りをすぐに使い尽くします。4 番目のエグゼキュターのメモリが不足し、ジョブは失敗します。その結果、そのメトリクスはすぐに報告されません。



次の図に示すように、AWS Glue コンソールのエラー文字列からは、ジョブが OOM 例外のために失敗したことを確認できます。

Run ID	Retry attempt	Run status	Error	Logs	Error logs	Execution time	Timeout	Delay	Triggered by	Start time	End time
j_f_4fc7d2723c5d834e90ac0a0215...	-	Failed	org.apache.spark.SparkException: Job aborted due to stage failure: Task 0 in stage 0.0 failed 4 times, most recent failure: Lost task 0.3 in stage 0.0 (TID 3, ip-10-1-2-96.ec2.internal, executor 4): ExecutorLostFailure (executor 4 exited caused by one of the running tasks) Reason: Container killed by YARN for exceeding memory limits. 5.5 GB of 5.5 GB physical memory used. Consider boosting spark.yarn.executor.memoryOverhead.			0 secs	2880 mins			13 June 2018 9:48 AM UT...	13 June 2018 9:50 AM UT...
j_f_d70a0e928d9e7399a8152db4...	-	Succeeded				2 mins	2880 mins			13 June 2018 9:32 AM UT...	13 June 2018 9:44 AM UT...
j_f_64c857823082befad919116a2...	-	Succeeded				2 mins	2880 mins			13 June 2018 8:57 AM UT...	13 June 2018 9:09 AM UT...
j_f_7a0d552d68b36bcd53bbe745...	-	Failed				1 hr, 8 mins	2880 mins			12 June 2018 5:15 PM UT...	12 June 2018 6:31 PM UT...

ジョブ出力ログ: エグゼキュター OOM 例外の結果をさらに確認するには、CloudWatch Logs を参照します。「**Error**」を検索すると、メトリクスダッシュボードに示されているように、4 つのエグゼキュターがほぼ同じ時間枠で停止されています。メモリ制限を超えると、すべて YARN により終了されます。

エグゼキュター 1

```
18/06/13 16:54:29 WARN YarnAllocator: Container killed by YARN for exceeding
memory limits. 5.5 GB of 5.5 GB physical memory used. Consider boosting
spark.yarn.executor.memoryOverhead.
```

```
18/06/13 16:54:29 WARN YarnSchedulerBackend$YarnSchedulerEndpoint: Container killed
by YARN for exceeding memory limits. 5.5 GB of 5.5 GB physical memory used. Consider
boosting spark.yarn.executor.memoryOverhead.
```

```
18/06/13 16:54:29 ERROR YarnClusterScheduler: Lost executor 1 on
ip-10-1-2-175.ec2.internal: Container killed by YARN for exceeding
```

```
memory limits. 5.5 GB of 5.5 GB physical memory used. Consider boosting
spark.yarn.executor.memoryOverhead.
18/06/13 16:54:29 WARN TaskSetManager: Lost task 0.0 in stage 0.0 (TID 0,
ip-10-1-2-175.ec2.internal, executor 1): ExecutorLostFailure (executor 1
exited caused by one of the running tasks) Reason: Container killed by YARN for
exceeding memory limits. 5.5 GB of 5.5 GB physical memory used. Consider boosting
spark.yarn.executor.memoryOverhead.
```

エグゼキュター 2

```
18/06/13 16:55:35 WARN YarnAllocator: Container killed by YARN for exceeding
memory limits. 5.8 GB of 5.5 GB physical memory used. Consider boosting
spark.yarn.executor.memoryOverhead.
18/06/13 16:55:35 WARN YarnSchedulerBackend$YarnSchedulerEndpoint: Container killed
by YARN for exceeding memory limits. 5.8 GB of 5.5 GB physical memory used. Consider
boosting spark.yarn.executor.memoryOverhead.
18/06/13 16:55:35 ERROR YarnClusterScheduler: Lost executor 2 on
ip-10-1-2-16.ec2.internal: Container killed by YARN for exceeding
memory limits. 5.8 GB of 5.5 GB physical memory used. Consider boosting
spark.yarn.executor.memoryOverhead.
18/06/13 16:55:35 WARN TaskSetManager: Lost task 0.1 in stage 0.0 (TID 1,
ip-10-1-2-16.ec2.internal, executor 2): ExecutorLostFailure (executor 2 exited
caused by one of the running tasks) Reason: Container killed by YARN for
exceeding memory limits. 5.8 GB of 5.5 GB physical memory used. Consider boosting
spark.yarn.executor.memoryOverhead.
```

エグゼキュター 3

```
18/06/13 16:56:37 WARN YarnAllocator: Container killed by YARN for exceeding
memory limits. 5.8 GB of 5.5 GB physical memory used. Consider boosting
spark.yarn.executor.memoryOverhead.
18/06/13 16:56:37 WARN YarnSchedulerBackend$YarnSchedulerEndpoint: Container killed
by YARN for exceeding memory limits. 5.8 GB of 5.5 GB physical memory used. Consider
boosting spark.yarn.executor.memoryOverhead.
18/06/13 16:56:37 ERROR YarnClusterScheduler: Lost executor 3 on
ip-10-1-2-189.ec2.internal: Container killed by YARN for exceeding
memory limits. 5.8 GB of 5.5 GB physical memory used. Consider boosting
spark.yarn.executor.memoryOverhead.
18/06/13 16:56:37 WARN TaskSetManager: Lost task 0.2 in stage 0.0 (TID 2,
ip-10-1-2-189.ec2.internal, executor 3): ExecutorLostFailure (executor 3
exited caused by one of the running tasks) Reason: Container killed by YARN for
```

```
exceeding memory limits. 5.8 GB of 5.5 GB physical memory used. Consider boosting spark.yarn.executor.memoryOverhead.
```

エグゼキュター 4

```
18/06/13 16:57:18 WARN YarnAllocator: Container killed by YARN for exceeding memory limits. 5.5 GB of 5.5 GB physical memory used. Consider boosting spark.yarn.executor.memoryOverhead.
18/06/13 16:57:18 WARN YarnSchedulerBackend$YarnSchedulerEndpoint: Container killed by YARN for exceeding memory limits. 5.5 GB of 5.5 GB physical memory used. Consider boosting spark.yarn.executor.memoryOverhead.
18/06/13 16:57:18 ERROR YarnClusterScheduler: Lost executor 4 on ip-10-1-2-96.ec2.internal: Container killed by YARN for exceeding memory limits. 5.5 GB of 5.5 GB physical memory used. Consider boosting spark.yarn.executor.memoryOverhead.
18/06/13 16:57:18 WARN TaskSetManager: Lost task 0.3 in stage 0.0 (TID 3, ip-10-1-2-96.ec2.internal, executor 4): ExecutorLostFailure (executor 4 exited caused by one of the running tasks) Reason: Container killed by YARN for exceeding memory limits. 5.5 GB of 5.5 GB physical memory used. Consider boosting spark.yarn.executor.memoryOverhead.
```

AWS Glue 動的フレームを使用してフェッチサイズ設定を修正する

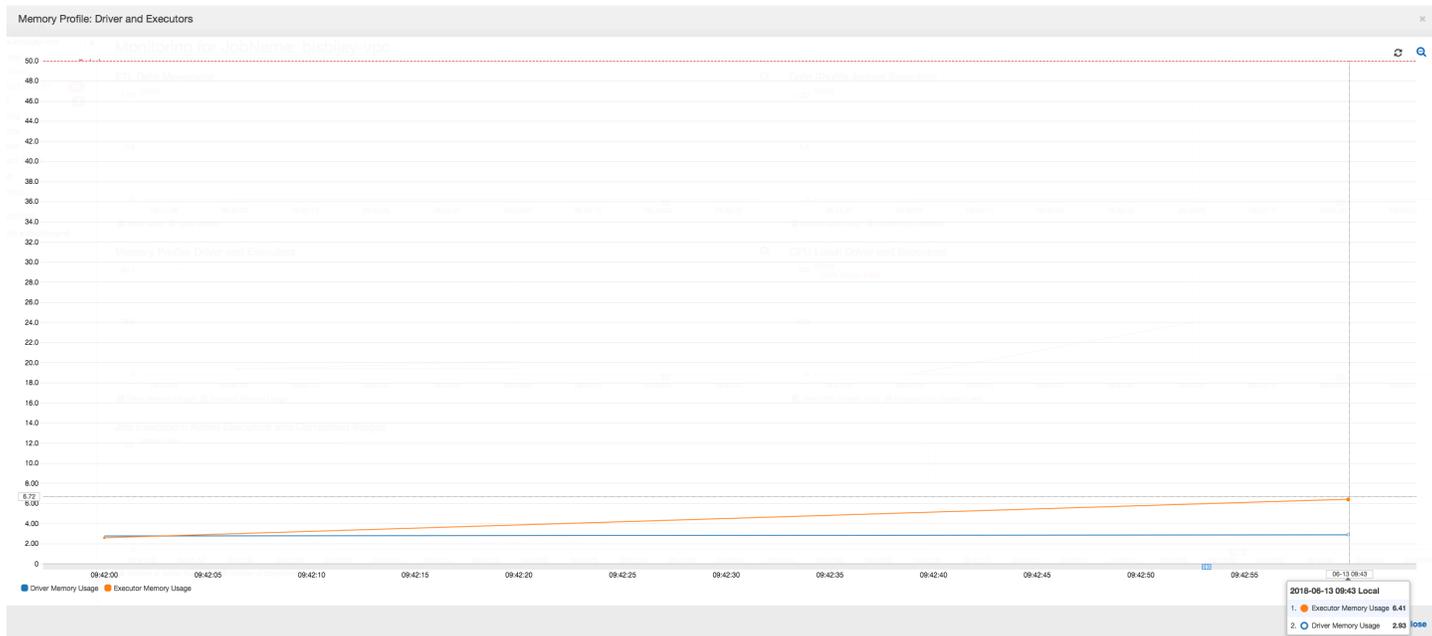
Spark JDBC フェッチサイズのデフォルトの設定が 0 であるため、エグゼキュターは JDBC テーブルの読み取り中にメモリ不足になります。つまり、Spark は行を一度に 1 つずつストリーミングしますが、Spark エグゼキュター上の JDBC ドライバーがデータベースから 3400 万行をまとめてフェッチして、それらをキャッシュしようとしています。Spark を使用すると、フェッチサイズパラメータを 0 以外のデフォルト値に設定することにより、このシナリオを回避できます。

代わりに AWS Glue 動的フレームを使用することで、この問題を解決することもできます。デフォルトでは、動的フレームは 1,000 行のフェッチサイズを使用します。これは通常十分な値です。このため、エグゼキュターが合計メモリの 7% 超を使用することはありません。AWS Glue ジョブは、エグゼキュターを 1 つだけ使用して 2 分未満で完了します。AWS Glue 動的フレームを使用することが推奨されるアプローチですが、Apache Spark の `fetchsize` プロパティを使用してフェッチサイズを設定することもできます。「[Spark SQL、DataFrames および Datasets ガイド](#)」を参照してください。

```
val (url, database, tableName) = {
  ("jdbc_url", "db_name", "table_name")
}
val source = glueContext.getSource(format, sourceJson)
```

```
val df = source.getDynamicFrame
glueContext.write_dynamic_frame.from_options(frame = df, connection_type = "s3",
connection_options = {"path": output_path}, format = "parquet", transformation_ctx =
"datasink")
```

プロファイルされた通常のメトリクス: 以下の図に示すように、AWS Glue 動的フレームを持つ [エグゼキューターメモリ](#) が安全しきい値を超えることはありません。データベースから行をストリーミングし、任意の時点で 1,000 行のみ JDBC ドライバーにキャッシュします。メモリ不足例外は発生しません。



要求の厳しいステージとストラグラータスクのデバッグ

AWS Glue ジョブプロファイリングを使用して、抽出、変換、ロード (ETL) ジョブで要求の厳しいステージとストラグラータスクを特定できます。ストラグラータスクは、AWS Glue ジョブのステージに含まれる他のタスクよりかなり時間がかかります。その結果、ステージの完了まで時間がかかり、ジョブの合計実行時間も遅延します。

小さい入力ファイルを大きい出力ファイルにまとめる

ストラグラータスクは、さまざまなタスク間の処理が均等に分散していないときや、データスキューによって特定のタスクが処理するデータが多くなった場合に発生する可能性があります。

Apache Spark の一般的なパターンである次のコードをプロファイルし、多数の小さいファイルをいくつかの大きい出力ファイルにまとめることができます。この例では、入力データセットは 32 GB の JSON Gzip 圧縮ファイルです。出力データセットには、約 190 GB の圧縮されていない JSON ファイルが含まれています。

プロファイルされたコードは次のとおりです。

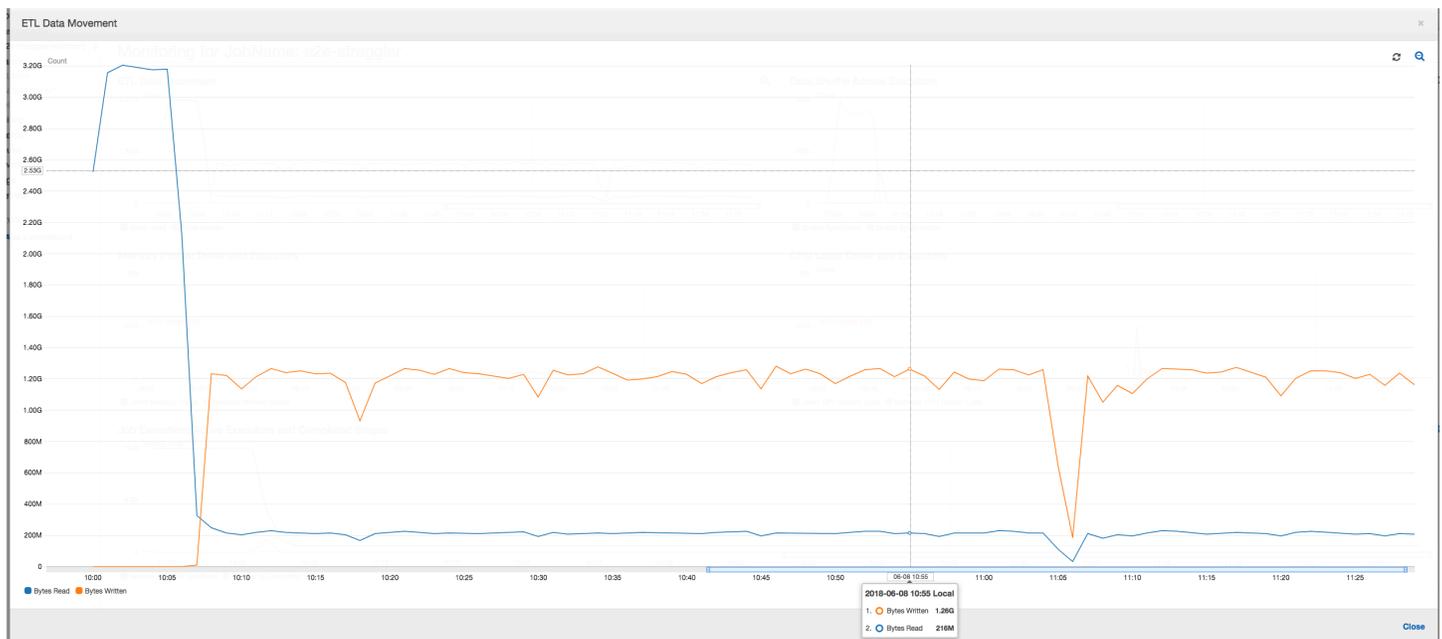
```
datasource0 = spark.read.format("json").load("s3://input_path")
df = datasource0.coalesce(1)
df.write.format("json").save(output_path)
```

AWS Glue コンソールでプロファイルされたメトリクスを可視化する

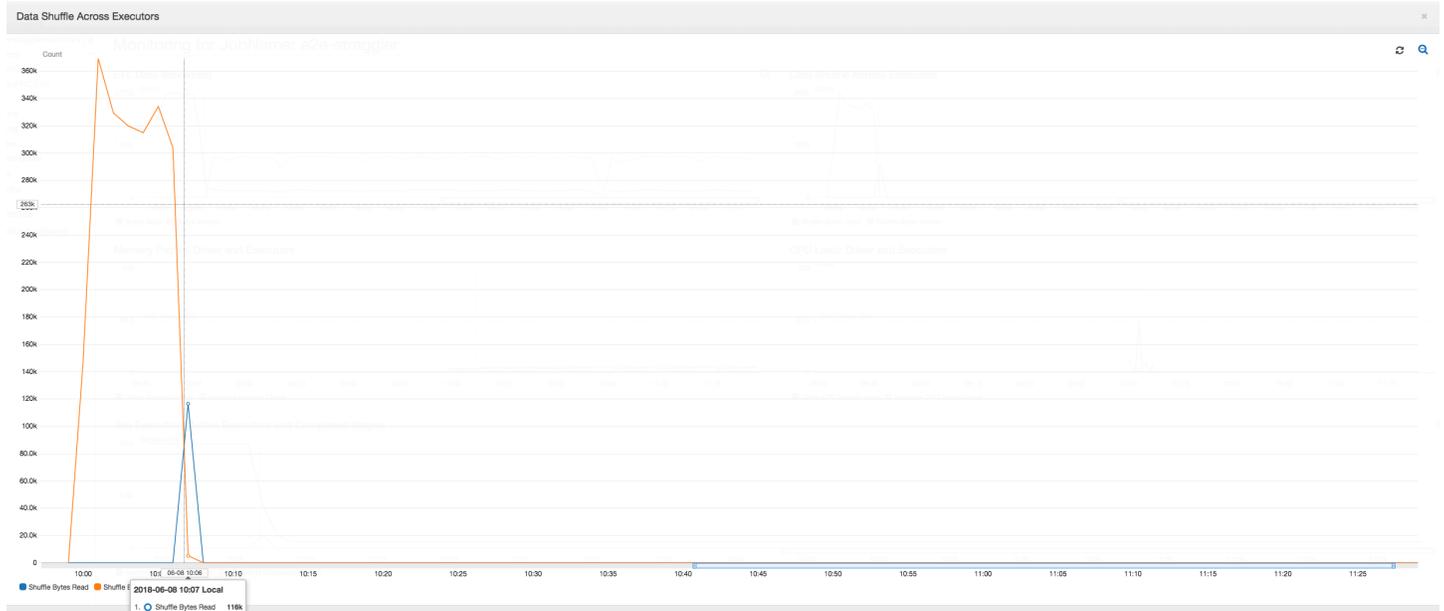
ジョブをプロファイルすると、4つのメトリクスセットを調べることができます。

- ETL データ移動
- エグゼキューター間のデータシャッフル
- ジョブの実行
- メモリプロファイル

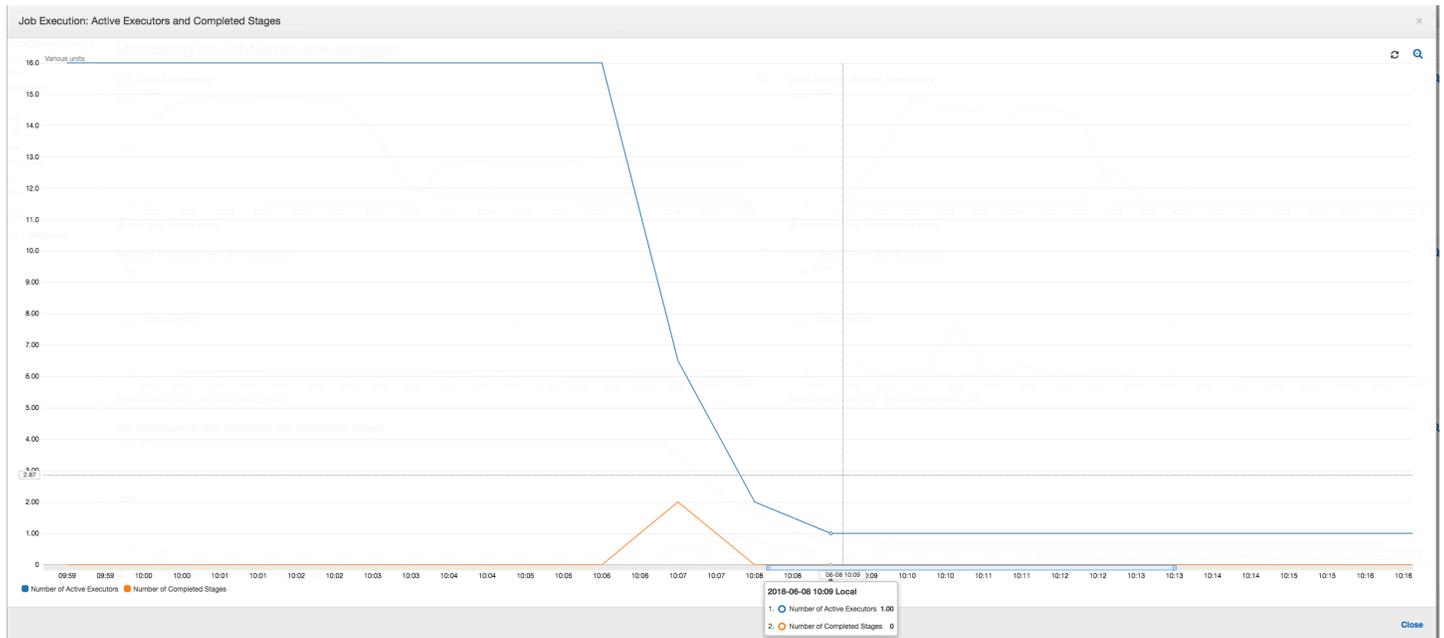
ETL データ移動: [ETL Data Movement (ETL データ移動)] プロファイルでは、バイトは、最初の 6 間に完了する最初のステージのすべてのエグゼキューターにより比較的早く 読み取られます。ただし、ジョブの合計実行時間が約 1 時間の場合、ほとんどがデータの 書き込み で構成されます。



エグゼキューター間のデータシャッフル: ジョブ実行メトリクスと データシャッフルメトリクスが示しているように、シャッフル中の読み取りおよび書き込みのバイト数も、ステージ 2 が終了する前に急増しています。すべてのエグゼキューターからのデータシャッフルの後、読み取りと書き込みはエグゼキューター 3 番のみから続行されます。



ジョブの実行: 以下のグラフに示すように、他のすべてのエグゼキュターはアイドル状態であり、最終的に時間 10:09 までに破棄されます。この時点で、エグゼキュターの合計数はわずか 1 に減ります。これは、エグゼキュター 3 番が、実行時間が最も長いストラグラータスクで構成されているため、ジョブの実行時間の大部分を占めていることを明確に示しています。



メモリプロファイル: 最初の 2 つのステージの後、[エグゼキュター 3 番](#)のみがメモリをアクティブに消費してデータを処理します。残りのエグゼキュターはアイドル状態のままであるか、最初の 2 つのステージの完了後すぐに破棄されます。



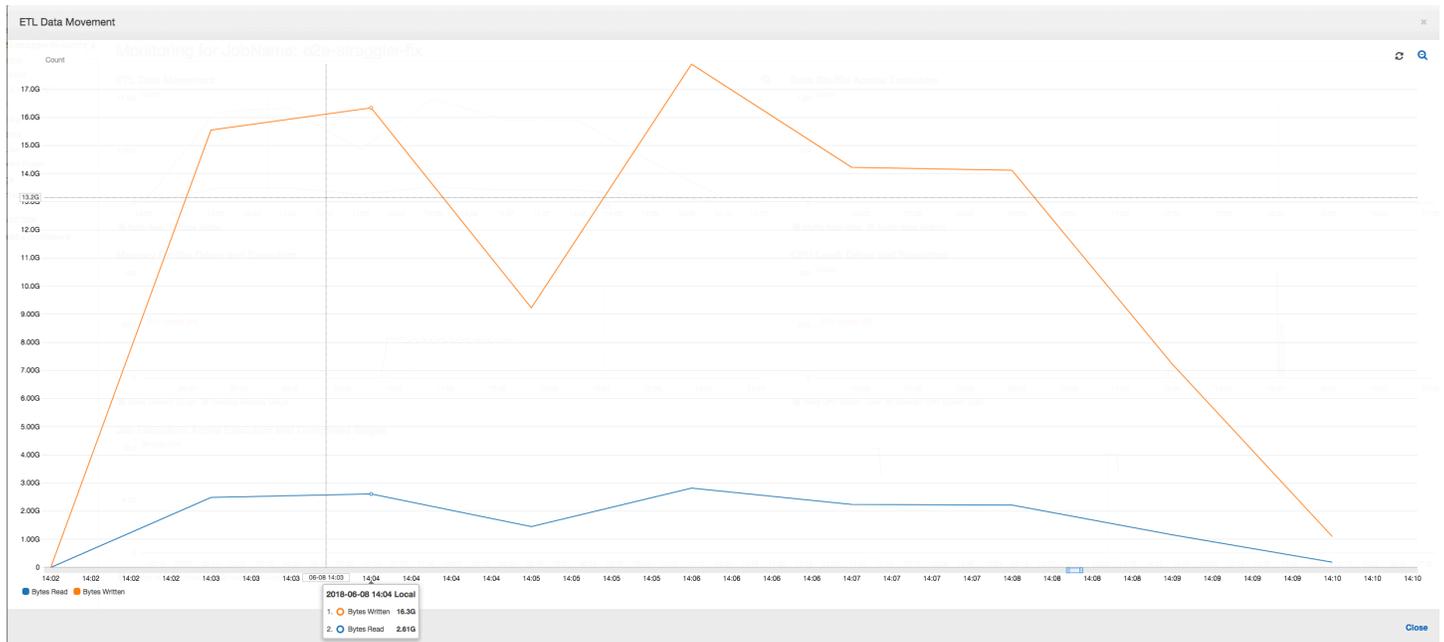
グループ化を使用してストラグラエグゼキュターを修正する

ストラグラエグゼキュターは、でグループ化AWS Glueを使用して回避できます。グループ化を使用してデータをすべてのエグゼキュター間に均等に分散し、クラスター内で利用可能なすべてのエグゼキュターを使用して大きいファイルにまとめます。詳細については、[大きなグループの入カファイルの読み取り](#)を参照してください。

AWS Glue ジョブで ETL データ移動を確認するには、グループ化を有効にして次のコードをプロファイルします。

```
df = glueContext.create_dynamic_frame_from_options("s3", {'paths': ["s3://input_path"],
"recurse":True, 'groupFiles': 'inPartition'}, format="json")
datasink = glueContext.write_dynamic_frame.from_options(frame = df, connection_type =
"s3", connection_options = {"path": output_path}, format = "json", transformation_ctx
= "datasink4")
```

ETL データ移動: データの書き込みは、ジョブの実行時間全体にわたって、データの読み取りと並行してストリーミングされるようになりました。その結果、ジョブは 8 分以内に終了します。以前よりもはるかに早いスピードです。



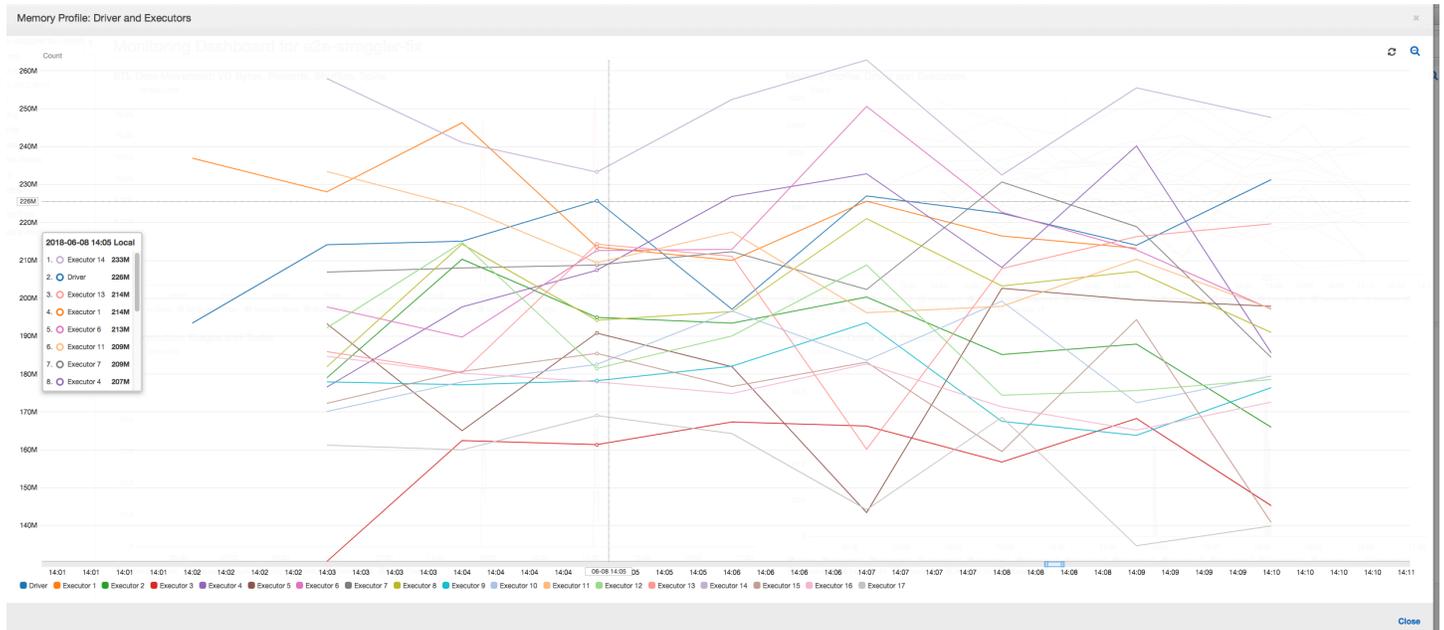
エグゼキューター間でのデータシャッフル: グループ化機能を使用して読み取り中に入力ファイルがまとめられると、データの読み取り後、コストのかかるデータシャッフルが行われなくなります。



ジョブの実行: ジョブ実行メトリクスは、実行されてデータを処理しているアクティブなエグゼキューターの合計数がかかり安定していることを示しています。ジョブに単一のストラグラーはありません。すべてのエグゼキューターがアクティブで、ジョブが完了するまでは破棄されません。エグゼキューター間でデータの間シャッフルは行われないため、ジョブにはステージが1つしかありません。



メモリプロファイル: このメトリクスは、すべてのエグゼキュター間の[アクティブなメモリ消費](#)を示しています。これは、すべてのエグゼキュター間にアクティビティがあることを再確認するものです。データは並行してストリーミングおよび書き出されるため、すべてのエグゼキュターの合計メモリ使用量はほぼ均等で、すべてのエグゼキュター安全しきい値をかなり下回っています。



複数のジョブの進行状況のモニタリング

複数の AWS Glue ジョブをまとめてプロファイルし、それらの間のデータフローを監視できます。これは、一般的なワークフローパターンであり、個々のジョブの進行状況、データ処理バックログ、データの再処理、ジョブのブックマークのモニタリングが必要です。

トピック

- [プロファイルされたコード](#)
- [AWS Glue コンソールでプロファイルされたメトリクスを可視化する](#)
- [ファイルの処理を修正する](#)

プロファイルされたコード

このワークフローには、入力ジョブと出力ジョブの 2 つがあります。入力ジョブは、定期的なトリガーを使用して 30 分ごとに実行するようにスケジュールされます。出力ジョブは、入力ジョブの実行が成功するたびに実行されるようにスケジュールされます。このようなスケジュールされたジョブは、ジョブトリガーを使用して制御されます。

Triggers A trigger starts a job when it fires.

Trigger name	Trigger type	Trigger status	Trigger parameters	Jobs to trigger
<input type="checkbox"/> e2e-bookmark-input	Schedule	ACTIVATED	Every 15 minutes	e2ebookmark-input
<input type="checkbox"/> e2e-bookmark-output	Job events	ACTIVATED	Job events: e2ebookmark-input	e2e-bookmark

入力ジョブ: このジョブは Amazon Simple Storage Service (Amazon S3) の場所からデータを読み取り、ApplyMapping を使用して変換し、Amazon S3 のステージング場所へ書き込みます。次のコードは、入力ジョブのプロファイルされたコードです。

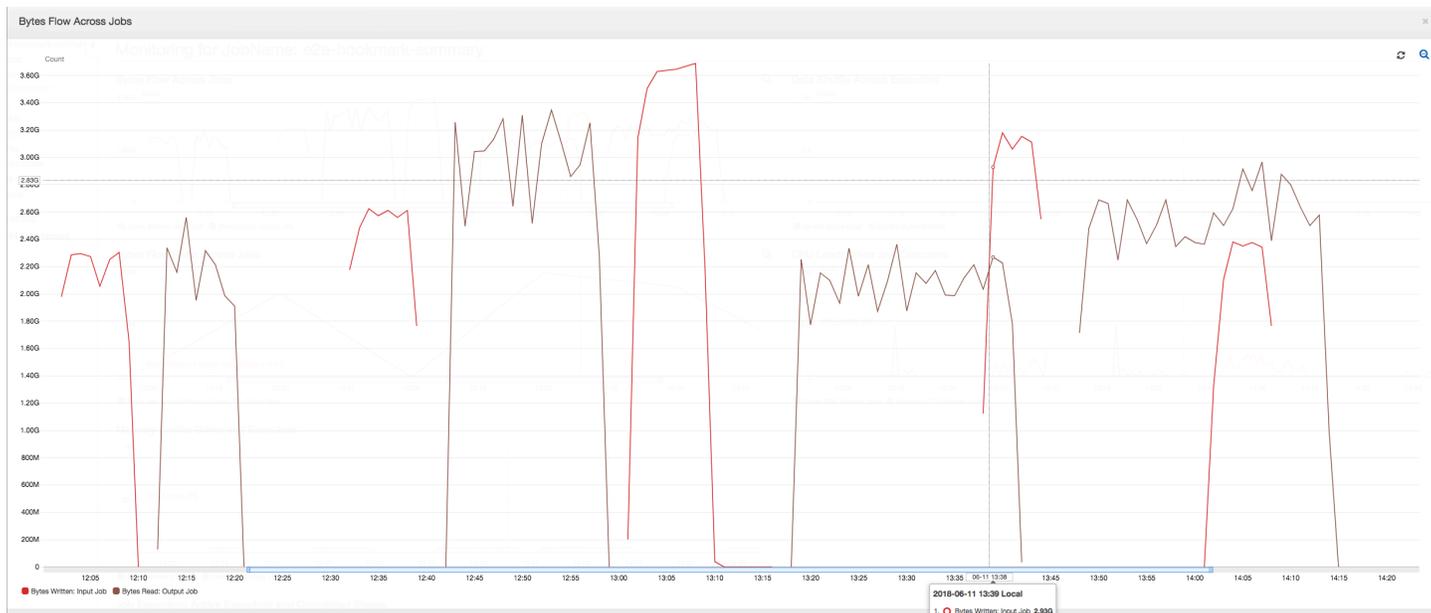
```
datasource0 = glueContext.create_dynamic_frame.from_options(connection_type="s3",
  connection_options = {"paths": ["s3://input_path"],
  "useS3ListImplementation":True,"recurse":True}, format="json")
applymapping1 = ApplyMapping.apply(frame = datasource0, mappings = [map_spec])
datasink2 = glueContext.write_dynamic_frame.from_options(frame = applymapping1,
  connection_type = "s3", connection_options = {"path": staging_path, "compression":
  "gzip"}, format = "json")
```

出力ジョブ: このジョブは、Amazon S3 内のステージング場所から入力ジョブの出力を読み取り、もう一度変換して、ターゲットの場所へ書き込みます。

```
datasource0 = glueContext.create_dynamic_frame.from_options(connection_type="s3",
  connection_options = {"paths": [staging_path],
  "useS3ListImplementation":True,"recurse":True}, format="json")
applymapping1 = ApplyMapping.apply(frame = datasource0, mappings = [map_spec])
datasink2 = glueContext.write_dynamic_frame.from_options(frame = applymapping1,
  connection_type = "s3", connection_options = {"path": output_path}, format = "json")
```

AWS Glue コンソールでプロファイルされたメトリクスを可視化する

次のダッシュボードでは、入力ジョブから取得された Amazon S3 バイトの書き込みメトリクスを、出力ジョブの同じタイムラインにある Amazon S3 バイトの読み取りメトリクスに重ね合わせています。タイムラインには、入力および出力ジョブのさまざまなジョブの実行が表示されます。入力ジョブ (赤で表示) は、30 分ごとに開始されます。出力ジョブ (茶色で表示) は、入力ジョブの完了時に開始し、最大同時実行数は 1 です。



この例では、[ジョブのブックマーク](#)が有効になっていません。スクリプトコードでジョブのブックマークを有効にするために使用される変換コンテキストはありません。

ジョブ履歴: [履歴] タブに示されているように、入力および出力ジョブには、午後 12 時に始まる複数の実行があります。

AWS Glue コンソールで入力ジョブは次のようになります。

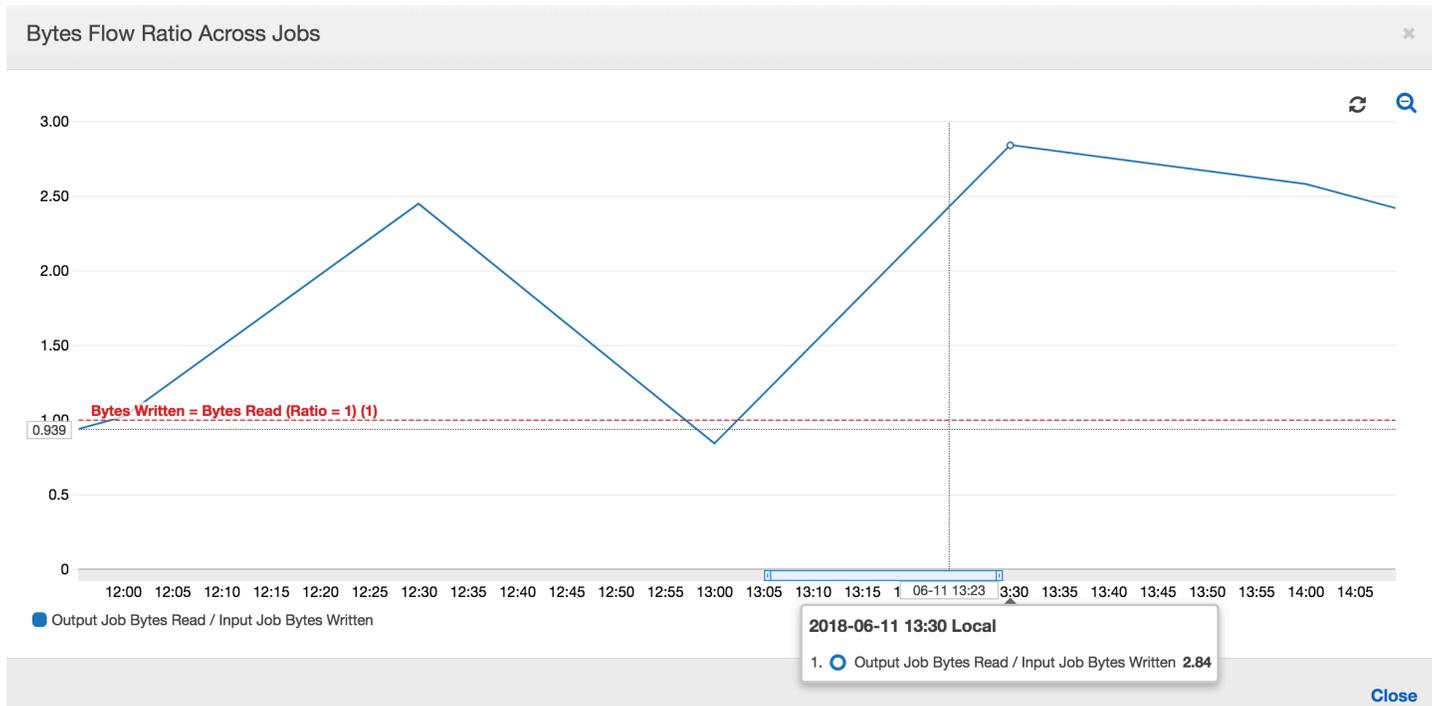
Run ID	Retry attempt	Run status	Error	Logs	Error logs	Execution time	Timeout	Delay	Triggered by	Start time	End time
j_0ce47b1a561051f06caae96e...	-	Succeeded		Logs		8 mins	2880 mins		e2e-bookmark-input	11 June 2018 2:30 PM UT...	11 June 2018 2:40 PM UT...
j_1b49ecdf73ad7614ccca2f4274...	-	Succeeded		Logs		8 mins	2880 mins		e2e-bookmark-input	11 June 2018 2:00 PM UT...	11 June 2018 2:10 PM UT...
j_07e4d5350ce516d89096821e...	-	Succeeded		Logs		7 mins	2880 mins		e2e-bookmark-input	11 June 2018 1:30 PM UT...	11 June 2018 1:46 PM UT...
j_fb9349097744be2arb655fb61...	-	Succeeded		Logs		15 mins	2880 mins		e2e-bookmark-input	11 June 2018 1:00 PM UT...	11 June 2018 1:16 PM UT...

以下のイメージに、出カジョブを示します。

Run ID	Retry attempt	Run status	Error	Logs	Error logs	Execution time	Timeout	Delay	Triggered by	Start time	End time
f_d2e5ba78770743d373d9dd63...	-	Failed	Max conc...	Logs	Error logs	0 secs	2880 mins		e2e-bookmark-output	11 June 2018 2:11 PM UT...	
f_3242babab08a6c6f0b5df2e3...	-	Succeeded		Logs		27 mins	2880 mins		e2e-bookmark-output	11 June 2018 1:47 PM UT...	11 June 2018 2:15 PM UT...
f_c98cc031be794a2b3a8047b...	-	Succeeded		Logs		24 mins	2880 mins		e2e-bookmark-output	11 June 2018 1:17 PM UT...	11 June 2018 1:43 PM UT...
f_0029a3c6f66c6395d59c9f965...	-	Succeeded		Logs		17 mins	2880 mins		e2e-bookmark-output	11 June 2018 12:41 PM U...	11 June 2018 12:59 PM U...

最初のジョブの実行: 以下のデータバイトの読み書きグラフに示されているように、12:00 から 12:30 の間に行われた入力および出カジョブの最初のジョブの実行は、曲線の下ほぼ同じ領域を示しています。これらの領域は、入力ジョブによって書き込まれた Amazon S3 バイトと、出カジョブによって読み取られた Amazon S3 バイトを表しています。このデータは、書き込まれた Amazon S3 バイトの比率によっても確認されます (30 分間の累計 - 入力ジョブのジョブトリガー頻度)。午後 12 時に開始された入力ジョブの実行の比率のデータポイントも 1 です。

以下のグラフは、すべてのジョブの実行におけるデータフロー比率を示しています。



2 回目のジョブの実行: 2 回目のジョブの実行では、入力ジョブによって書き込まれるバイト数と比較して、出カジョブにより読み取られるバイト数に明確な違いがあります。(出カジョブの 2 つのジョブの実行にまたがる曲線の下領域を比較するか、入力および出カジョブの 2 回目の実行の領域を比較します)。読み取りバイトと書き込みバイトの比率は、12:30 から 13:00 までの 2 番目の 30 分スパンにおいて入力ジョブによって書き込まれるデータの約 2.5 倍のデータを出カジョブが読み取ること示しています。ジョブのブックマークが有効になっていなかったため、出カジョブが入力

上部による最初のジョブの実行の出力を再処理したことがその理由です。1 より大きい比率は、出力ジョブにより処理された追加のデータバックログがあることを示しています。

3 回目のジョブの実行: 入力ジョブは、書き込まれるバイト数の観点でほぼ一定です (赤色の曲線の下にある領域を参照)。ただし、入力ジョブの 3 回目の実行には、予想されるよりも時間がかかります (赤色の曲線が長く続いていることを確認)。その結果、出力ジョブの 3 回目のジョブの実行が遅れて開始します。3 回目のジョブの実行では、13:00 から 13:30 までの残りの 30 分においてステージング場所に蓄積されたデータのごく一部のみ処理されました。バイトフローの比率は、入力ジョブの 3 回目の実行により書き込まれたデータの 0.83 のみ処理されたことを示しています (13:00 の比率を参照)。

入力ジョブと出力ジョブの重複: 入力ジョブの 4 回目のジョブの実行は、スケジュールに従って 13:30 に開始されました。これは、出力ジョブの 3 回目のジョブの実行が終了する前です。これらの 2 つのジョブの実行は一部の重複しています。ただし、出力ジョブの 3 回目の実行では、13:17 前後に開始されたとき、Amazon S3 のステージング場所にリストされたファイルのみキャプチャされます。これは、入力ジョブの最初のジョブの実行から取得されたすべてのデータ出力で構成されています。13:30 時点での実際の比率は約 2.75 です。出力ジョブの 3 回目のジョブの実行では、13:30 から 14:00 に入力ジョブの 4 回目のジョブの実行により書き込まれたデータの約 2.75 倍のデータが処理されました。

これらのイメージが示すように、出力ジョブは、入力ジョブのそれ以前のすべてのジョブの実行から取得されたステージング場所からのデータを再処理します。その結果、出力ジョブの 4 番目のジョブの実行が最も長くなり、入力ジョブの 5 回目のジョブの実行全体と重複します。

ファイルの処理を修正する

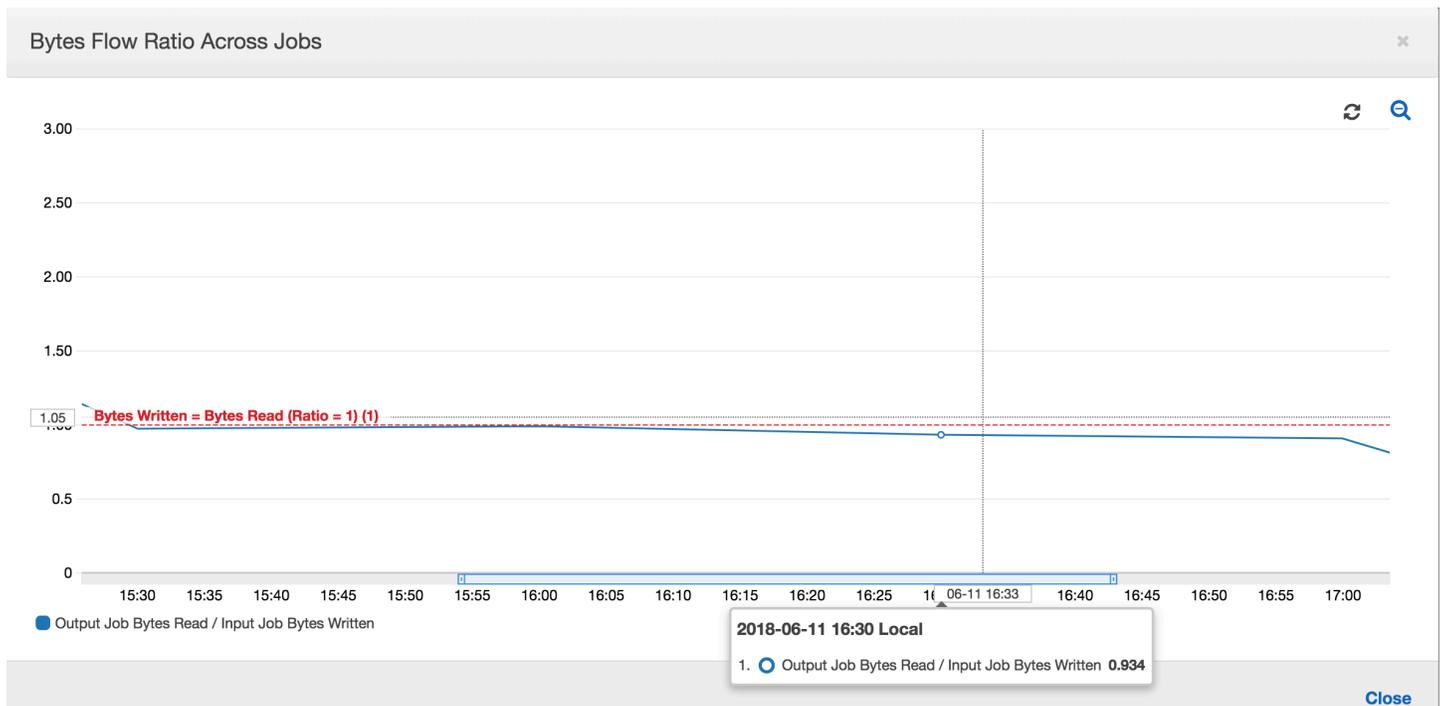
出力ジョブが、出力ジョブの以前のジョブの実行で処理されていないファイルのみ処理するようになる必要があります。これを行うには、次のように、ジョブのブックマークを有効にし、ジョブ出力に変換コンテキストを設定します。

```
datasource0 = glueContext.create_dynamic_frame.from_options(connection_type="s3",
  connection_options = {"paths": [staging_path],
  "useS3ListImplementation":True,"recurse":True}, format="json", transformation_ctx =
  "bookmark_ctx")
```

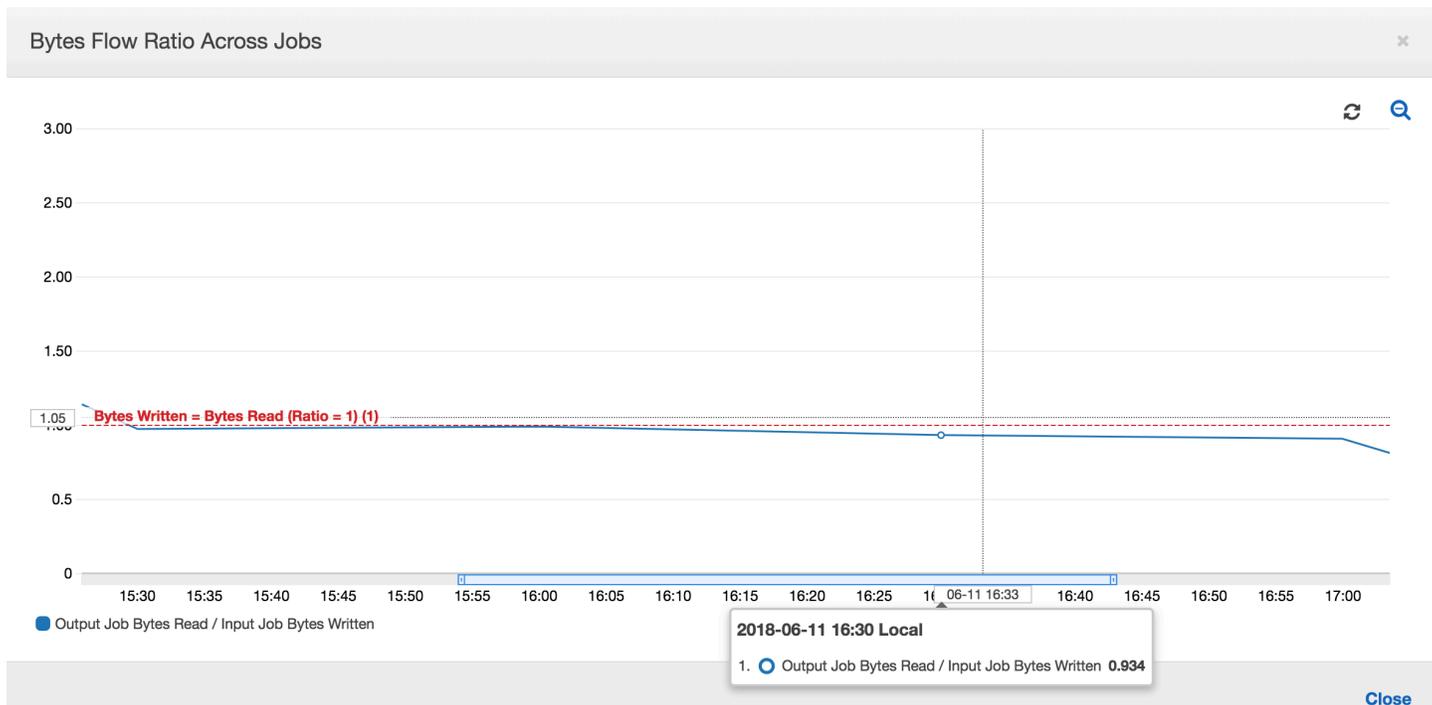
ジョブのブックマークが有効な場合、出力ジョブは、入力ジョブのそれまでのすべてのジョブの実行から取得されたステージング場所のデータを再処理しません。読み書きされたデータを示す以下のイメージでは、茶色の曲線の下にある領域はほぼ一定であり、赤色の曲線と似ています。



追加のデータが処理されないため、バイトフローの比率もほぼ 1 に維持されます。



出カジョブのジョブの実行が開始され、次の入カジョブの実行がデータをさらにステージング場所に置き始める前にステージング場所のファイルをキャプチャします。これを続けている限り、それまでの入カジョブの実行からキャプチャされたファイルのみ処理され、比率はほぼ 1 に維持されます。



入力ジョブに予想以上に時間がかかるため、2 回の入力ジョブの実行から取得されたステージング場所のファイルが出力ジョブによってキャプチャされるとします。その出力ジョブの実行では、比率が 1 よりも大きくなります。ただし、出力ジョブのそれ以降のジョブの実行では、出力ジョブのそれまでのジョブの実行により既に処理されているファイルは処理されません。

DPU の容量計画のモニタリング

AWS Glue でジョブメトリクスを使用すると、AWS Glue ジョブをスケールアウトするために使用できるデータ処理単位 (DPU) の数を予測できます。

Note

このページは AWS Glue バージョン 0.9 および 1.0 にのみ適用できます。AWS Glue 以降のバージョンには、容量計画時に追加の考慮事項を導入するコスト削減機能が含まれていません。

トピック

- [プロファイルされたコード](#)
- [AWS Glue コンソールでプロファイルされたメトリクスを可視化する](#)
- [最適な DPU 容量を決定する](#)

プロファイルされたコード

次のスクリプトは、428 個の gzip で圧縮された JSON ファイルを含む Amazon Simple Storage Service (Amazon S3) パーティションを読み取ります。このスクリプトは、マッピングを適用してフィールド名を変更し、Apache Parquet 形式に変換して Amazon S3 に書き込みます。デフォルトに従って 10 個の DPU をプロビジョニングし、このジョブを実行します。

```
datasource0 = glueContext.create_dynamic_frame.from_options(connection_type="s3",
  connection_options = {"paths": [input_path],
  "useS3ListImplementation":True,"recurse":True}, format="json")
applymapping1 = ApplyMapping.apply(frame = datasource0, mappings = [(map_spec)])
datasink2 = glueContext.write_dynamic_frame.from_options(frame = applymapping1,
  connection_type = "s3", connection_options = {"path": output_path}, format =
  "parquet")
```

AWS Glue コンソールでプロファイルされたメトリクスを可視化する

ジョブの実行 1: このジョブの実行では、プロビジョニングが不足している DPU がクラスターにあるかどうかを調べる方法を示します。AWS Glue でのジョブ実行機能は、[アクティブに実行されているエグゼキュターの合計数](#)、[完了したステージの数](#)、[必要なエグゼキュターの最大数](#)を表示します。

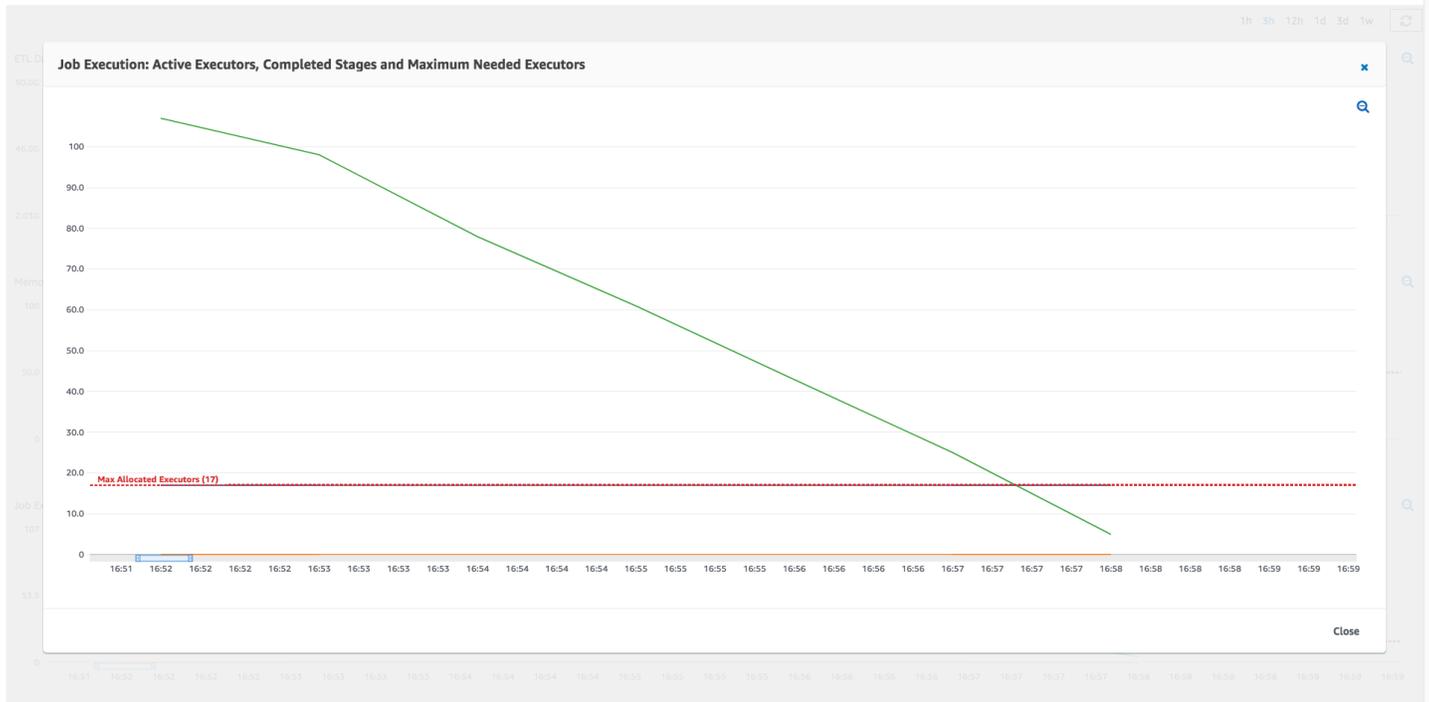
必要なエグゼキュターの最大数は、実行中のタスクと保留中のタスクの合計数を加算し、エグゼキュターごとのタスクで除算することによって算出されます。この結果は、現在の負荷に対応するために必要なエグゼキュターの合計数の尺度となります。

一方、アクティブに実行されているエグゼキュターの数、アクティブな Apache Spark タスクで実行されているエグゼキュターの数、を測定します。ジョブが進行するにつれて、必要なエグゼキュターの最大数は変化し、保留中のタスクキューが減るため通常はジョブの終わりに向かって減少します。

次のグラフの横方向の赤色の線は、割り当てられたエグゼキュターの最大数を示しています。これは、ジョブに割り当てる DPU の数によって異なります。この場合、ジョブの実行に対して 10 個の DPU を割り当てます。1 つの DPU が管理用に予約されています。9 個の DPU はそれぞれ 2 つのエグゼキュターを実行し、1 つのエグゼキュターは Spark ドライバー用に予約されています。Spark ドライバーはプライマリアプリケーション内で実行されます。そのため、割り当てられるエグゼキュターの最大数は、 $2 \times 9 - 1 = 17$ です。

Jobs > e2e-dpus

Detailed job metrics

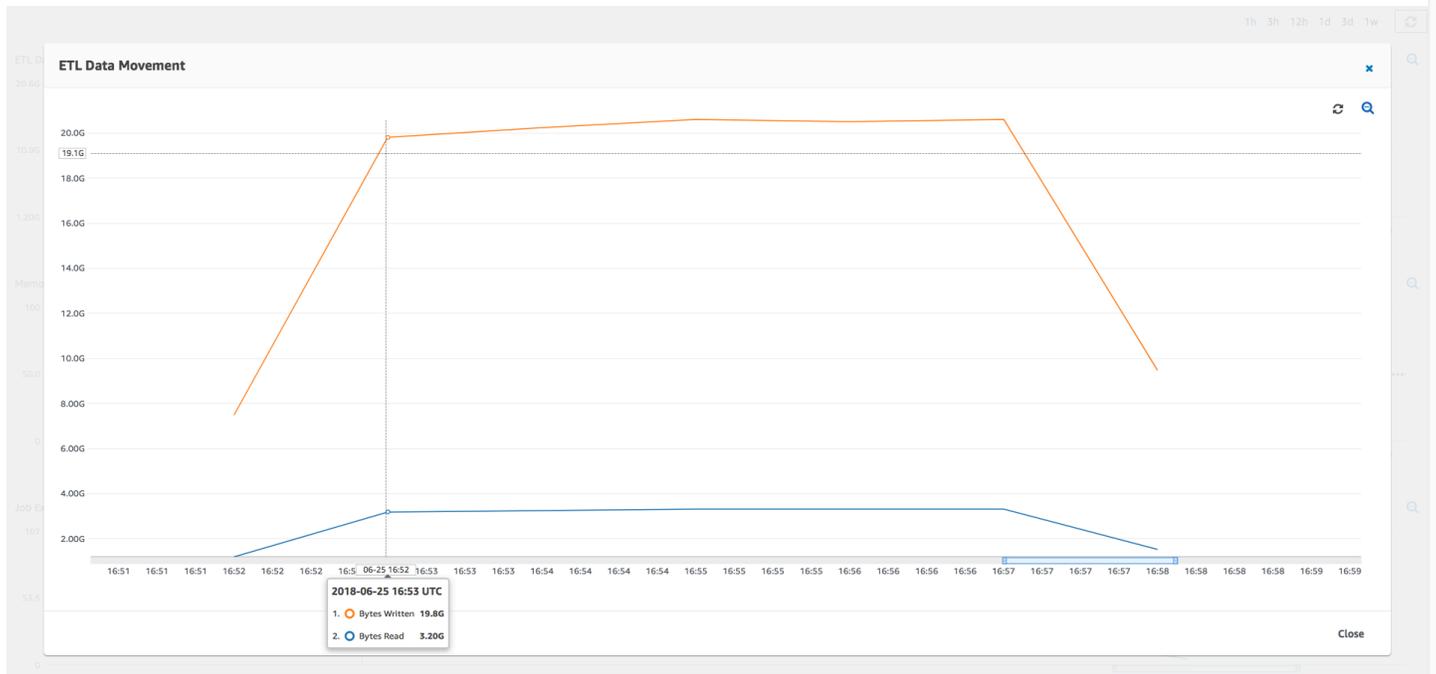


グラフが示すとおり、必要なエグゼキュターの最大数はジョブの開始時に 107 から始まりますが、アクティブなエグゼキュターの数 は 17 のままです。これは、10 個の DPU を持つ割り当てられるエグゼキュターの最大数と同じです。必要なエグゼキュターの最大数と割り当てられるエグゼキュターの最大数の比率 (Spark ドライバーでは両方に 1 を加算) から、プロビジョニングが不足している係数が $108/18 = 6$ 倍であるとわかります。6 (プロビジョニング率未満) * 9 (現在の DPU 容量 - 1) + 1 DPU = 55 個の DPU をプロビジョニングしてジョブをスケールアウトし、最大限の並列処理で実行することで処理時間を短縮することができます。

AWS Glue コンソールは、詳細なジョブメトリクスを、元の割り当てられるエグゼキュターの最大数を表す静的な線として表示します。コンソールは、メトリクスのジョブ定義から割り当てられるエグゼキュターの最大数を計算します。対照的に、詳細なジョブ実行メトリクスについては、コンソールはジョブ実行設定から割り当てられるエグゼキュターの最大数、特にジョブ実行に割り当てられた DPU を計算します。個々のジョブ実行のメトリクスを表示するには、ジョブ実行を選択して、[実行メトリクスの表示] を選択します。

Jobs > e2e-dpus

Detailed job metrics



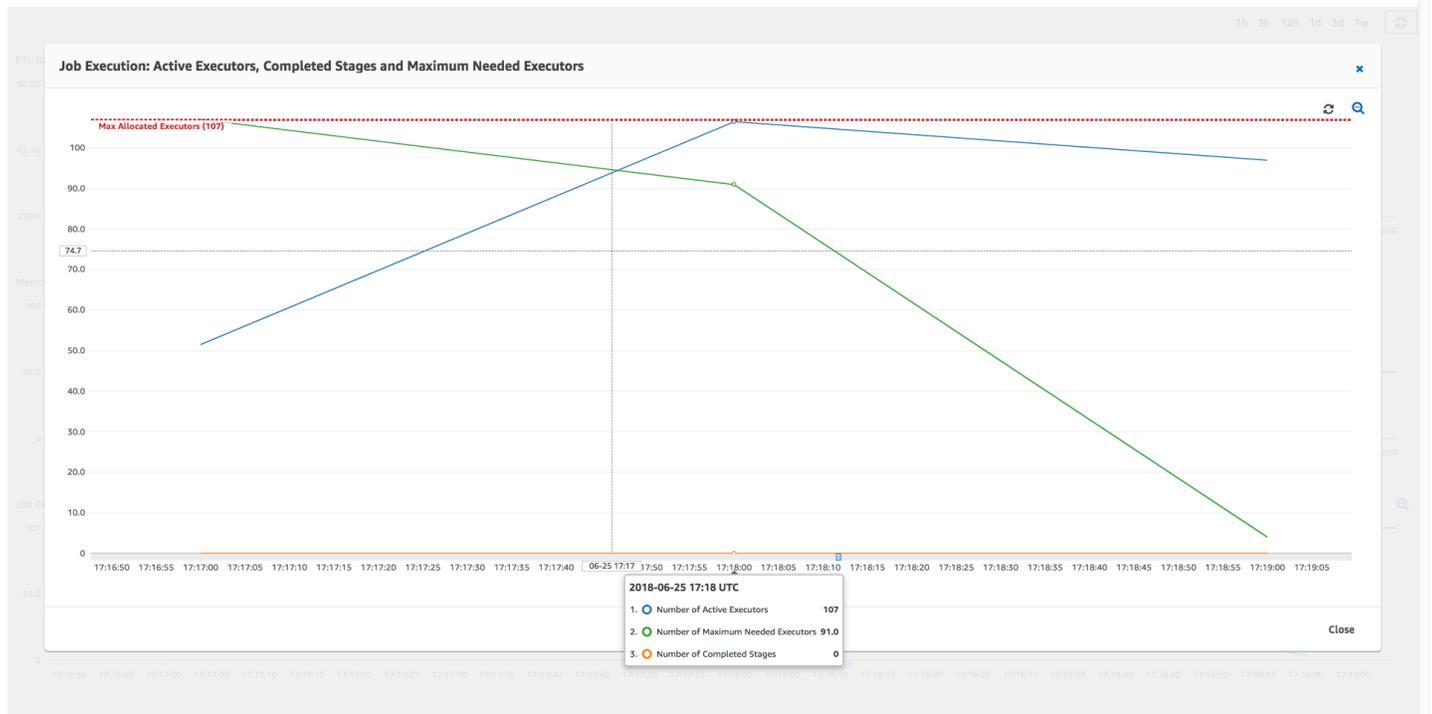
読み取りおよび書き込みされた Amazon S3 バイトを見ると、ジョブが Amazon S3 からのデータをストリーミングして並列に書き出すことに 6 分すべてを消費していることがわかります。割り当てられた DPU 上のすべてのコアは、Amazon S3 に対して読み取りと書き込みを行います。必要なエグゼキュターの最大数 107 は、入力 Amazon S3 パス内のファイル数 428 とも整合しています。各エグゼキュターは、4 つの Spark タスクを起動し、4 つの入力ファイルを処理します (gzip で圧縮された JSON)。

最適な DPU 容量を決定する

前のジョブの実行の結果に基づいて、割り当てられる DPU の合計数を 55 に増やし、ジョブがどのように実行されるかを確認できます。ジョブは 3 分未満で完了します。以前に必要な時間の半分です。この場合、実行時間の短いジョブであるため、ジョブのスケールアウトは直線的ではありません。存続期間の長いタスクまたは多数のタスクを持つジョブ (必要なエグゼキュターの最大数が大きい) は、直線的に近い DPU スケールアウトのパフォーマンス向上から恩恵を受けます。

Jobs > e2e-dpus

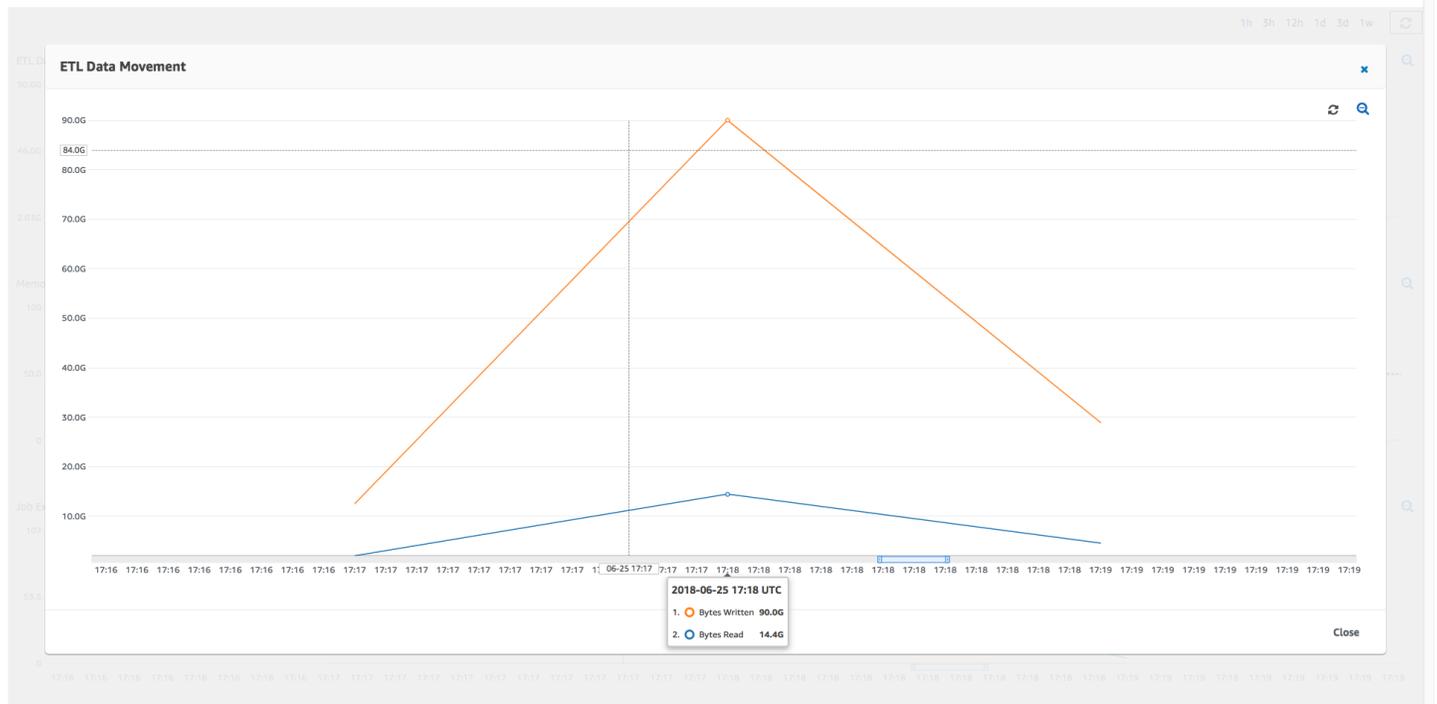
Detailed job metrics



上の図に示すように、アクティブなエグゼキューターの合計数が最大割り当て数 107 個のエグゼキューターに到達します。同様に、必要なエグゼキューターの最大数が割り当てられるエグゼキューターの最大数を上回ることはありません。必要なエグゼキューターの最大数は、アクティブに実行されているタスクと保留中のタスクの数から計算されるため、アクティブなエグゼキューターの数よりも小さい可能性があります。これは、短期間部分的または完全にアイドル状態になり、まだ停止されていないエグゼキューターが存在する可能性があるためです。

Jobs > e2e-dpus

Detailed job metrics



このジョブの実行では、Amazon S3 から並列で読み取りと書き込みをするために 6 倍のエグゼキューターが使用されます。その結果、このジョブの実行では、読み取りと書き込みの両方に使用される Amazon S3 帯域幅が多くなり、早く完了します。

過度にプロビジョニングされた DPU を識別する

次に、100 個の DPU (99 * 2 = 198 個のエグゼキューター) でジョブをスケールアウトするとこれ以上スケールアウトしなくてよくなるかどうかを調べることができます。次のグラフが示すように、ジョブが完了するまでまだ 3 分かかります。同様に、ジョブは 107 個を超えてスケールアウトしないため (55 DPU 構成)、残りの 91 個のエグゼキューターは過剰にプロビジョニングされてまったく使用されません。必要なエグゼキューターの最大数からわかるように、これは DPU の数を増やしても必ずパフォーマンスが向上するわけではないことを示しています。

Jobs > e2e-dpus

Detailed job metrics



時間の差を比較する

次の表に示す 3 つのジョブが実行は、10 個の DPU、55 個の DPU、100 個の DPU のジョブの実行時間をまとめています。最初のジョブの実行をモニタリングすることで確立した推定値を使用して、ジョブの実行時間を向上させることができる DPU 容量を調べることができます。

ジョブ ID	DPU の数	実行時間
jr_c894524c8ef5048a4d9...	10	6 分。
jr_1a466cf2575e7ffe6856...	55	3 分。
jr_34fa1ed4c6aa9ff0a814...	100	3 分。

AWS Glue でのストリーミング ETL ジョブ

連続的に実行されるストリーミング抽出/変換/ロード (ETL) ジョブを作成し、Amazon Kinesis Data Streams、Apache Kafka、Amazon Managed Streaming for Apache Kafka (Amazon MSK) などのストリーミングソースからのデータを使用できます。ジョブはデータをクレンジングして変換し、その結果を Amazon S3 データレイクまたは JDBC データストアにロードします。

さらに、Amazon Kinesis Data Streams ストリーム用のデータを生成できます。この機能は、AWS Glue スクリプトを記述する場合にのみ使用できます。詳細については、「[the section called “Kinesis 接続”](#)」を参照してください。

デフォルトでは、AWS Glue によるデータ処理と書き出しは 100 秒ウィンドウ単位で行われます。これにより、データを効率的に処理しつつ、想定より遅く到着したデータに対する集計を実行できます。このウィンドウサイズを変更して、タイムリーで高精度の集計にできます。AWS Glue ストリーミングジョブでは、ジョブブックマークではなくチェックポイントを使用して、読み取られたデータを追跡します。

Note

AWS Glue では、実行中の ETL ジョブのストリーミングに対して 1 時間ごとに課金します。

この動画では、ストリーミング ETL のコスト課題と、 のコスト削減機能について説明します AWS Glue。

ストリーミング ETL ジョブを作成するには、次の手順を実行します。

1. Apache Kafka ストリーミングソースの場合は、Kafka ソースまたは Amazon MSK クラスターへの AWS Glue 接続を作成します。
2. ストリーミングソースの Data Catalog テーブルを手動で作成します。
3. ストリーミングデータソースの ETL ジョブを作成します。ストリーミング固有のジョブプロパティを定義し、独自のスクリプトを指定するか、生成されたスクリプトを必要に応じて変更します。

詳細については、「[AWS Glue でのストリーミング ETL](#)」を参照してください。

Amazon Kinesis Data Streams にストリーミング ETL ジョブを作成する場合、AWS Glue 接続を作成する必要はありません。ただし、Kinesis Data Streams をソースとする AWS Glue ストリーミング ETL ジョブへの接続がある場合は、 の場合は、Kinesis への Virtual Private Cloud (VPC) エンドポイントが必要です。詳細については、Amazon VPC ユーザーガイドの [インターフェイスエンドポイントの作成](#) を参照してください。別のアカウントで Amazon Kinesis Data Streams ストリームを指定する場合は、アカウント間アクセスを許可するようにロールとポリシーを設定する必要があります。詳細については、「[Example: Read From a Kinesis Stream in a Different Account](#)」を参照してください。

AWS Glue ストリーミング ETL ジョブにより、圧縮データの自動検出が行えます。また、ストリーミングデータを透過的に解凍し、入力ソースに対して一般的な変換を実行して、出力ストアにロードすることができます。

AWS Glue は、入力形式に対して指定された、以下の圧縮タイプのための自動解凍をサポートしています。

圧縮タイプ	Avro ファイル	Avro データム	JSON	CSV	Grok
BZIP2	はい	はい	はい	はい	はい
GZIP	いいえ	はい	はい	はい	はい
SNAPPY	はい (生の Snappy)	はい (フレーム付 Snappy)	はい (フレーム付 Snappy)	はい (フレーム付 Snappy)	はい (フレーム付 Snappy)
XZ	はい	はい	はい	はい	はい
ZSTD	はい	いいえ	いいえ	いいえ	いいえ
DEFLATE	はい	はい	はい	はい	はい

トピック

- [Apache Kafka データストリームの AWS Glue 接続の作成](#)
- [ストリーミングソースの Data Catalog テーブルの作成](#)
- [Avro ストリーミングソースに関する注意事項と制約事項](#)
- [ストリーミングソースへの grok パターンの適用](#)
- [ストリーミング ETL ジョブのジョブプロパティの定義](#)
- [ストリーミング ETL に関する注意と制限](#)

Apache Kafka データストリームの AWS Glue 接続の作成

Apache Kafka ストリームから読み取りを行うには、AWS Glue 接続を作成する必要があります。

Kafka ソース用の AWS Glue 接続を作成するには (コンソール)

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソールを開きます。
2. ナビゲーションペインの [Data catalog] で [Connections] (接続) を選択します。
3. [接続の追加] を選択し、[接続のプロパティを設定します] ページで接続名を入力します。

Note

接続プロパティ指定の詳細については、「[AWS Glue connection properties.](#)」(の接続プロパティ) を参照してください。

4. [接続タイプ] で、[Kafka] を選択します。
5. Kafka のブートストラップサーバーの URL として、Amazon MSK クラスターまたは Apache Kafka クラスターのブートストラップブローカーのホストとポート番号を入力します。Kafka クラスターへの初期接続を確立するには、Transport Layer Security (TLS) エンドポイントのみを使用します。Plaintext エンドポイントはサポートされていません。

Amazon MSK クラスターのホスト名とポート番号のペアのリスト例を次に示します。

```
myserver1.kafka.us-east-1.amazonaws.com:9094,myserver2.kafka.us-east-1.amazonaws.com:9094,  
myserver3.kafka.us-east-1.amazonaws.com:9094
```

ブートストラップブローカー情報の取得の詳細については、Amazon Managed Streaming for Apache Kafka デベロッパーガイドの「[Getting the Bootstrap Brokers for an Amazon MSK Cluster](#)」を参照してください。

6. Kafka データソースへのセキュアな接続が必要な場合は、[Require SSL connection] (SSL 接続が必要) を選択し、[Kafka private CA certificate location] (Kafka のプライベート CA 証明書の場合) に、カスタム SSL 証明書への有効な Amazon S3 パスを入力します。

自己管理型 Kafka への SSL 接続の場合、カスタム証明書は必須です。Amazon MSK ではオプションです。

Kafka のカスタム証明書を指定する方法の詳細については、「[the section called “SSL 接続プロパティ”](#)」を参照してください。

7. AWS Glue Studio または AWS CLI を使用して Kafka クライアント認証方法を指定します。左側のナビゲーションペインの ETL メニュー AWS Glue Studio AWS Glue から選択するには

Kafka クライアントの認証方法の詳細については、「[クライアント認証用の AWS Glue Kafka 接続プロパティ](#)」を参照してください。

- 必要に応じて説明を入力し、[Next] (次へ) をクリックします。
- Amazon MSK クラスターの場合は、その Virtual Private Cloud (VPC)、サブネット、およびセキュリティグループを指定します。VPC 情報は、自己管理型 Kafka の場合はオプションです。
- [Next] (次へ) をクリックして、すべての接続プロパティを確認し、[Finish] (完了) をクリックします。

AWS Glue 接続の詳細については、「[データへの接続](#)」を参照してください。

クライアント認証用の AWS Glue Kafka 接続プロパティ

SASL/GSSAPI (Kerberos) 認証

この認証方法を選択すると、Kerberos のプロパティを指定できます。

Kerberos のキータブ

キータブファイルの場所を選択します。キータブには、1 つ以上のプリンシパルの長期キーが保存されます。詳細については、[MIT Kerberos ドキュメント: キータブ](#)を参照してください。

Kerberos の krb5.conf ファイル

krb5.conf ファイルを選択します。これには、デフォルトの領域 (同じ KDC のシステムのグループを定義する論理ネットワークで、ドメインに類似しています) と KDC サーバーの場所が含まれます。詳細については、[MIT Kerberos ドキュメント: krb5.conf](#)を参照してください。

Kerberos のプリンシパルと Kerberos のサービス名

Kerberos のプリンシパルとサービス名を入力します。詳細については、「[MIT Kerberos Documentation: Kerberos principal](#)」(MIT Kerberos ドキュメント: Kerberos のプリンシパル) を参照してください。

SASL/SCRAM-SHA-512 認証

この認証方法を選択すると、認証情報を指定することができます。

AWS シークレットマネージャー

[Search] (検索) ボックスに名前または ARN を入力して、トークンを検索します。

ユーザー名とパスワードを直接提供する

[Search] (検索) ボックスに名前または ARN を入力して、トークンを検索します。

SSL クライアント認証

この認証方法を選択すると、Amazon S3 を参照することによって Kafka クライアントキーストアの場所を選択できます。オプションで、Kafka クライアントキーストアのパスワードと Kafka クライアントキーのパスワードを入力できます。

IAM 認証

この認証方法は追加の仕様を必要とせず、ストリーミングソースが MSK Kafka の場合にのみ適用されます。

SASL/PLAIN 認証

この認証方法を選択すると、認証情報を指定することができます。

ストリーミングソースの Data Catalog テーブルの作成

データスキーマなど、ソースデータストリームのプロパティを指定するデータカタログテーブルを手動でストリーミングソースに対して作成できます。このテーブルは、ストリーミング ETL ジョブのデータソースとして使用されます。

ソースデータストリーム内のデータのスキーマがわからない場合は、スキーマなしでテーブルを作成できます。次に、ストリーミング ETL ジョブを作成するときに、AWS Glue スキーマ検出機能を有効にできます。AWS Glue では、ストリーミングデータからスキーマを決定します。

[AWS Glue コンソール](#)、AWS Command Line Interface (AWS CLI)、または AWS Glue API を使用してテーブルを作成します。AWS Glue コンソールを使用して手動でテーブルを作成する方法については、「[the section called “テーブルの作成”](#)」を参照してください。

Note

AWS Lake Formation コンソールを使用してテーブルを作成することはできません。AWS Glue コンソールを使用する必要があります。

また、Avro 形式のストリーミングソースや Grok パターンを適用できるログデータについては、次の情報を考慮します。

- [the section called “Avro ストリーミングソースに関する注意事項と制約事項”](#)
- [the section called “ストリーミングソースへの grok パターンの適用”](#)

トピック

- [Kinesis データソース](#)
- [Kafka データソース](#)
- [AWS Glue スキーマレジストリテーブルのソース](#)

Kinesis データソース

テーブルを作成するときは、次のストリーミング ETL プロパティを設定します (コンソール)。

ソースのタイプ

Kinesis

同じアカウント内の Kinesis ソースの場合:

リージョン

Amazon Kinesis Data Streams サービスが存在する AWS リージョン。リージョンと Kinesis ストリーム名は一緒にストリーム ARN に変換されます。

例: <https://kinesis.us-east-1.amazonaws.com>

Kinesis ストリーム名

Amazon Kinesis Data Streams デベロッパーガイドの「[Creating a Stream](#)」に記載されているストリーム名。

別のアカウントの Kinesis ソースについては、[この例](#) を参照して、アカウント間アクセスを許可するようにロールとポリシーを設定します。次の設定をします。

ストリーム ARN

コンシューマーを登録する Kinesis データストリームの ARN。詳細については、「[」の ARNs](#)」および AWS 「[サービス名前空間](#)」を参照してくださいAWS 全般のリファレンス。

引き受けたロールの ARN

引き受けるロールの Amazon リソースネーム (ARN)。

セッション名 (オプション)

引き受けたロールセッションの識別子。

異なるプリンシパルによって同じロールが引き継がれた場合にセッションを一意に識別するために、またはその他の理由で、ロールセッション名を使用します。クロスアカウントシナリオでは、ロールセッション名が表示され、ロールを所有するアカウントでログ記録できます。引

き受けたロールプリンシパルの ARN では、ロールセッション名も使用されません。つまり、一時的なセキュリティ認証情報を使用する後続のクロスアカウント API リクエストは、ロールセッション名を AWS CloudTrail ログ内の外部アカウントに公開しません。

Amazon Kinesis Data Streams のストリーミング ETL プロパティを設定するには (AWS Glue API または AWS CLI)

- 同じアカウント内の Kinesis ソースのストリーミング ETL プロパティを設定するには、CreateTable API オペレーションまたは create_table CLI コマンドの StorageDescriptor データ構造の streamName および endpointUrl パラメータを指定します。

```
"StorageDescriptor": {  
  "Parameters": {  
    "typeOfData": "kinesis",  
    "streamName": "sample-stream",  
    "endpointUrl": "https://kinesis.us-east-1.amazonaws.com"  
  }  
  ...  
}
```

または、streamARN を指定します。

Example

```
"StorageDescriptor": {  
  "Parameters": {  
    "typeOfData": "kinesis",  
    "streamARN": "arn:aws:kinesis:us-east-1:123456789:stream/sample-stream"  
  }  
  ...  
}
```

- 別のアカウントの Kinesis ソースのストリーミング ETL プロパティを設定するには、CreateTable API オペレーションまたは create_table CLI コマンドの StorageDescriptor データ構造の streamARN、awsSTSRoleARN、および awsSTSSessionName (オプション) パラメータを指定します。

```
"StorageDescriptor": {
```

```
"Parameters": {
  "typeOfData": "kinesis",
  "streamARN": "arn:aws:kinesis:us-east-1:123456789:stream/sample-stream",
  "awsSTSRoleARN": "arn:aws:iam::123456789:role/sample-assume-role-arn",
  "awsSTSSessionName": "optional-session"
}
...
}
```

Kafka データソース

テーブルを作成するときは、次のストリーミング ETL プロパティを設定します (コンソール)。

ソースのタイプ

Kafka

Kafka ソースの場合:

トピック名

Kafka で指定されたトピック名。

Connection

Kafka ソースを参照する AWS Glue 接続 ([「the section called “Kafka データストリームの接続の作成”」](#)を参照)。

AWS Glue スキーマレジストリテーブルのソース

AWS Glue スキーマレジストリをストリーミングジョブに使用するには、[ユースケース: AWS Glue Data Catalog](#) スキーマレジストリのテーブルを作成または更新する手順に従ってください。

現在、AWS Glue ストリーミングは、false にスキーマ推論が設定された Glue スキーマレジストリ Avro 形式のみをサポートします。

Avro ストリーミングソースに関する注意事項と制約事項

Avro 形式のストリーミングソースには、次の注意事項と制限事項が適用されます。

- スキーマの検出が有効になっている場合、Avro スキーマをペイロードに含める必要があります。無効にすると、ペイロードにはデータのみが含まれます。

- 一部の Avro データ型は、ダイナミックフレームでサポートされていません。AWS Glue コンソールの [create table] (テーブル作成) ウィザードの [Define a schema] (スキーマの定義) でスキーマを定義するときに、これらのデータ型を指定することはできません。スキーマの検出中に、Avro スキーマでサポートされていない型は、次のようにサポートされている型に変換されます。
 - EnumType => StringType
 - FixedType => BinaryType
 - UnionType => StructType
- コンソールの [Define a schema] (スキーマの定義) ページでテーブルスキーマを定義する場合は、スキーマの暗黙のルート要素タイプは record です。record 以外のルート要素タイプ (array や map など) が必要な場合、[Define a schema] (スキーマの定義) ページでスキーマを指定することはできません。代わりに、そのページをスキップして、スキーマをテーブルプロパティとして、または ETL スクリプト内で指定する必要があります。
- テーブルのプロパティでスキーマを指定するには、テーブル作成ウィザードを完了し、テーブルの詳細を編集し、[Table properties] (テーブルのプロパティ) に新しいキーと値のペアを追加します。avroSchema キーを使用して、次のスクリーンショットに示すように、値のスキーマ JSON オブジェクトを入力します。

Edit table details

Key

Value

Description

Table properties

Key	Value	
classification	avro	✕
avroSchema	{"type": "array", "items": "strin	✕
<input type="text"/>	<input type="text"/>	

- ETL スクリプトでスキーマを指定するには、次の Python と Scala の例に示すように、datasource0 代入ステートメントを修正し、avroSchema キーを additional_options 引数に追加します。

Python

```
SCHEMA_STRING = '{"type": "array", "items": "string"}'
datasource0 = glueContext.create_data_frame.from_catalog(database =
    "database", table_name = "table_name", transformation_ctx = "datasource0",
    additional_options = {"startingPosition": "TRIM_HORIZON", "inferSchema":
    "false", "avroSchema": SCHEMA_STRING})
```

Scala

```
val SCHEMA_STRING = """"{"type": "array", "items": "string"}""""
val datasource0 = glueContext.getCatalogSource(database = "database", tableName
    = "table_name", redshiftTmpDir = "", transformationContext = "datasource0",
```

```
additionalOptions = JsonOptions(s""{"startingPosition": "TRIM_HORIZON",
"inferSchema": "false", "avroSchema": "$SCHEMA_STRING"}""").getDataFrame()
```

ストリーミングソースへの grok パターンの適用

ログデータソース用のストリーミング ETL ジョブを作成し、Grok パターンを使用してログを構造化データに変換できます。その後、ETL ジョブは、データを構造化データソースとして処理します。ストリーミングソースの Data Catalog テーブルを作成するときに適用する Grok パターンを指定します。

Grok パターンとカスタムパターン文字列値については、「[Grok カスタム分類子の書き込み](#)」を参照してください。

Grok パターンを Data Catalog テーブルに追加するには (コンソール)

- [create table] (テーブル作成) ウィザードを使用し、[the section called “ストリーミングソースの Data Catalog テーブルの作成”](#) に指定されたパラメータでテーブルを作成します。データ形式を Grok として指定し、[Grok pattern] (Grok パターン) フィールドに入力し、必要に応じてカスタムパターンを [Custom patterns (optional)] (カスタムパターン (オプション)) に追加します。

The screenshot shows a wizard titled "Choose a data format". Under the "Classification" section, "Grok" is selected with a radio button. Below this, there is a text input field for the "Grok pattern" containing the placeholder text "%{PATTERN-NAME:field-name:optional-data-type}". A note below the input field states: "Built-in and custom named patterns used to parse your data into a structured schema. For more information, see the list of built-in patterns." Under the "Custom patterns" section, there is a table with one column and one row containing the number "1". Below the table, it says "Optional custom building blocks for the grok pattern." At the bottom of the wizard, there are "Back" and "Next" buttons.

各カスタムパターンの後に、Enter を押します。

Grok パターンを Data Catalog テーブルに追加するには (AWS Glue API または AWS CLI)

- GrokPattern パラメータと、必要に応じて CustomPatterns パラメータを CreateTable API オペレーションまたは create_table CLI コマンドに追加します。

```
"Parameters": {  
  ...  
  "grokPattern": "string",  
  "grokCustomPatterns": "string",  
  ...  
},
```

grokCustomPatterns を文字列として表現し、パターン間の区切りとして「\n」を使用します。

これらのパラメータの指定の例を次に示します。

Example

```
"parameters": {  
  ...  
  "grokPattern": "%{USERNAME:username} %{DIGIT:digit:int}",  
  "grokCustomPatterns": "digit \d",  
  ...  
}
```

ストリーミング ETL ジョブのジョブプロパティの定義

AWS Glue コンソールでストリーミング ETL ジョブを定義する場合は、次に示すストリーム固有のプロパティを指定します。その他のジョブプロパティの説明については、「[Spark ジョブのジョブプロパティの定義](#)」を参照してください。

IAM ロール

ジョブの実行、ストリーミングソースへのアクセス、ターゲットデータストアへのアクセスに使用されるリソースの承認に使用される AWS Identity and Access Management (IAM) ロールを指定します。

Amazon Kinesis Data Streams にアクセスするには、AmazonKinesisFullAccess AWS 管理ポリシーをロールにアタッチするか、よりきめ細かなアクセスを許可する同様の IAM ポリシー

をアタッチします。サンプルポリシーについては、「[Controlling Access to Amazon Kinesis Data Streams Resources Using IAM](#)」を参照してください。

AWS Glue でジョブを実行するためのアクセス権限の詳細については、[AWS Glue の Identity and Access Management](#) を参照してください。

タイプ

[Spark Streaming] を選択します。

AWS Glue バージョン

AWS Glue バージョンによって、ジョブで使用できる Apache Spark、および Python または Scala のバージョンが決まります。ジョブで使用可能な Python または Scala のバージョンを指定する選択肢を選択します。AWS Glueバージョン 2.0 (Python 3 をサポート) は、ストリーミング ETL ジョブのデフォルトです。

メンテナンスウィンドウ

ストリーミングジョブを再起動できるウィンドウを指定します。[the section called “メンテナンスウィンドウ”](#) を参照してください。

ジョブのタイムアウト

必要に応じて、長さを分単位で入力します。デフォルト値は空白です。

- ストリーミングジョブのタイムアウト値は 7 日または 10080 分未満である必要があります。
- 値を空白のままにすると、メンテナンスウィンドウを設定していない場合、ジョブは 7 日後に再開されます。メンテナンスウィンドウを設定している場合、ジョブはメンテナンスウィンドウ中に 7 日後に再起動されます。

データソース

「[the section called “ストリーミングソースの Data Catalog テーブルの作成”](#)」で作成したテーブルを指定します。

データターゲット

次のいずれかを行います。

- [データターゲットでテーブルを作成する] を選択し、次に示すデータターゲットプロパティを指定します。

データストア

Amazon S3 または JDBC を選択します。

[形式]

任意の形式を選択します。ストリーミングではすべてがサポートされています。

- [Use tables in the data catalog and update your data target] (データカタログのテーブルを使用し、データターゲットを更新する) を選択し、JDBC データストアのテーブルを選択します。

出力スキーマ定義

次のいずれかを行います。

- [Automatically detect schema of each record] (各レコードのスキーマを自動的に検出) を選択してスキーマ検出を有効にします。AWS Glue は、ストリーミングデータからスキーマを決定します。
- [Specify output schema for all records] (すべてのレコードの出力スキーマを指定する) を選択して、[Apply Mapping] (マッピングの適用) 変換を使用して出力スキーマを定義します。

スクリプト

必要に応じて、独自のスクリプトを指定するか、生成されたスクリプトを変更して、Apache Spark 構造化ストリーミングエンジンでサポートされる操作を実行します。使用可能なオペレーションの詳細については、[「ストリーミング DataFrames/データセットのオペレーション」](#)を参照してください。

ストリーミング ETL に関する注意と制限

注意事項と制限事項を次に示します。

- AWS Glue ストリーミング ETL ジョブでの自動解凍は、サポートされた圧縮タイプに対してのみ使用できます。また、以下の点にも注意してください。
 - フレーム付 Snappy とは、公式の[フレーミング形式](#)を使用する Snappy を指します。
 - Deflate がサポートされるのは、Glue バージョン 3.0 です (Glue バージョン 2.0 ではサポートされません)。
- スキーマ検出を使用する場合、ストリーミングデータの結合を実行できません。
- AWS Glue のストリーミング ETL ジョブは、Avro 形式の AWS Glue スキーマレジストリに対する Union データ型をサポートしません。
- ETL スクリプトでは、AWS Glue の組み込み変換と Apache Spark 構造化ストリーミングのネイティブな変換を使用できます。詳細については、Apache Spark [ウェブサイトの「ストリーミング DataFrames/データセットの操作」](#)または「」を参照してください[AWS Glue PySpark 変換リファレンス](#)。

- AWS Glue ストリーミング ETL ジョブでは、チェックポイントを使用して、読み取ったデータの追跡が行われます。このため、停止して再起動されたジョブは、ストリーム内で中断した位置から開始されます。データを再処理する場合は、スクリプト内で参照されているチェックポイントフォルダを削除することができます。
- ジョブのブックマークはサポートされていません。
- ジョブで Kinesis Data Streams の拡張ファンアウト機能を使用するには、「[the section called “Kinesis ストリーミングジョブでの拡張ファンアウトの使用”](#)」を参照してください。
- AWS Glue スキーマレジストリ から作成されたデータカタログテーブルを使用する場合、新しいスキーマバージョンが使用可能になったら、新しいスキーマを反映するには、次の操作を行う必要があります。
 1. テーブルに関連するジョブを停止します。
 2. [データカタログ] テーブルのスキーマを更新します。
 3. テーブルに関連するジョブを再起動します。

AWS Lake Formation FindMatches によるレコードのマッチング

Note

レコードの照合は現在、次のリージョンのAWS Glue コンソールではご利用いただけません: 中東 (アラブ首長国連邦)、欧州 (スペイン) (eu-south-2)、欧州 (チューリッヒ (eu-central-2))。

AWS Lake Formation には、カスタム変換を作成するための機械学習機能が用意されています。現在、FindMatches という名前の変換が 1 つあります。FindMatches 変換を使用すると、レコードに共通の一意の識別子がなく、正確に一致するフィールドがない場合でも、データセット内の重複レコードまたは一致するレコードを識別できます。そのため、コードを記述したり、機械学習の仕組みを知ったりする必要はありません。FindMatches は、次のようなさまざまな問題で役立ちます。

- 一致する顧客: 多くの顧客フィールドがデータベース間で正確に一致しない場合でも (名前のスペルや住所の違い、データの欠落や不正確ななど)、異なる顧客データベース間で顧客レコードをリンクします。
- 製品のマッチング: カタログ内の製品を、競合他社のカタログに対する製品カタログなど、他の製品ソースと照合します。エントリの構造は異なります。

- 不正検出の向上: 重複した顧客アカウントを特定し、新しく作成したアカウントが、以前に知られている不正ユーザーと一致する (またはその可能性のある) タイミングを判断します。
- マッチングに関するその他の問題: マッチアドレス、映画、パーツリストなど。一般的に、人間がデータベース行を見て、それらが一致していると判断した場合、FindMatches 変換が役に立つ可能性は非常にあります。

これらの変換は、ジョブの作成時に作成できます。作成する変換は、ラベル付けするソースデータセットからのストアスキーマとサンプルデータに基づいています (このプロセスを変換の「トレーニング」と呼びます)。ラベルを付けるレコードは、ソースデータセットに存在する必要があります。このプロセスでは、ラベル付けしたファイルを生成し、変換が学習する方法でアップロードし直します。変換を教えたら、Spark ベースの AWS Glue ジョブ (PySpark または Scala Spark) から呼び出し、互換性のあるソースデータストアを持つ他のスクリプトで使用できます。

作成した変換は AWS Glue に保存されます。AWS Glue コンソールで、作成した変換を管理できます。[データ統合とETL] のナビゲーションペインで [データ分類ツール] > [レコードのマッチング] の順に移動して、機械学習変換を編集し、トレーニングを継続できます。コンソールで変換を管理する方法の詳細については、「[AWS Glue コンソールでの機械学習変換の使用](#)」を参照してください。

Note

AWS Glue version 2.0 FindMatches ジョブは、transform がデータを処理している間、aws-glue-temp-*<accountID>*-*<region>* 一時ファイルを保存するために Amazon S3 バケットを使用します。このデータは、手動または Amazon S3 ライフサイクルルールを設定して、実行の完了後に削除する事ができます。

機械学習変換のタイプ

機械学習変換を作成して、データを最適化できます。これらの変換は、ETL スクリプトから呼び出すことができます。データは、変換から変換へと DynamicFrame というデータ構造で渡されます。これは、Apache Spark SQL DataFrame を拡張したものです。DynamicFrame にはデータが含まれており、データを処理するためにそのスキーマを参照します。

次のタイプの機械学習変換を利用できます。

Find matches (一致の検索)

ソースデータの重複したレコードを見つけます。この機械学習変換をトレーニングするには、サンプルデータセットにラベルを付けて、どの行が一致するかを示します。機械学習変換は、ラ

ベル付けしたサンプルデータを使用してトレーニングするほど、どの行が一致するかを学習します。変換の設定方法に応じて、出力は次のいずれかになります。

- 入力テーブルのコピーと、一致するレコードのセットを示す値が入力された `match_id` 列。 `match_id` 列は任意の識別子です。同じ `match_id` を持つレコードは、互いに一致するレコードとして識別されています。 `match_id` が異なるレコードが一致しません。
- 重複する行が削除された入力テーブルのコピー。複数の重複が見つかった場合、最も低いプライマリーキーを持つレコードが保持されます。

インクリメンタルマッチを検索する

Find matches トランスフォームは、既存のフレームとインクリメンタルフレーム間で一致を検索し、一致グループごとに一意の ID を含む列を出力として返すように構成することもできます。

詳細については、「[インクリメンタルマッチを検索する。](#)」を参照してください。

FindMatches 変換の使用

FindMatches 変換を使用してソースデータの重複したレコードを見つけることができます。変換をトレーニングするために役立つラベリングファイルが生成または提供されます。

Note

現在、カスタム暗号化キーを使用する FindMatches 変換は、次のリージョンではサポートされていません。

- アジアパシフィック (大阪) - ap-northeast-3

FindMatches 変換の使用を開始する場合は、以下の手順に従ってください。より高度で詳細な例については、AWS ビッグデータブログの「[Harmonize data using AWS Glue and AWS Lake Formation FindMatches ML to build a customer 360 view](#)」を参照してください。

Find Matches Transform (一致の検索変換) の使用開始

FindMatches 変換の使用を開始するには、次の手順に従います。

1. 最適化するソースデータ用のテーブルを AWS Glue Data Catalog に作成します。クローラーの作成方法については、「[Working with Crawlers on the AWS Glue Console](#)」を参照してください。

ソースデータがカンマ区切り値 (CSV) ファイルなどのテキストベースのファイルである場合は、以下の点を考慮してください。

- 入力レコード CSV ファイルとラベリングファイルを別個のフォルダに保持します。そうしないと、AWS Glue クローラーはこれらのファイルを同じテーブルの複数の部分と見なして、データカタログに誤ったテーブルを作成する場合があります。
- CSV ファイルに ASCII 文字のみが含まれている場合を除き、CSV ファイルには必ず BOM (バイトオーダーマーク) なしの UTF-8 エンコードを使用します。Microsoft Excel は、UTF-8 CSV ファイルの先頭に BOM を追加することがよくあります。これを削除するには、CSV ファイルをテキストエディタで開き、ファイルを BOM なしの UTF-8 として保存し直します。

2. AWS Glue コンソールでジョブを作成し、[Find matches (一致の検索)] 変換タイプを選択します。

 Important

ジョブ用に選択するデータソーステーブルの列数は 100 を超えることはできません。

3. [Generate labeling file] (ラベルファイルを生成) を選択して、ラベル付けファイルを生成するように AWS Glue に伝えます。AWS Glue は、各 `labeling_set_id` に対して類似のレコードをグループ化する最初のパスを実行し、グループ化を確認できるようにします。label 列で一致のラベルを付けます。

- ラベル付けファイルが既にある場合 (つまり、一致する行を示すサンプルレコードがある場合) は、そのファイルを Amazon Simple Storage Service (Amazon S3) にアップロードします。ラベリングファイルの形式の詳細については、「[ラベリングファイルの形式](#)」を参照してください。ステップ 4 に進みます。

4. ラベル付けファイルをダウンロードし、「[ラベリング](#)」セクションの説明に従ってファイルにラベルを付けます。

5. 修正したラベル付きファイルをアップロードします。AWS Glue は一致の検索方法について変換をトレーニングするためのタスクを実行します。

[Machine learning transforms (機会学習変換)] リストのページで、[History (履歴)] タブを選択します。このページは、AWS Glue で以下のタスクを実行するタイミングを指定します。

- ラベルのインポート
- ラベルのエクスポート
- ラベルの生成
- 品質の推定

- より適切な変換を作成するには、ラベル付きファイルのダウンロード、ラベル付け、アップロードを繰り返すことができます。初期の実行では、レコードのミスマッチがより多く発生しがちです。ただし、ラベリングファイルを検証して継続的にトレーニングすることで、AWS Glue は学習します。
- 変換を評価して微調整するには、一致の検索のパフォーマンスと結果を評価します。詳細については、「[AWS Glue での機械学変換の調整](#)」を参照してください。

ラベリング

FindMatches がラベル付けファイルを生成すると、ソーステーブルからレコードが選択されます。これまでのトレーニングに基づき、FindMatches は最も学習価値が高いレコードを特定します。

ラベル付けでは、ラベリングファイル (Microsoft Excel などのスプレッドシートの使用を推奨) を編集し、一致するレコードと一致しないレコードを識別する label 列に識別子またはラベルを追加します。ソースデータの一致について一貫した明確な定義を持つことが重要です。FindMatches は、ユーザーが一致 (または不一致) として指定したレコードから学習します。また、ユーザーの判断を使用して重複するレコードを見つける方法を学習します。

ラベリングファイルが FindMatches で生成されると、約 100 個のレコードが生成されます。これらの 100 個のレコードは通常 10 個のラベリングセットに分割され、各ラベリングセットは FindMatches により生成される一意の labeling_set_id により識別されます。各ラベリングセットは、他のラベリングセットとは独立した個別のラベリングタスクとして表示する必要があります。タスクは、各ラベリングセット内の一致レコードと非一致レコードを識別することです。

スプレッドシートでラベル付けファイルを編集するためのヒント

スプレッドシートアプリケーションでラベリングファイルを編集する場合は、以下の点を考慮してください。

- ファイルを開いたときに、列フィールドが完全に展開されていない場合があります。必要に応じて labeling_set_id 列と label 列を展開し、これらのセルの内容を表示してください。
- プライマリキー列が数値 (long データ型など) である場合、スプレッドシートはこれを数値として解釈し、値を変更する場合があります。このキー値は、テキストとして扱う必要があります。この問題を修正するには、プライマリキー列のすべてのセルをテキストデータとしてフォーマットします。

ラベリングファイルの形式

FindMatches 変換をトレーニングするために AWS Glue により生成されたラベリングファイルは、次の形式を使用します。AWS Glue 用に独自のファイルを作成する場合は、次の形式にも従う必要があります。

- これは、カンマ区切り値 (CSV) ファイルです。
- UTF-8 でエンコードする必要があります。Microsoft Windows を使用してファイルを編集すると、ファイルは cp1252 でエンコードされる場合があります。
- このファイルを AWS Glue に渡すには、ファイルを Amazon S3 の場所に配置する必要があります。
- 各ラベル付けタスクに使用する行は多すぎないようにします。タスクごとに許容される行数は 2~30 行ですが、10~20 行をお勧めします。50 行を超えるタスクは推奨されません。粗悪な結果やシステム障害を引き起こす可能性があります。
- 「一致」または「不一致」とラベリングされたレコードのペアで構成される、ラベリングされたデータがすでにある場合、これは問題ありません。これらのラベリングされたペアは、サイズ 2 のラベリングセットとして表すことができます。この場合、両方のレコードに一致する場合は文字「A」とラベリングし、一致しない場合は「B」とラベリングします。

Note

ラベリングファイルには追加の列が含まれているため、ソースデータが含まれているファイルとはスキーマが異なります。ラベリングファイルは、他のすべての変換入力 CSV ファイルとは異なるフォルダに配置し、データカタログのテーブルの作成時に AWS Glue クローラーによって同種のファイルと見なされないようにします。そうしないと、AWS Glue クローラーで作成されたテーブルではデータが正しく表現されない場合があります。

- 最初の 2 つの列 (labeling_set_id、label) は AWS Glue の必須列です。残りの列は、処理するデータのスキーマと一致する必要があります。
- labeling_set_id ごとに、同じラベルを使用してすべての一致するレコードを特定します。ラベルは label 列に配置されている一意の文字列です。シンプルな文字 (A、B、C など) を含むラベルを使用することをお勧めします。ラベルは、大文字と小文字が区別され、label 列に入力されます。
- 同じ labeling_set_id と同じラベルを含む行は、一致とみなされ、ラベリングされます。
- 同じ labeling_set_id と異なるラベルを含む行は、不一致とみなされ、ラベリングされます。

- 異なる `labeling_set_id` を含む行は、一致に関するいかなる情報も伝達しないとみなされま

以下に示しているのは、データのラベル付けの例です。

labeling_set_id	ラベル	first_name	last_name	誕生日
ABC123	A	ジョン	Doe	04/01/1980
ABC123	B	Jane	Smith	04/03/1980
ABC123	A	Johnny	Doe	04/01/1980
ABC123	A	Jon	Doe	04/01/1980
DEF345	A	Richard	Jones	12/11/1992
DEF345	A	Rich	Jones	11/12/1992
DEF345	B	Sarah	Jones	12/11/1992
DEF345	C	Richie	Jones Jr.	05/06/2017
DEF345	B	Sarah	Jones-Walker	12/11/1992
GHI678	A	Robert	Miller	1/3/1999
GHI678	A	Bob	Miller	1/3/1999
XYZABC	A	William	Robinson	2/5/2001
XYZABC	B	Andrew	Robinson	2/5/1971

- 上記の例では、John/Johnny/Jon Doe を一致と認識し、これらのレコードは Jane Smith と一致しないことをシステムにトレーニングします。これとは別に、Richard と Rich Jones は同一人物であるが、これらのレコードは Sarah Jones/Jones-Walker および Jones Jr. とは一致しないとシステムをトレーニングします。
- ご覧のとおり、ラベルの範囲は `labeling_set_id` に制限されています。したがって、ラベルは `labeling_set_id` 境界を超えません。たとえば、`labeling_set_id 1` のラベル「A」は、`labeling_set_id 2` のラベル「A」とは何の関係もありません。

- ラベリングセット内に一致するものがほかにないレコードには、一意のラベルを割り当てます。たとえば、Jane Smith はラベリングセット ABC123 内のどのレコードとも一致しないため、そのラベリングセット内で唯一の、ラベル B を持つレコードです。
- ラベリングセット「GHI678」は、ラベリングセットが、同じラベルが付与された 2 つのレコードのみで構成されることを表示し、一致していることを示します。同様に、「XYZABC」は、異なるラベルが付与された 2 つのレコードを表示し、一致しないことを示します。
- ラベリングセットには一致が含まれていないこと（ラベリングセット内のレコードそれぞれに異なるラベルを付ける）も、ラベリングセットがすべて「同じ」こと（すべてに同じラベルを付けた）もあります。これは、ラベリングセットに、基準によって「同じ」レコードと「同じでない」レコードの例がまとめて含まれている限り、問題ありません。

Important

AWS Glue に渡す IAM ロールに、ラベリングファイルが格納されている Amazon S3 バケットへのアクセス権があることを確認してください。慣例として、AWS Glue ポリシーでは、名前に「aws-glue-」というプレフィックスが付いている Amazon S3 バケットまたはフォルダにアクセス許可を付与します。ラベリングファイルが別の場所にある場合、IAM ロールで、その場所に対してアクセス許可を追加します。

AWS Glue での機械学変換の調整

AWS Glue で機械学習変換を調整することで、目的に応じてデータクレンジングジョブの結果を改善できます。変換を改善するには、変換をトレーニングします。そのためには、ラベリングセットを生成し、ラベルを追加します。これらのステップを、目的の結果が得られるまで何回も繰り返します。いくつかの機械学習パラメータを変更することで調整することもできます。

機械学習変換の詳細については、「[AWS Lake Formation FindMatches によるレコードのマッチング](#)」を参照してください。

トピック

- [機械学習の測定](#)
- [適合率と再現率のトレードオフ](#)
- [精度とコストのトレードオフ](#)
- [一致信頼度スコアを使用して一致の品質を見積もります。](#)
- [「一致の検索」変換のトレーニング](#)

機械学習の測定

機械学習変換の調整に使用する測定を理解するには、以下の用語を理解しておく必要があります。

真陽性 (TP)

変換が正しく検出したデータ内の一致。ヒットとも呼ばれます。

真陰性 (TN)

変換が正しく拒否したデータ内の不一致。

偽陽性 (FP)

変換が誤って一致として分類したデータ内の不一致。誤警告とも呼ばれます。

偽陰性 (FN)

変換が検出しなかったデータ内の一致。ミスとも呼ばれます。

機械学習で使用される用語の詳細については、Wikipedia の「[混同行列](#)」を参照してください。

機械学習変換を調整するには、変換の [詳細プロパティ] で以下の測定値を変更できます。

- 適合率は、変換で要請と判定されたレコードの総数 (真陽性と偽陽性) のうち、どの程度真陽性を検出するかの尺度です。詳細については、Wikipedia の「[適合率と再現率](#)」を参照してください。
- 再現率は、変換がソースデータのレコード総数から真陽性を見つける割合を測定します。詳細については、Wikipedia の「[適合率と再現率](#)」を参照してください。
- 精度は、変換が真陽性と真陰性を見つける割合を測定します。精度を高めるには、マシンリソースとコストを増やす必要があります。ただし、再現率も高くなります。詳細については、Wikipedia の「[適合率と再現率](#)」を参照してください。
- コストは、変換を実行するために消費されるコンピューティングリソースの数およびそれに伴う金額を測定します。

適合率と再現率のトレードオフ

各 FindMatches 変換には、precision-recall パラメータが含まれています。このパラメータを使用して、次のいずれかを指定します。

- 変換の結果で、2 つのレコードが誤って一致すると報告され、実際には一致していない場合は、適合率を高めます。

- 変換の結果で、2つのレコードが実際には一致しているにもかかわらず、一致として検出されない場合は、再現率を高めます。

このトレードオフを行うには、AWS Glue コンソールまたは AWS Glue 機械学習 API オペレーションを使用できます。

適合率を優先する場合

FindMatches で実際には一致しないレコードのペアが一致と見なされるリスクに対処するには、適合率を優先します。適合率を優先するには、適合率-再現率トレードオフとしてより高い値を選択します。値を高くするほど、FindMatches 変換でレコードのペアを一致と見なす根拠がさらに必要となります。変換は調整されて、レコードのペアが一致しないと判断する傾向が強化されます。

たとえば、FindMatches を使用してビデオカタログ内の重複する商品を検出する際に、適合率-再現率のより高い値を変換に指定したとします。変換で「スターウォーズ: 新たなる希望」と「スターウォーズ: 帝国の逆襲」が誤って同じものとして検出された場合、「新たなる希望」を欲しいユーザーに「帝国の逆襲」が提供される可能性があります。これは、カスタマーエクスペリエンスを悪化させることとなります。

一方、変換で「スターウォーズ: 新たなる希望」と「スターウォーズ エピソード 4/新たなる希望」が同じ商品として検出されない場合、利用者は最初は混乱するとしても、最終的には同じものとして認識する可能性があります。これは誤検出ではあっても、前の例ほど深刻ではありません。

再現率を優先する場合

FindMatches 変換の結果で、実際には一致するレコードのペアが一致として検出されないリスクに対処するには、再現率を優先します。再現率を優先するには、適合率-再現率トレードオフとしてより低い値を選択します。値を低くするほど、FindMatches 変換でレコードのペアを一致と見なす根拠が少なくて済みます。変換は調整されて、レコードのペアが一致すると判断する傾向が強化されます。

たとえば、これはセキュリティを重視する組織で優先される場合があります。利用者を詐欺行為者のリストと照合する場合、利用者が詐欺行為者であるかどうかを判断することが重要です。FindMatches を使用して利用者リストと詐欺行為者リストを照合するとします。FindMatches において2つのリスト間で一致が検出されるたびに、該当する利用者が実際に詐欺行為者であるかどうかを人間の監査担当者が確認します。このような組織では、適合率よりも再現率を優先できます。つまり、利用者が実際に詐欺行為者リストに該当することを見逃すよりは、利用者が詐欺行為者ではないことを監査担当者が手動で確認して除外することを優先します。

適合率と再現率の両方を優先する方法

適合率と再現率の両方を改善する最適な方法は、ラベル付けするデータを増やすことです。より多くのデータをラベル付けすると、FindMatches 変換全体の精度が向上するため、適合率と再現率の両方が向上します。ただし、最も正確な変換であっても、適合率の優先、再現率の優先、または中間値の選択を試す必要があるグレー領域が常に存在します。

精度とコストのトレードオフ

各 FindMatches 変換には、accuracy-cost パラメータが含まれています。このパラメータを使用して、次のいずれかを指定できます。

- 変換で 2 つのレコードの一致が正確に報告されることをより重視する場合は、精度を強調します。
- 変換を実行するコストまたはスピードをより重視する場合は、より低いコストを強調します。

このトレードオフを行うには、AWS Glue コンソールまたは AWS Glue 機械学習 API オペレーションを使用できます。

精度を優先する場合

find matches の結果に一致が含まれないリスクに対処するには、精度を優先します。精度を優先するには、精度-コストトレードオフのより高い値を選択します。値が高いほど、FindMatches 変換でレコードを正しく一致させるために、より詳細な検索を行う時間を増やす必要があります。このパラメータは、一致しないレコードのペアを誤って一致と判断する可能性を減らすものではありません。変換は調整されて、一致の検索に費やす時間が増加されます。

コストを優先する場合

一致の検索数よりも find matches 変換を実行するコストをより重視する場合は、コストを優先します。コストを優先するには、精度-コストトレードオフのより低い値を選択します。値を低くするほど、FindMatches 変換で実行する必要があるリソース数が減ります。一致の検索数を減らすように、変換が調整されます。より低いコストを優先したときの結果が許容できるものであれば、この設定を使用します。

精度とより低いコストの両方を優先する方法

レコードが一致するかどうかを判断するためにより多くのレコードのペアを調査するには、より多くのマシン時間が必要になります。品質を下げずにコストを削減するには、以下を実行できます。

- 一致の対象としないレコードをデータソースから除外します。

- 一致/不一致の判断に役立たないことが確実な列をデータソースから除外します。これを決める適切な方法としては、一連のレコードが「同じ」であるかどうかの判断に影響しないと思われる列を除外します。

一致信頼度スコアを使用して一致の品質を見積もります。

一致信頼度スコアは、機械学習モデルの信頼度が高い、不確か、またはありそうにない一致レコードを区別するために、findMatches によって検出された一致の品質を推定します。一致信頼度スコアは 0 から 1 の間で、スコアが高いほど類似度が高くなります。一致信頼度スコアを調べると、システムが非常に信頼できる一致のクラスター (マージすることを決定する可能性があります)、システムが不確実であるクラスター (人間がレビューしたと決定する可能性がある)、および一致のクラスターを区別できます。システムは可能性が低いと見なします (拒否することを決定する場合があります)。

一致信頼度スコアが高いが一致がないと判断した場合、またはスコアが低い但实际上に一致しているかどうかを判断する状況では、トレーニングデータを調整できます。

信頼度スコアは、すべての FindMatches 決定を確認することが不可能な、大規模な産業用データセットがある場合に特に役立ちます。

一致信頼度スコアは AWS Glue バージョン 2.0 以降。

マッチ信頼度スコアの生成

ブール値を computeMatchConfidenceScores True に設定するか、FindMatches または FindIncrementalMatches API を呼び出すときに一致信頼スコアを生成できます。

AWS Glue 新規に column match_confidence_score を出力に追加します。

スコアリングの例に一致

たとえば、次のマッチングレコードを考えてみましょう。

スコア ≥ 0.9

一致したレコードのサマリー:

```
primary_id | match_id | match_confidence_score
3281355037663    85899345947    0.9823658302132061
```

1546188247619 85899345947 0.9823658302132061

詳細:

raw_id	phone source	website	poi_id	display_position	primary_name locale_name	street1 street2 street3
city state country postal_code street_in_one_line	primary_id	match_id match_confidence_score				
aeJq8SD0iCbIqHFPPL1jg +43262681168 yelp http://www.commerzbank.at yelp::aeJq8SD0iCbIqHFPPL1jg geo:47.711590000,16.344020000 Commerzbank Mattersburg en_US Hauptstr. 59 null null Forchtenstein 1 AT 7212	Hauptstr. 59 1546188247619 85899345947 0.9823658302132061	uHhQK6v2j5lZ4N8lXm-QQ +43268747266 yelp http://www.commerzbank.at yelp::uHhQK6v2j5lZ4N8lXm-QQ geo:47.787420000,16.455440000 Commerzbank Mattersburg en_US Hauptstr. 9 null null Hirm 1 AT 7024	Hauptstr. 9 3281355037663 85899345947 0.9823658302132061			

この例から、2つのレコードが非常によく似ていて、`display_position`、`primary_name`、および `street name` を共有していることがわかります。

スコア ≥ 0.8 、得点 < 0.9

一致したレコードのサマリー:

primary_id	match_id	match_confidence_score
309237680432	85899345928	0.8309852373674638
3590592666790	85899345928	0.8309852373674638
343597390617	85899345928	0.8309852373674638
249108124906	85899345928	0.8309852373674638
463856477937	85899345928	0.8309852373674638

詳細:

raw_id	phone source	website	poi_id	display_position	primary_name locale_name	street1 street2 street3	city state country postal_code street_in_one_line	
primary_id	match_id match_confidence_score							
NVIMVA35Tm41mnaokyvr_w null yelp null yelp::NVIMVA35Tm41mnaokyvr_w geo:50.541800000,7.102920000 Eiscafe Dolomiten en_US Ahrhutstr. 49 null null Bad Neuenahr-Ahrweiler RP DE 53474 Ahrhutstr. 49	343597390617 85899345928 0.8309852373674638	S3HnQeSvjkcIsh9XQFpeQ +49557465221 yelp null yelp::S3HnQeSvjkcIsh9XQFpeQ geo:51.447337266,9.414379060 Eiscafe Dolomiten en_US Markt 5 null null Grebenstein HE DE 34393 Markt 5	153HnQeSvjkcIsh9XQFpeQ +49557465221 yelp null yelp::S3HnQeSvjkcIsh9XQFpeQ geo:51.447337266,9.414379060 Eiscafe Dolomiten en_US Markt 5 null null Grebenstein HE DE 34393 Markt 5	463856477937 85899345928 0.8309852373674638	o6f-p0xtJmI9PIKpsjx5CQ +493691744935 yelp null yelp::o6f-p0xtJmI9PIKpsjx5CQ geo:50.976200000,10.324000000 Eiscafe Dolomiten en_US Alexanderstr. 105 null null Eisenach TH DE 99817 Alexanderstr. 105	309237680432 85899345928 0.8309852373674638	D10Q21YDNXonoG52royfjw +4926445735 yelp null yelp::D10Q21YDNXonoG52royfjw geo:50.565900000,7.280050000 Eiscafe Dolomiten en_US Rheinstr. 15 null null Linz RP DE 53545 Rheinstr. 15	15 3590592666790 85899345928 0.8309852373674638

この例から、これらのレコードが同じ `primary_name` および `country` を共有していることがわかります。

スコア ≥ 0.6 、得点 < 0.7

一致したレコードのサマリー:

primary_id	match_id	match_confidence_score
2164663519676	85899345930	0.6971099896480333
317827595278	85899345930	0.6971099896480333

```

472446424341 85899345930 0.6971099896480333
3118146262932 85899345930 0.6971099896480333
214748380804 85899345930 0.6971099896480333

```

詳細:

primary_id	match_id	match_confidence_score	raw_id	phone	source	website	poi_id	display_position	primary_name	locale_name	street1	street2	street3	city	state	country	postal_code	street_in_one_line	
[10T_R8tk4ngTFXhpy88w]	[+33490963451]	[yelp]	[null]	[yelp::1QT_R8tk4ngTFXhpy88w]	[geo:43.675559000,4.626792000]	[Le Vésuve]	[en_US]	[15 Rue de la Rotonde]	[null]	[null]	[Arles]	[13]	[FR]	[13200]	[15 Rue de la Rotonde]				
[17827595278]	[85899345930]	[0.6971099896480333]	[b8cCaxbvEcug27QmQYjQ]	[null]	[yelp]	[null]	[yelp::b8cCaxbvEcug27QmQYjQ]	[geo:50.631700000,3.067300000]	[Le Vésuve]	[en_US]	[30 ave du President Kennedy]	[null]	[null]	[Lille]	[59]	[FR]	[59800]	[30 ave du President]	
[Kennedy]	[3118146262932]	[85899345930]	[0.6971099896480333]	[dJOC4FZmXSlwEnFB6vJ5g]	[+33442750804]	[yelp]	[null]	[yelp::dJOC4FZmXSlwEnFB6vJ5g]	[geo:43.427710000,5.236950000]	[le Vesuve]	[en_US]	[24 ave Bruxelles]	[null]	[null]	[Vitrolles]	[13]	[FR]	[13127]	[24 ave]
[Bruxelles]	[214748380804]	[85899345930]	[0.6971099896480333]	[u859qGa561Cj1t4wypkg]	[+33297251001]	[yelp]	[null]	[yelp::u859qGa561Cj1t4wypkg]	[geo:48.071493200,-2.963742000]	[Le Vésuve]	[en_US]	[49 Rue Gén de Gaulle]	[null]	[null]	[Pontivy]	[56]	[FR]	[56300]	[49 Rue Gén de Gaulle]
[472446424341]	[85899345930]	[0.6971099896480333]	[3wH0MEhra3DUUgF_YcoTA]	[+33164069200]	[yelp]	[null]	[yelp::3wH0MEhra3DUUgF_YcoTA]	[geo:48.618984000,2.888184000]	[Le Vesuve]	[en_US]	[59 Avenue Charles de Gaulle]	[null]	[null]	[Mormant]	[77]	[FR]	[77720]	[59 Avenue Charles de Gaulle]	
[214748380804]	[85899345930]	[0.6971099896480333]																	

この例から、これらのレコードが同じものの `primary_name` を共有しているだけであることがわかります。

詳細については、以下を参照してください。

- [ステップ 5: 機械学習変換でジョブを追加して実行する](#)
- PySpark: [FindMatches クラス](#)
- PySpark: [FindIncrementalMatches クラス](#)
- Scala: [FindMatches クラス](#)
- Scala: [FindIncrementalMatches クラス](#)

「一致の検索」変換のトレーニング

一致と不一致を見分けられるように各 `FindMatches` 変換をトレーニングする必要があります。変換をトレーニングするには、ファイルにラベルを追加して選択結果を AWS Glue にアップロードします。

このラベル付けは、AWS Glue コンソールまたは AWS Glue 機械学習 API オペレーションを使用して調整できます。

ラベルを追加する回数は? 必要なラベルの数は?

これらの質問に対する回答は、主にユーザーが決定します。`FindMatches` で必要な精度レベルを達成できるかどうか、さらにラベル付けの作業を追加するメリットがあるかどうかは、ユーザーが評価する必要があります。これを決定する最善の方法は、AWS Glue コンソールで `[Estimate quality]` (品質の推定) を選択して、「適合率」、「再現率」、「適合率再現率曲線の下での面積」のメトリクスを調べることです。さらに多くのタスクセットにラベル付けをしたら、これらのメトリクスを再実行し

て結果が改善されたかどうかを確認します。いくつかのタスクセットをラベル付けをした後で、注目しているメトリクスに改善が見られなかった場合、変換の品質はそれ以上向上しない状態に達している可能性があります。

真陽性と真陰性の両方のラベルが必要な理由は？

FindMatches 変換は、ユーザーが何を一致と見なしているかを学習するために真陽性と真陰性の両方のサンプルを必要とします。FindMatches で生成されるトレーニングデータ ([I do not have labels (ラベルがありません)] オプションなどを使用した場合) にラベル付けを行う場合、FindMatches はユーザーに代わって「ラベルセット ID」を生成しようとします。各タスク内で、一部のレコードに同じ「ラベル」を付け、他のレコードに異なる「ラベル」を付けます。つまり、通常、タスクがすべて同じものであったり、すべて異なるものであったりすることはありません (ただし、特定のタスク内のレコードがすべて「同じ」ものであったり、すべて「異なる」ものであったりすることは問題ありません)。

[Upload labels from S3] (S3 からラベルをアップロードする) オプションを使用して FindMatches 変換をトレーニングする場合は、一致レコードと不一致レコードの両方を含めるようにしてください。1つのタイプのみを使用することもできます。これらのラベルはより正確な FindMatches 変換の構築に役立ちますが、[Generate labeling file (ラベリングファイルの生成)] オプションを使用して生成する一部のレコードには依然としてラベル付けを行う必要があります。

どうしたら、ユーザーがトレーニングしたとおりのマッチングを変換に実行させることができるか？

FindMatches 変換は、ユーザーが指定したラベルから学習するため、指定したラベルを優先しないレコードのペアが生成される場合があります。ユーザーのラベルを優先するように FindMatches 変換を強制するには、[FindMatchesParameter] の [EnforceProvidedLabels] を選択します。

ML 変換が真の一致ではない項目を一致として識別した場合、どのように対策できますか？

以下の対策を実行できます。

- precisionRecallTradeoff の値をより高くします。これにより、検出される一致の数が減りますが、十分に高い値に達すると、大きなクラスターの分割も行われます。
- 不正確な結果に対応する出力行を取得し、これらの行をラベリングセットとして再フォーマットします (match_id 列を削除し、labeling_set_id 列および label 列を追加します)。必要に応じて、複数のラベリングセットに分割 (再分割) し、ラベラーが各ラベリングセットに留意しながらラベルの割り当てができるようにします。次に、一致するセットに正しくラベル付けを行い、ラベルファイルをアップロードして既存のラベルに付加します。これにより、パターンを理解するためにどのようなデータが必要であるかについてトランスフォーマーが十分にトレーニングされます。

- (高度) 最後に、そのデータを確認し、システムが気付いていないパターンを検出できるかどうかを判断します。標準の AWS Glue 関数を使用してそのデータを事前処理し、データを正規化します。データの種類が十分に異なることが判明している場合は、データを種類別に分離し、アルゴリズムに学習させるデータを明確にします。または、複数の列のデータが関連していることが判明している場合は、これらの列を結合します。

AWS Glue コンソールでの機械学習変換の使用

AWS Glue を使用して、データのクレンジングに使用できるカスタムの機械学習トランスフォームを作成できます。これらの変換は、AWS Glue コンソールでのジョブの作成時に使用できます。

機械学習変換の作成方法の詳細については、「[AWS Lake Formation FindMatches によるレコードのマッチング](#)」を参照してください。

トピック

- [変換のプロパティ](#)
- [機械学習変換の追加と編集](#)
- [変換の詳細の確認](#)
- [ラベルを使って変換を教える](#)

変換のプロパティ

[既存の機械学習トランスフォームを表示するには、にログインし AWS Management Console、https://console.aws.amazon.com/glue/ AWS Glue のコンソールを開きます。](#) [データ統合と ETL] のナビゲーションペインで、[データ分類ツール] > [レコードのマッチング] の順に選択します。

各変換のプロパティ:

変換名

変換の作成時に付けた一意の変換名。

ID

変換の一意の識別子。

[Label count] (ラベル数)

変換のトレーニングに役立てるために指定されたラベリングファイル内のラベルの数。

ステータス

変換が [Ready] (使用可能) であるか、[Needs training] (トレーニングが必要) であるかを示します。ジョブで機械学習変換を正常に実行するには、ジョブのステータスが [Ready] (使用可能) である必要があります。

作成

変換の作成日。

変更済み

変換が最後に更新された日付。

説明

変換の説明 (説明が入力された場合)。

AWS Glue バージョン

使用する AWS Glue のバージョン。

実行 ID

変換の作成時に付けた一意の変換名。

タスクタイプ

機械学習変換のタイプ (一致するレコードの検索など)。

ステータス

タスク実行のステータスを示します。次のようなステータスがあります。

- スタート
- 実行中
- 停止中
- 停止
- 成功
- [失敗]
- タイムアウト

エラー

ステータスが [失敗] の場合、失敗の理由を説明するエラーメッセージが表示されます。

機械学習変換の追加と編集

AWS Glue コンソールで変換の表示、削除、セットアップ、トレーニング、および調整を行うことができます。リストの変換の横にあるチェックボックスをオンにし、[アクション] を選択して、実行するアクションを選択します。

新しい ML 変換の作成

新しい機械学習変換を追加するには、[変換の作成] を選択します。次に、[ジョブの追加] ウィザードの手順に従います。詳細については、「[AWS Lake Formation FindMatches によるレコードのマッチング](#)」を参照してください。

Step 1. 変換のプロパティを設定する。

1. 名前と説明 (オプション) を入力します。
2. オプションで、セキュリティ設定を設定します。[機械学習変換でのデータ暗号化の使用](#) を参照してください。
3. オプションで、タスク実行の設定を設定します。タスク実行の設定では、タスクの実行方法をカスタマイズできます。ワーカータイプ、ワーカーの数、タスクのタイムアウト (分単位)、再試行回数、AWS Glue バージョンを選択します。
4. 必要に応じて、タグを設定します。AWS タグはリソースに割り当てることができるラベルです。各タグは、キー、および値 (オプション) で構成されます。タグは、リソースを検索して絞り込んだり、AWS コストを追跡したりするために使用できます。

Step 2. テーブルとプライマリキーを選択する。

1. AWS Glue カタログデータベースおよびテーブルを選択します。
2. 選択したテーブルからプライマリキーを選択します。プライマリキー列には通常、データソース内のすべてのレコードの固有の識別子が含まれます。

ステップ 3. 調整オプションを選択する。

1. [再現率と適合率] では、再現率または適合率を優先するように変換を調整する調整値を選択します。デフォルトでは [バランス] が選択されていますが、再現率または適合率を優先するように選択するか、[カスタム] を選択して 0.0~1.0 (両端を含む) の値を入力できます。
2. [コスト削減と正確度] では、コスト削減または正確度を優先するように調整値を選択するか、[カスタム] を選択して 0.0~1.0 (両端を含む) の値を入力できます。

3. 使用するラベルと出力を一致させることで ML 変換を教える場合は、[一致のエンフォースメント] で、[出力をラベルに一致させる] を選択します。

ステップ 4 確認して作成します。

1. ステップ 1~3 のオプションを確認します。
2. 変更が必要な手順については、[編集] を選択します。[変換の作成] を選択して、変換の作成ウィザードを完了します。

機械学習変換でのデータ暗号化の使用

機械学習変換を AWS Glue に追加する際に、データソースまたはデータターゲットに関連付けられたセキュリティ設定をオプションで指定できます。データの格納に使用される Amazon S3 バケットがセキュリティ設定で暗号化されている場合は、変換の作成時に同じセキュリティ設定を指定します。

また、(SSE-KMS) によるサーバー側の暗号化 AWS KMS (SSE-KMS) を使用してモデルとラベルを暗号化し、権限のない人が検査できないようにすることもできます。このオプションを選択すると、AWS KMS key 名前で選択するように求めるプロンプトが表示されるか、[Enter a key ARN] を選択できます。KMS キーの ARN を入力することを選択した場合、KMS キー ARN を入力できる 2 番目のフィールドが表示されます。

Note

現在、次のリージョンではカスタム暗号化キーを使用する ML 変換がサポートされていません。

- アジアパシフィック (大阪) - ap-northeast-3

変換の詳細の確認

変換のプロパティを確認する

[変換プロパティ] ページには、変換の属性が含まれます。変換定義に関する詳細として、以下の内容が表示されます。

- [変換名] は変換の名前を示します。
- [タイプ] は変換のタイプを一覧表示します。

- [ステータス] は、変換がスクリプトまたはジョブで使用可能かどうかを示します。
- [Force output to match labels (出力をラベルに一致させる)] は、ユーザーから提供されたラベルに出力を一致させるかどうかを示します。
- [Spark version] (Spark バージョン) は、変換を追加したときに [Task run properties] (タスク実行のプロパティ) で選択した AWS Glue のバージョンに関連しています。ほとんどのお客様に AWS Glue 1.0 および Spark 2.4 が推奨されます。詳細については、「[AWS Glue Versions](#)」を参照してください。

[履歴]、[品質を推定する]、[タグ] タブ

変換の詳細には、変換の作成時に定義した情報が含まれます。変換の詳細を確認するには、機械学習変換のリストで変換を選択し、以下のタブの情報を確認します。

- 履歴
- 品質の推定
- タグ

履歴

[履歴] タブには、変換のタスク実行の履歴が表示されます。変換をトレーニングするには、複数のタイプのタスクを実行します。タスクごとに、実行メトリクスとして以下が含まれます。

- [実行 ID] は、このタスクの実行ごとに AWS Glue によって作成される識別子です。
- [Task type (タスクタイプ)] は、タスク実行のタイプを示します。
- [実行ステータス] は、各タスクの成功した実行を一覧表示します。最新の実行が一番上に表示されます。
- [Error (エラー)] には、実行が正常に行われなかった場合のエラーメッセージの詳細が表示されます。
- [開始時刻] は、タスクの開始日時 (現地時間) を示します。
- [終了時刻] には、タスクの終了日時 (現地時間) が表示されます。
- [ログ] は、このジョブ実行の stdout に書き込まれたログにリンクされています。

「ログ」リンクをクリックすると、Amazon CloudWatch ログに移動します。ここで、AWS Glue Data Catalog で作成されたテーブルの詳細と、発生したエラーを確認できます。CloudWatch ログの保持期間はコンソールで管理できます。デフォルトのログ保持期間は Never Expire で

す。保持期間を変更する方法の詳細については、Amazon CloudWatch Logs ユーザーガイドの「[CloudWatch ログのログデータ保持期間の変更](#)」を参照してください。

- [ラベルファイル] には、生成されたラベリングファイル用の Amazon S3 へのリンクが示されません。

品質の推定

[Estimate Quality (品質の推定)] タブは、変換の品質を測定するために使用するメトリクスを示します。推定値は、ラベル付きデータのサブセットを使用して、指定したラベルに変換一致予測を比較することによって計算されます。これらの推定値はおおよその値です。このタブから [Estimate quality (品質の推定)] タスク実行を呼び出すことができます。

Estimate quality タブには、次のプロパティを含む最後の Estimate quality の実行のメトリクスが表示されます。

- Precision-Recall 曲線の下にある領域は、変換の全体的な品質の上限を推定する単一の数値です。これは適合率-再現率パラメータ用に行った選択とは関係ありません。値が高いほど、適合率と再現率のトレードオフが適切であることを示します。
- 適合率は、変換が一致を推定して、その推定が正確である度合いを見積もります。
- 再現率の上限は、実際的一致件数に対して、変換が推定した一致件数の割合の見積もりです。
- [F1] は、変換の正確性 (1~0) を示します。1 は正確性が最適であることを表します。詳細については、Wikipedia の「[F1 スコア](#)」を参照してください。
- [Column importance] (列の重要度) テーブルには、各列の列名と重要度スコアが表示されます。列の重要度は、レコード内のどの列がマッチングを行うために最も使用されているかを識別することで、列がモデルにどの程度寄与しているかを理解するのに役立ちます。列の重要度を上げたり下げたりするために、このデータをラベルセットに追加したりラベルセットを変更したりするように促されることがあります。

[Importance] (重要度) 列には、各列の数値スコアが 1.0 以下の 10 進数で示されます。

品質の推定と真の品質の比較については、「[品質推定値と end-to-end \(真の\) 品質](#)」を参照してください。

変換の調整の詳細については、「[AWS Glue での機械学変換の調整](#)」を参照してください。

品質推定値と end-to-end (真の) 品質

AWS Glue は、変換の品質を推定するために、内部の機械学習モデルに対してレコードのペアをいくつか提示します。これらのレコードのペアは、一致のラベルを指定したレコードですが、モデルは初見のものです。これらの品質の推定は、機械学習モデルの品質の関数です (モデルは、ユーザーが変換を「トレーニング」するためにラベル付けするレコードの数から影響を受けます)。または真のリコール (によって自動的に計算されるわけではない ML transform) は end-to-end、ML transform 機械学習モデルに一致する可能性のあるさまざまなものを提案するフィルタリングメカニズムの影響も受けます。

このフィルタリング方法を調整するには、主に [低コスト - 正確性] の調整値を指定します。[正確性] を優先するようにこの調整値を近づけるほど、システムは一致する可能性のあるレコードのペアをより深く、広く検索します。機械学習モデルに送られるレコードのペアが増え、ML transform end-to-end ユーザーのリコールまたは真の想起値が推定リコール指標に近づきます。その結果、end-to-end マッチのコストと精度のトレードオフの変化によるマッチの質の変化は、通常、品質推定には反映されません。

タグ

タグはリソースに割り当てることができるラベルです。AWS 各タグは、キー、および値 (オプション) で構成されます。タグは、リソースを検索して絞り込んだり、AWS コストを追跡したりするために使用できます。

ラベルを使って変換を教える

ML 変換の詳細ページから [変換を教える] を選択すると、ラベル (例) を使用して ML 変換を教えることができます。例 (ラベルと呼ばれます) を提供して機械学習のアルゴリズムを教える場合、使用する既存のラベルを選択したり、ラベリングファイルを作成できます。

Teach the transform using labels

Labeling

Teach your machine learning algorithms by providing examples, called labels. For your transform, provide examples of matching and nonmatching records.

I do not have labels

I have labels

▶ How to label

Generate labeling file

AWS Glue extracts records from your source data and suggests potential matching records. The file will contain approximately 100 data samples for you to work with. You can download the file once it has been generated.

S3 path to store the generated label file

🔍 s3://bucket/prefix/object

View [🔗](#)

Browse S3

Generate labeling file

Download labeling file

Upload labels from S3

The completed labeling file must be in the correct format and in Amazon S3.

S3 path where the label file is stored

🔍 s3://bucket/prefix/object

View [🔗](#)

Browse S3

Existing labels

Append to my existing labels

Overwrite my existing labels

Upload labeling file from S3

- [ラベリング] — ラベルがある場合は、[ラベルがあります] を選択します。ラベルがない場合でも、ラベリングファイルを生成する次の手順に進むことができます。
- [ラベリングファイルを生成する] — AWS Glue により、ソースデータからレコードが抽出され、一致する可能性のあるレコードが提案されます。生成されたラベルファイルを保存する Amazon S3 バケットを選択します。[ラベリングファイルを生成する] を選択してプロセスを開始します。完了したら、[ラベリングファイルをダウンロードする] を選択します。ダウンロードしたファイルには、ラベルを入力できるラベル用の列があります。
- [Amazon S3 からラベルをアップロードする] — ラベルファイルが保存されている Amazon S3 バケットから、完成したラベリングファイルを選択します。その後、既存のラベルにラベルを追加するか、既存のラベルを上書きするかを選択します。[Amazon S3 からラベリングファイルをアップロードする] を選択します。

チュートリアル: AWS Glue で機械学習変換を作成する

このチュートリアルでは、AWS Glue を使用して機械学習 (ML) 変換を作成および管理するためのアクションについてガイドします。このチュートリアルを使用するには、AWS Glue コンソールを使用してクローラとジョブを追加し、スクリプトを編集する方法をよく理解している必要があります。Amazon Simple Storage Service (Amazon S3) で、ファイルを検索してダウンロードする方法についてもよく理解している必要があります。

次の例では、一致するレコードを見つけるための FindMatches 変換を作成します。また、一致する/一致しないレコードを識別する方法を、この変換にトレーニングして AWS Glue ジョブで使用します。AWS Glue ジョブは、match_id という別の列を追加して新しい Amazon S3 ファイルを作成します。

このチュートリアルで使用されているソースデータは、dblp_acm_records.csv という名前のファイルです。このファイルは、元の [DBLP ACM データセット](#) から入手できる学術刊行物 (DBLP および ACM) の修正バージョンです。dblp_acm_records.csv ファイルは、BOM (バイトオーダーマーク) なしの UTF-8 形式でエンコードされたカンマ区切り値 (CSV) ファイルです。

2 番目のファイル dblp_acm_labels.csv はサンプルのラベリングファイルです。このファイルには、チュートリアルの一環として変換をトレーニングするために使用する、一致する/一致しないレコードの両方が含まれています。

トピック

- [ステップ 1: ソースデータをクローラする](#)
- [ステップ 2: 機械学習変換を追加する](#)
- [ステップ 3: 機械学習変換をトレーニングする](#)
- [ステップ 4: 機械学習変換の品質を推定する](#)
- [ステップ 5: 機械学習変換でジョブを追加して実行する](#)
- [ステップ 6: Amazon S3 の出力データを確認する](#)

ステップ 1: ソースデータをクローラする

まず、ソースである Amazon S3 CSV ファイルをクローラして、対応するメタデータテーブルを Data Catalog に作成します。

⚠ Important

CSV ファイルのみのテーブルを作成するようクローラに指示するために、CSV ソースデータを他のファイルとは異なる Amazon S3 フォルダに保存します。

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインで、[クローラ]、[クローラの追加] の順に選択します。
3. ウィザードに従って demo-crawl-dblp-acm というクローラを作成して実行し、出力をデータベース demo-db-dblp-acm に保存します。ウィザードの実行時に、データベース demo-db-dblp-acm がまだ存在していない場合は、これを作成します。現在の AWS リージョンでのサンプルデータへの Amazon S3 インクルードパスを選択します。例えば、us-east-1 の場合、ソースファイルへの Amazon S3 インクルードパスは s3://ml-transforms-public-datasets-us-east-1/dblp-acm/records/dblp_acm_records.csv です。

完了すると、クローラによってテーブル dblp_acm_records_csv が作成され、テーブルの列として、ID、タイトル、作成者、場所、年、およびソースが設定されます。

ステップ 2: 機械学習変換を追加する

次に、demo-crawl-dblp-acm というクローラで作成されたデータソーステーブルのスキーマに基づく機械学習変換を追加します。

1. AWS Glue コンソールにある [データ統合と ETL] のナビゲーションペインで、[データ分類ツール] > [レコードのマッチング] の順に選択し、[変換の追加] を選択します。ウィザードに従って、以下のプロパティを持つ Find matches 変換を作成します。
 - a. [変換の名前] に、「**demo-xform-dblp-acm**」と入力します。これは、ソースデータの一致を見つけるために使用する変換の名前です。
 - b. [IAM role] (IAM ロール) で、Amazon S3 ソースデータ、ラベリングファイル、および AWS Glue API オペレーションへのアクセス許可を持つ IAM ロールを選択します。詳細については、AWS Glue デベロッパーガイドの「[Create an IAM Role for AWS Glue](#)」を参照してください。
 - c. [Data source] (データソース) に、demo-db-dblp-acm というデータベースの dblp_acm_records_csv という名前のテーブルを選択します。

- d. [Primary key (プライマリキー)] で、テーブルのプライマリキー列である [ID] を選択します。
2. ウィザードで、[完了] を選択し、[ML 変換] リストに戻ります。

ステップ 3: 機械学習変換をトレーニングする

次に、チュートリアルサンプルのラベリングファイルを使用して、機械学習変換をトレーニングします。

機械学習変換は、そのステータスが [Ready for use (使用可能)] になるまで、ETL (抽出、変換、ロード) ジョブで使用できません。変換を使用可能にするには、一致するレコードと一致しないレコードのサンプルを提供し、レコードの一致/不一致を見分ける方法について変換をトレーニングする必要があります。変換をトレーニングするには、ラベルファイルを生成して、ラベルを追加し、ラベルファイルをアップロードします。このチュートリアルでは、dblp_acm_labels.csv というサンプルラベリングファイルを使用できます。ラベリングプロセスの詳細については、「[ラベリング](#)」を参照してください。

1. AWS Glue コンソールのナビゲーションペインで、[レコードのマッチング] を選択します。
2. demo-xform-dblp-acm 変換を選択し、次に [アクション]、[Teach (トレーニング)] の順に選択します。ウィザードに従って、Find matches 変換をトレーニングします。
3. 変換のプロパティページで、[I have labels (ラベルがあります)] を選択します。現在の AWS リージョンでサンプルラベリングファイルへの Amazon S3 パスを選択します。例えば、us-east-1 の場合、指定したラベリングファイルを Amazon S3 パス s3://ml-transforms-public-datasets-us-east-1/dblp-acm/labels/dblp_acm_labels.csv からアップロードします。この際に、既存のラベルを上書きするオプションを使用します。ラベリングファイルは、AWS Glue コンソールと同じリージョンの Amazon S3 に配置する必要があります。

ラベリングファイルをアップロードすると、データソースの処理方法について変換をトレーニングするために使用するラベルを追加または上書きするためのタスクが AWS Glue で開始されます。

4. ウィザードの最後のページで、[完了] を選択し、[ML 変換] リストに戻ります。

ステップ 4: 機械学習変換の品質を推定する

次に、機械学習変換の品質を推定できます。品質は、どれだけのラベル付けを実施したかに応じて異なります。品質の推定の詳細については、「[品質の推定](#)」を参照してください。

1. AWS Glue コンソールにある [データ統合と ETL] のナビゲーションペインで、[データ分類ツール] > [レコードのマッチング] の順に選択します。
2. demo-xform-dblp-acm 変換を選択し、[Estimate quality (品質の推定)] タブをクリックします。このタブには、変換の現在の品質の推定 (使用可能な場合) が表示されます。
3. [Estimate quality (品質の推定)] を選択して、変換の品質を推定するためのタスクを開始します。品質の推定の精度はソースデータのラベル付けに基づきます。
4. [履歴] タブに移動します。このペインには、変換のタスク実行が一覧表示されます。これには、品質の推定タスクも含まれます。実行の詳細を参照するには、[ログ] を選択します。実行が完了したときの実行ステータスが [成功] になっていることを確認します。

ステップ 5: 機械学習変換でジョブを追加して実行する

このステップでは、機械学習変換を使用して AWS Glue にジョブを追加して実行します。変換 demo-xform-dblp-acm は、[Ready for use (使用可能)] である場合に、ETL ジョブで使用できません。

1. AWS Glue コンソールのナビゲーションペインで、[ジョブ] を選択します。
2. [Add job (ジョブの追加)] を選択し、ウィザードの手順に従い、生成されたスクリプトを使用して ETL Spark ジョブを作成します。変換のプロパティ値として以下を選択します。
 - a. [名前] で、このチュートリアルのサンプルジョブである demo-etl-dblp-acm を選択します。
 - b. [IAM role] (IAM ロール) に、Amazon S3 ソースデータ、ラベリングファイル、および AWS Glue API オペレーションへのアクセス許可を持つ IAM ロールを選択します。詳細については、AWS Glue デベロッパーガイドの「[Create an IAM Role for AWS Glue](#)」を参照してください。
 - c. [ETL 言語] で、[Scala] を選択します。これは ETL スクリプトのプログラミング言語です。
 - d. [スクリプトファイル名] で、[demo-etl-dblp-acm] を選択します。これは Scala スクリプトのファイル名です (ジョブ名と同じです)。
 - e. [データソース] で、[dblp_acm_records_csv] を選択します。選択したデータソースは、機械学習変換のデータソーススキーマと一致する必要があります。
 - f. [Transform type (変換タイプ)] で、[Find matching records (一致するレコードの検索)] を選択し、機械学習変換を使用してジョブを作成します。
 - g. [Remove duplicate records (重複するレコードの削除)] をオフにします。重複するレコードを削除しない理由は、書き込まれる出力レコードに別の match_id フィールドが追加されるためです。

- h. [変換] で、ジョブによって使用される機械学習変換である `demo-xform-dblp-acm` を選択します。
- i. [データターゲットでテーブルを作成する] で、次のプロパティを使用してテーブルを作成することを選択します。
 - [Data store type] (データストアのタイプ) – **Amazon S3**
 - [Format] (形式) – **CSV**
 - [Compression type] (圧縮タイプ) – **None**
 - [Target path] (ターゲットパス) – ジョブの出力が書き込まれる Amazon S3 パス (現在のコンソールの AWS リージョン)
3. [ジョブを保存してスクリプトを編集する] を選択して、スクリプトエディタページを表示します。
4. スクリプトを編集し、ターゲットパスへのジョブ出力を単一のパーティションファイルに書き込ませるためのステートメントを追加します。このステートメントは、`FindMatches` 変換を実行するステートメントの直後に追加します。このステートメントは次のようになります。

```
val single_partition = findmatches1.repartition(1)
```

出力を `.writeDynamicFrame(single_partition)` として書き込むには、`.writeDynamicFrame(findmatches1)` ステートメントを変更する必要があります。

5. スクリプトを編集したら、[保存] を選択します。変更後のスクリプトは次のコードのようになりますが、使用環境向けにカスタマイズされます。

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.errors.CallSite
import com.amazonaws.services.glue.ml.FindMatches
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._

object GlueApp {
  def main(sysArgs: Array[String]) {
    val spark: SparkContext = new SparkContext()
    val glueContext: GlueContext = new GlueContext(spark)
    // @params: [JOB_NAME]
```

```

val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
Job.init(args("JOB_NAME"), glueContext, args.asJava)
// @type: DataSource
// @args: [database = "demo-db-dblp-acm", table_name = "dblp_acm_records_csv",
transformation_ctx = "datasource0"]
// @return: datasource0
// @inputs: []
val datasource0 = glueContext.getCatalogSource(database = "demo-db-dblp-acm",
tableName = "dblp_acm_records_csv", redshiftTmpDir = "", transformationContext =
"datasource0").getDynamicFrame()
// @type: FindMatches
// @args: [transformId = "tfm-123456789012", emitFusion = false,
survivorComparisonField = "<primary_id>", transformation_ctx = "findmatches1"]
// @return: findmatches1
// @inputs: [frame = datasource0]
val findmatches1 = FindMatches.apply(frame = datasource0, transformId
= "tfm-123456789012", transformationContext = "findmatches1",
computeMatchConfidenceScores = true)

// Repartition the previous DynamicFrame into a single partition.
val single_partition = findmatches1.repartition(1)

// @type: DataSink
// @args: [connection_type = "s3", connection_options = {"path": "s3://aws-
glue-ml-transforms-data/sal"}, format = "csv", transformation_ctx = "datasink2"]
// @return: datasink2
// @inputs: [frame = findmatches1]
val datasink2 = glueContext.getSinkWithFormat(connectionType =
"s3", options = JsonOptions("{"path": "s3://aws-glue-ml-transforms-
data/sal"}"), transformationContext = "datasink2", format =
"csv").writeDynamicFrame(single_partition)
Job.commit()
}
}

```

6. [ジョブの実行] を選択してジョブの実行を開始します。ジョブリストでジョブのステータスを確認します。ジョブが完了すると、[ML 変換] の [履歴] タブに、[ETL ジョブ] タイプの新しい [実行 ID] 行が追加されます。
7. [ジョブ]、[履歴] タブの順に移動します。このペインには、ジョブの実行が一覧表示されます。実行の詳細を参照するには、[ログ] を選択します。実行が完了したときの実行ステータスが [成功] になっていることを確認します。

ステップ 6: Amazon S3 の出力データを確認する

このステップでは、ジョブの追加時に選択した Amazon S3 バケットのジョブ実行の出力を確認します。出力ファイルをローカルマシンにダウンロードして、一致するレコードが識別されていることを確認します。

1. <https://console.aws.amazon.com/s3/> で Amazon S3 コンソールを開きます。
2. ジョブ demo-etl-dblp-acm のターゲット出力ファイルをダウンロードします。このファイルをスプレッドシートアプリケーションで開きます (ファイルを正常に開くには、必要に応じてファイル拡張子 .csv を追加します)。

次の図は、Microsoft Excel で開いた出力の抜粋を示しています。

	B	C	D	E	F	G	H	I
	title	authors	venue	year	source	primary_id	match_id	match_confidence_score
1	Semantic Integration of Environmental Models for Application to Global Information S	D. Scott Mackay	SIGMOD Record	1999	DBLP	3092	0	0.830985237
2	Semantic integration of environmental models for application to global information s	D. Scott Mackay	ACM SIGMOD Recor	1999	ACM	3590	0	0.830985237
3	Estimation of Query-Result Distribution and Its Application in Parallel-Join Load	Balan Viswanath Poosala, Yannis E. I	VLDB	1996	DBLP	3435	1	0.801848258
4	Estimation of Query-Result Distribution and Its Application in Parallel-Join Load	Balan Viswanath Poosala, Yannis E. I	Very Large Data Bas	1996	ACM	2491	1	0.801848258
5	Incremental Maintenance for Non-Distributive Aggregate Functions	Themistoklis Palpanas, Richar	VLDB	2002	DBLP	4638	2	0.697109993
6	Cost-based Selection of Path Expression Processing Algorithms in Object-Oriented Da	Zhao-Hui Tang, Georges Gardin	Very Large Data Bas	1996	DBLP	3768	3	0.791241276
7	Cost-based Selection of Path Expression Processing Algorithms in Object-Oriented Da	Georges Gardin, Jean-Rober	Very Large Data Bas	1996	ACM	5926	3	0.791241276
8	Benchmarking Spatial Join Operations with Spatial Output	Erik G. Hoel, Hanan Samet	Very Large Data Bas	1995	ACM	9739	4	0.723535024
9	Benchmarking Spatial Join Operations with Spatial Output	Erik G. Hoel, Hanan Samet	VLDB	1995	DBLP	8124	4	0.723535024
10	Efficient geometry-based similarity search of 3D spatial databases	Daniel A. Keim	International Confe	1999	ACM	5647	5	0.786350237
11	Efficient Geometry-based Similarity Search of 3D Spatial Databases	Daniel A. Keim	SIGMOD Conference	1999	DBLP	3432	5	0.786350237
12	Mining the World Wide Web: An Information Search Approach - Book Review	Aris M. Oukel	SIGMOD Record	2002	DBLP	6790	6	0.697109993
13	Enhanced Abstract Data Types in Object-Relational Databases	Praveen Seshadri	VLDB J.	1998	DBLP	3617	7	0.827350237
14	Enhanced abstract data types in object-relational databases	Praveen Seshadri	The VLDB Journal &	1998	ACM	4906	7	0.827350237
15	Report on DART '96: Databases: Active and Real-Time (Concepts meet Practice)	Nandit Soparkar, Krithi Ramam	SIGMOD Record	1997	DBLP	7937	8	0.708350237
16	Report on DART '96: databases: active and real-time (concepts meet practice)	Krithi Ramamritham, Nandit S	ACM SIGMOD Recor	1997	ACM	8193	8	0.708350237
17	UNISQL's next-generation object-relational database management system	Albert D'Andrea, Phil Janus	ACM SIGMOD Recor	1996	ACM	8491	9	0.818340237
18	UNISQL's Next-Generation Object-Relational Database Management System	Phil Janus, Albert D'Andrea	SIGMOD Record	1996	DBLP	4869	9	0.818340237

データソースおよびターゲットファイルの両方に 4,911 個のレコードがあります。ただし、Find matches 変換には、出力内の一致するレコードを識別するために match_id という別の列が追加されています。同じ match_id の行は、一致するレコードと見なされます。match_confidence_score は 0~1 の間の数値で、Find matches によって検出された次の値で一致の品質を推定します。

3. 出力ファイルを match_id で並べ替えると、どのレコードが一致しているか簡単にわかります。他の列の値と比較して、Find matches 変換の結果が適切であるかどうかを判断します。適切でない場合は、さらにラベルを追加することで、引き続き変換をトレーニングします。

ファイルをソートするために別のフィールド (title など) を使用して、同じようなタイトルのレコードが同じ match_id を持っているかどうかを確認することもできます。

インクリメンタルマッチを検索する。

FindMatches 変換を使用すると、レコードに共通の一意の識別子がなく、正確に一致するフィールドがない場合でも、データセット内の重複レコードまたは一致するレコードを識別できます。検索マッチ変換の最初のリリースでは、単一のデータセット内で一致するレコードが識別されました。データ

セットに新しいデータを追加する場合、既存のクリーンなデータセットとマージし、マージされた完全なデータセットに対してマッチングを再実行する必要があります。

インクリメンタルマッチング機能を使用すると、既存の一致したデータセットに対するインクリメンタルレコードとの照合が簡単になります。既存の顧客データセットと見込み客のデータを照合するとします。差分一致機能を使用すると、結果を単一のデータベースまたはテーブルにマージすることで、数十万の新規プロスペクトを既存のプロスペクトと顧客の既存のデータベースと照合できる柔軟性が得られます。新規データセットと既存のデータセット間のみ照合を行うことで、差分一致の検索最適化によって計算時間が短縮され、コストも削減されます。

インクリメンタルマッチングの使用法は、[チュートリアル: AWS Glue で機械学習変換を作成する](#) で説明されているように、「一致を検索」と似ています。このトピックでは、インクリメンタルマッチングの相違点のみについて説明します。

詳細については、[インクリメンタルデータマッチング](#)のブログ記事を参照してください。

インクリメンタルマッチングジョブの実行

次の手順では、次の様に仮定します。

- 既存のデータセット内をクローलして first_records テーブルに到達しました。first_records データセットは、一致したデータセット、または一致したジョブの出力であるはずです。
- AWS Glue バージョン 2.0 を使用して、一致検索変換の作成とトレーニングが完了しました。このバージョンが唯一インクリメンタルマッチを AWS Glue サポートしています。
- ETL 言語は Scala です。Python もサポートされていることに注意してください。
- 既に生成されたモデルは demo-xform と呼ばれます。

1. インクリメンタルデータセットを、テーブル second records にクローलする。
2. AWS Glue コンソールのナビゲーションペインで、[ジョブ] を選択します。
3. [Add job (ジョブの追加)] を選択し、ウィザードの手順に従い、生成されたスクリプトを使用して ETL Spark ジョブを作成します。変換のプロパティ値として以下を選択します。
 - a. 名前で、デモetlを使用する場合。
 - b. IAM ロール に、Amazon S3 ソースデータ、ラベリングファイル、および [AWS GlueAPI オペレーション](#)へのアクセス許可を持つ IAM ロールを選択します。
 - c. [ETL 言語] で、[Scala] を選択します。
 - d. [スクリプトファイル名] で、[demo-etl] を選択します。これは Scala スクリプトのファイル名です。

- e. データソースはfirst recordsを選択します。選択したデータソースは、機械学習変換のデータソーススキーマと一致する必要があります。
 - f. [Transform type (変換タイプ)] で、[Find matching records (一致するレコードの検索)] を選択し、機械学習変換を使用してジョブを作成します。
 - g. インクリメンタルマッチングのオプションを選択し、データ・ソースにsecond_recordsという名前のテーブルを選択します。
 - h. [トランスフォーム] で、ジョブによって使用される機械学習変換である demo-xform-dblp-acm を選択します。
 - i. データターゲットにテーブルを作成するか、データカタログのテーブルを使用してデータターゲットを更新するを選択します。
4. [ジョブを保存してスクリプトを編集する] を選択して、スクリプトエディタページを表示します。
 5. [ジョブの実行] を選択してジョブの実行を開始します。

ビジュアルジョブでの FindMatches の使用

AWS Glue Studio で FindMatches 変換を使用するには、FindMatches API を呼び出す カスタム変換ノードを使用します。カスタム変換の使用に関する詳細については、「[カスタム変換を作成する](#)」を参照してください。

Note

現在、FindMatches API は、Glue 2.0 でのみ動作します。FindMatches API を呼び出すカスタム変換を使用してジョブを実行するには、AWS Glue のバージョンが [ジョブの詳細] タブの Glue 2.0 に表示されていることを確認してください。AWS Glue のバージョンが Glue 2.0 でない場合、ジョブは実行時に失敗し、「'awsglueml.transforms' から 'FindMatches' という名前の変換をインポートできません」というエラーメッセージが表示されます。

前提条件

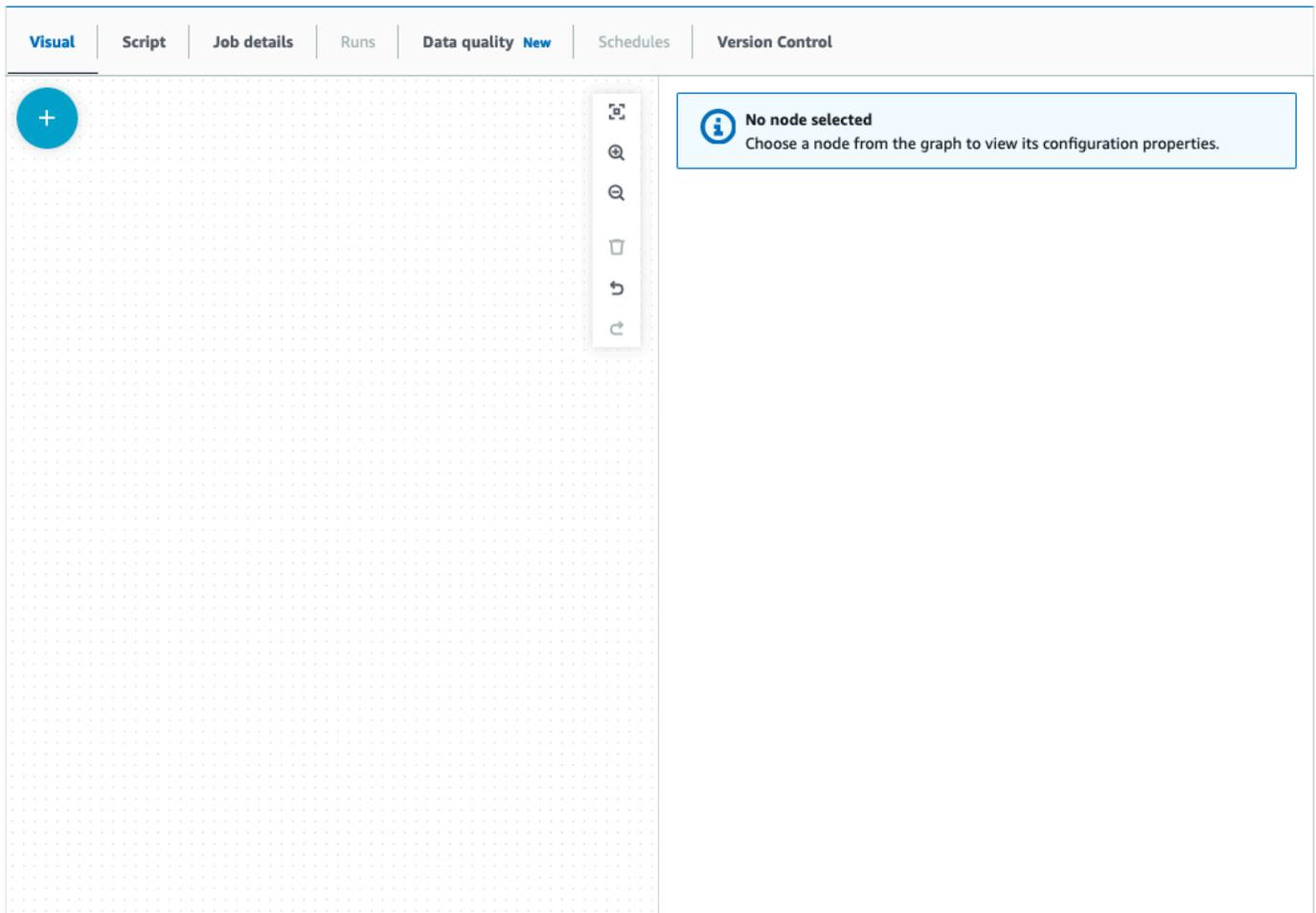
- 一致の検索 変換を使用するには、AWS Glue Studio で <https://console.aws.amazon.com/gluestudio/> コンソールを開きます。

- 機械学習変換を作成します。作成すると、transformId が生成されます。この ID は以下の手順で必要になります。機械学習変換の作成方法の詳細については、「[機械学習変換の追加と編集](#)」を参照してください。

FindMatches 変換の追加

FindMatches 変換を追加するには

- AWS Glue Studio ジョブエディターで、ビジュアルジョブグラフの左上隅にある十字記号をクリックして [リソース] パネルを開き、[データ] タブを選択してデータソースを選択します。これは、一致を確認するデータソースです。



- データソースノードを選択し、ビジュアルジョブグラフの左上隅にある十字記号をクリックして [リソース] パネルを開き、「custom transform」を検索します。カスタム変換ノードを選択してグラフに追加します。カスタム変換はデータソースノードにリンクされています。リンクされていない場合は、カスタム変換ノードをクリックして [ノードプロパティ] タブを選択し、[ノードの親] でデータソースを選択します。

3. ビジュアルグラフの[カスタム変換]ノードをクリックし、[ノードプロパティ]タブを選択し、カスタム変換に名前を付けます。ビジュアルグラフに表示される変換名は、わかりやすい名前に変更することをお勧めします。
4. コードブロックを編集するには、[変換] タブを選択します。ここで、FindMatches API を呼び出すコードを追加できます。

The screenshot displays the AWS Glue console interface. At the top, there are tabs for 'Visual', 'Script', 'Job details', 'Runs', 'Data quality New', 'Schedules', and 'Version Control'. The 'Visual' tab is active, showing a graph with two nodes: 'Data source - S3 bucket Amazon S3' and 'Transform - Custom code ml transform'. A blue arrow points from the data source to the transform node. To the right, the 'Node properties' panel is open, showing the 'Transform' tab. The 'Code block' section contains the following code:

```
1 - def MyTransform (glueContext, dfc) -> DynamicFrameCollection:
2
```

コードブロックには、すぐに始めることができるようにコードが予め入力されています。予め入力されているコードを下記のテンプレートで上書きします。テンプレートには、入力できる transformId のプレースホルダーが表示されています。

```
def MyTransform (glueContext, dfc) -> DynamicFrameCollection:
    dynf = dfc.select(list(dfc.keys())[0])
    from awsglueml.transforms import FindMatches
    findmatches = FindMatches.apply(frame = dynf, transformId = "<your id>")
    return(DynamicFrameCollection({"FindMatches": findmatches}, glueContext))
```

5. ビジュアルグラフの[カスタム変換]ノードをクリックし、ビジュアルジョブグラフの左上隅にある十字記号をクリックして[リソース]パネルを開き、「Select From Collection」を検索します。コレクションには DynamicFrame が 1 つしかないため、デフォルトの選択を変更する必要はありません。
6. 続けて変換を追加したり、結果を保存したりすることができます。結果には、[一致の検索] 列が追加されました。これらの新しい列を下流の変換で参照する場合は、それらを変換出力スキーマに追加する必要があります。それを行う最も簡単な方法は、[データプレビュー] タブを選択し、[スキーマ] タブで [データプレビュースキーマの使用] を選択することです。
7. FindMatches をカスタマイズするには、「apply」メソッドに渡すパラメーターを追加できます。[FindMatches クラス](#)を参照してください。

「FindMatches」段階的な変換の追加

インクリメンタルマッチングの場合、プロセスは FindMatches 変換の追加と同じですが、次の点が異なります。

- カスタム変換の親ノードの代わりに、2 つの親ノードが必要です。
- 最初の親ノードはデータセットでなければなりません。
- 2 番目の親ノードはインクリメンタルデータセットでなければなりません。

テンプレートコードブロックで transformId を incrementalTransformId に置き換えてください。

```
def MyTransform (glueContext, dfc) -> DynamicFrameCollection:
    dfs = list(dfc.values())
    dynf = dfs[0]
    inc_dynf = dfs[1]
    from awsglueml.transforms import FindIncrementalMatches
    findmatches = FindIncrementalMatches.apply(existingFrame = dynf, incrementalFrame
    = inc_dynf,
                                           transformId = "<your id>")
    return(DynamicFrameCollection({"FindMatches": findmatches}, glueContext))
```

- オプションのパラメータについては、「[FindIncrementalMatches クラス](#)」を参照してください。

Apache Spark プログラムを AWS Glue に移行する

Apache Spark は、大規模なデータセットで実行される分散コンピューティングワークロード向けのオープンソースプラットフォームです。AWS Glue は、Spark の機能を活用して ETL に最適化されたエクスペリエンスを提供します。Spark プログラムは、AWS Glue に移行することで機能を最大限に活用できます。AWS Glue は、Amazon EMR の Apache Spark と同等のパフォーマンス向上を実現します。

Spark コードを実行する

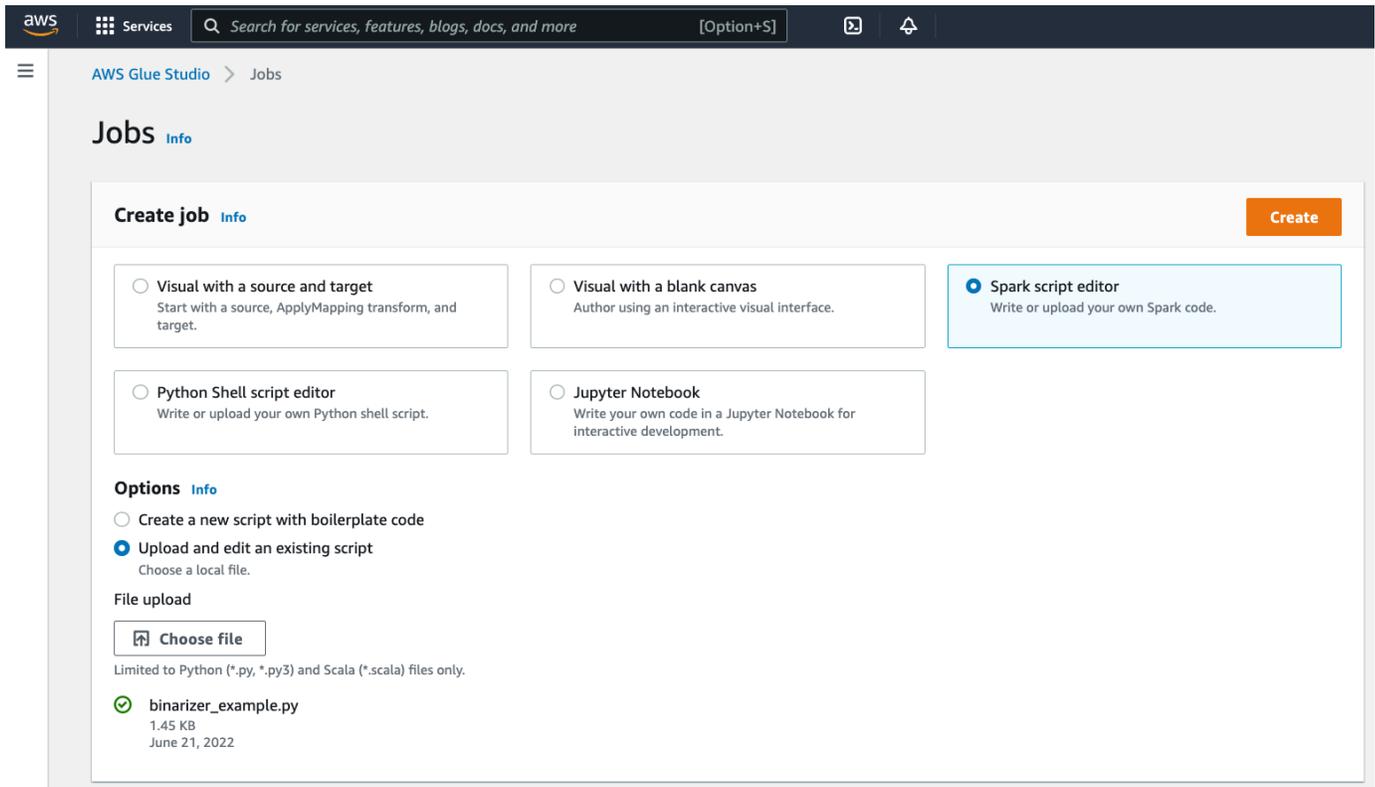
ネイティブ Spark コードは、AWS Glue 環境で追加設定なしですぐに使用できます。スクリプトは多くの場合、インタラクティブセッションに適したワークフローであるコードを繰り返し変更することによって開発されます。ただし、既存のコードの方が AWS Glue ジョブでの実行に適しています。これにより、スクリプトを実行するたびにログとメトリクスをスケジュールし、一貫して取得できます。コンソールから既存のスクリプトをアップロードして編集できます。

1. スクリプトのソースを取得します。この例では、Apache Spark リポジトリのサンプルスクリプトを使用します。[バイナライザーの例](#)
2. AWS Glue コンソールで、左側のナビゲーションペインを展開し、[ETL] > [Jobs] (ジョブ) を選択します。

[Create job] (ジョブの作成) パネルで、[Spark script editor] (Spark スクリプトエディタ) を選択します。[Options] (オプション) セクションが表示されます。[Options] (オプション) で [Upload and edit an existing script] (既存のスクリプトのアップロードと編集) を選択します。

[File upload] (ファイルのアップロード) セクションが表示されます。[File upload] (ファイルのアップロード) で [Choose file] (ファイルを選択) をクリックします。システムファイル選択ダイアログが表示されます。binarizer_example.py を保存した場所に移動してこの場所を選択し、選択を確定します。

[Create job] (ジョブの作成) パネルのヘッダーに [Create] (作成) ボタンが表示されます。このボタンをクリックします。



The screenshot shows the AWS Glue Studio interface. At the top, there is a search bar and navigation icons. The main heading is 'Jobs Info'. Below this, there is a 'Create job Info' section with a 'Create' button. The 'Create job' section contains five options, each with a radio button:

- Visual with a source and target: Start with a source, ApplyMapping transform, and target.
- Visual with a blank canvas: Author using an interactive visual interface.
- Spark script editor: Write or upload your own Spark code.
- Python Shell script editor: Write or upload your own Python shell script.
- Jupyter Notebook: Write your own code in a Jupyter Notebook for interactive development.

Below the options, there is an 'Options Info' section with two radio buttons:

- Create a new script with boilerplate code
- Upload and edit an existing script: Choose a local file.

Under 'File upload', there is a 'Choose file' button. Below it, a list of uploaded files is shown:

- binarizer_example.py
1.45 KB
June 21, 2022

A note below the file list states: 'Limited to Python (*.py, *.py3) and Scala (*.scala) files only.'

3. ブラウザがスクリプトエディタに移動します。ヘッダーで、[Job details] (ジョブの詳細) タブをクリックします。名前と IAM ロールを設定します。AWS Glue IAM ロールに関するガイダンスについては、「[the section called “IAM アクセス許可のセットアップ”](#)」を参照してください。

必要に応じて、[Requested number of workers] (要求されたワーカー数) を 2 に、[Number of retries] (再試行回数) を 1 に設定します。これらのオプションは本番ジョブを実行する場合に役立ちますが、無効にすると、機能をテストする際のエクスペリエンスが合理化されます。

タイトルバーで、[Save] (保存) をクリックし、続いて [Run] (実行) をクリックします。

The screenshot shows the AWS Glue console interface. At the top, there is a navigation bar with the AWS logo, 'Services', a search bar, and utility icons. Below this, the job name 'Binarizer Example' is displayed with a share icon. To the right are buttons for 'Save', 'Delete', 'Actions', and 'Run'. The main content area has tabs for 'Script', 'Job details' (which is selected), 'Runs', and 'Schedules'. The 'Job details' tab shows a 'Basic properties' section with the following fields:

- Name:** Binarizer Example
- Description - optional:** (Empty text area)
- IAM Role:** AWSGlueServiceRole (with a refresh icon)
- Type:** Spark
- Glue version:** Glue 3.0 - Supports spark 3.1, Scala 2, Python 3

4. [Runs] (実行) タブに移動します。ジョブの実行に対応するパネルが表示されます。数分間待機した後、ページが自動的に更新され、[Run status] (実行ステータス) の下に[Succeeded] (成功) が表示されます。

The screenshot shows the AWS Glue console interface for a job named "Binarizer Example". The "Runs" tab is selected, displaying details for a recent job run that succeeded on July 13, 2022 at 12:24:58 PM. The job ID is jr_EXAMPLEID. The console provides a comprehensive overview of the job's execution, including its status, version, start and end times, execution duration, and resource allocation. It also offers links to view logs and performance recommendations.

Job name	Id	Run status	Glue version
Binarizer Example	jr_EXAMPLEID	✔ Succeeded	3.0
Retry attempt number	Start time	End time	Start-up time
Initial run	July 13, 2022 12:24:58 PM	July 13, 2022 12:25:36 PM	7 seconds
Execution time	Last modified on	Trigger name	Security configuration
30 seconds	July 13, 2022 12:25:36 PM	-	-
Timeout	Max capacity	Number of workers	Worker type
2880 minutes	2 DPUs	2	G.1X
Execution class	Log group name	Cloudwatch logs	Performance and debugging recommendations
-	/aws-glue/jobs	<ul style="list-style-type: none"> All logs Output logs Error logs 	<ul style="list-style-type: none"> View in CloudWatch

Input arguments (10)
Arguments used when this job run was executed.

- 出力を調べて、Spark スクリプトが意図したとおりに実行されたことを確認する必要があります。この Apache Spark サンプルスクリプトは、出力ストリームに文字列を書き込む必要があります。これは、実行が成功したジョブのパネル内の [Cloudwatch Logs] の下にある [Output logs] (出力ログ) に移動して確認できます。ジョブ実行 ID をメモします。これは、jr_で始まる [ID] ラベルの下に生成されます。

これにより、CloudWatch コンソールが開き、デフォルトの AWS Glue ロググループ /aws-glue/jobs/output のコンテンツを視覚化するように設定されます。これで、ジョブ実行 ID のログストリームのコンテンツに絞り込まれます。各ワーカーはログストリームを生成し、[Log streams] (ログストリーム) の下の行として表示されます。1 人のワーカーが要求されたコードを実行しているはずですが、正しいワーカーを特定するには、すべてのログストリームを開く必要があります。適切なワーカーが見つかったら、次の画像に示すように、スクリプトの出力が表示されます。

The screenshot shows the AWS CloudWatch console interface. The breadcrumb navigation is: CloudWatch > Log groups > /aws-glue/jobs/output > jr_EXAMPLEID. The main content area is titled "Log events" and includes a search bar, a "View as text" checkbox, and a "Create Metric Filter" button. Below this is a table of log events with columns for "Timestamp" and "Message".

Timestamp	Message
2022-07-13T13:27:33.060-07:00	No older events at this moment. Retry
2022-07-13T13:27:33.062-07:00	2022-07-13 20:27:33,058 main WARN JNDI lookup class is not available because...
2022-07-13T13:27:54.066-07:00	2022-07-13 20:27:33,062 main INFO Log4j appears to be running in a Servlet e... Binarizer output with Threshold = 0.500000
2022-07-13T13:28:02.833-07:00	+++++-----+-----+-----+ id feature binarized_feature +-----+... +-----+-----+-----+ id feature binarized_feature +-----+-----+-----+ 0 0.1 0.0 1 0.8 1.0 2 0.2 0.0 +-----+-----+-----+

At the bottom of the log events list, it says: "No newer events at this moment. Auto retry paused. [Resume](#)".

Spark プログラムの移行に必要な一般的な手順

Spark バージョンサポートの評価

AWS Glue リリースバージョンは、AWS Glue ジョブで使用できる Apache Spark と Python のバージョンを決定します。AWS Glue バージョンとそのサポート対象については、「[the section called “AWS Glue バージョン”](#)」をご覧ください。特定の AWS Glue 機能にアクセスするには、Spark プログラムを新しいバージョンの Spark と互換性があるように更新する必要がある場合があります。

サードパーティライブラリを含める

既存の Spark プログラムの多くは、プライベートアーティファクトとパブリックアーティファクトの両方に依存しています。AWS Glue は、Scala ジョブの JAR スタイルの依存関係と、Python ジョブのホイールとソース Pure-Python 依存関係をサポートします。

Python-Python の依存関係については、「[the section called “Python ライブラリ”](#)」を参照してください。

一般的な Python の依存関係は、一般的に要求される [Pandas](#) ライブラリなどを含めて AWS Glue 環境で提供されています。依存関係は AWS Glue バージョン 2.0 以降に含まれています。提供されるモジュールの詳細については、「[the section called “AWS Glue で提供済みの Python モジュール”](#)」

を参照してください。デフォルトで含まれている依存関係の別のバージョンをジョブに提供する必要がある場合は、`--additional-python-modules` を使用します。ジョブ引数の詳細については、「[the section called “ジョブのパラメータ”](#)」を参照してください。

追加の Python 依存関係は、`--extra-py-files` ジョブ引数で指定できます。Spark プログラムからジョブを移行する場合、このパラメータは良いオプションです。PySpark の `--py-files` フラグと機能的に同等で、同じ制限が適用されるからです。`--extra-py-files` パラメータの詳細については、「[the section called “PySpark ネイティブ機能を備えた Python ファイルを含める”](#)」を参照してください。

新しいジョブについては、`--additional-python-modules` 引数を使用して Python の依存関係を管理できます。この引数を使用すると、より徹底した依存関係管理が可能になります。このパラメータは、Amazon Linux 2 と互換性のあるネイティブコードバインディングを含むホイールスタイルの依存関係をサポートします。

Scala

追加の Scala 依存関係は、`--extra-jars` ジョブ引数で指定できます。依存関係は Amazon S3 でホストされている必要があり、引数の値は、スペースを含まない Amazon S3 パスのカンマ区切りリストである必要があります。依存関係をホストして設定する前に再バンドルすると、設定を管理しやすくなる場合があります。AWS GlueJAR 依存関係には Java バイトコードが含まれており、任意の JVM 言語から生成できます。Java などの他の JVM 言語を使用して、カスタムの依存関係を記述できます。

データソース認証情報の管理

既存の Spark プログラムには、データソースからデータを取得するための複雑な構成またはカスタム構成が付属している場合があります。一般的なデータソース認証フローは、AWS Glue 接続でサポートされます。AWS Glue 接続の詳細については、「[データへの接続](#)」を参照してください。

AWS Glue 接続は、主に次の 2 つの方法 (ライブラリへのメソッド呼び出しと、AWS コンソールで [Additional network connection] (追加のネットワーク接続) を設定する) を通じて、ジョブをさまざまなタイプのデータストアに接続しやすくします。また、接続から情報を取得するため、ジョブ内からの AWS SDK を呼び出すこともできます。

メソッドの呼び出し – AWS Glue 接続は、AWS Glue データカタログと密接に統合されています。これは、データセットに関する情報と、この情報が反映された AWS Glue 接続の操作に使用できるメソッドの管理に利用できます。JDBC 接続で再利用したい既存の認証設定がある場合は、GlueContext の `extract_jdbc_conf` メソッドを介して AWS Glue 接続の設定にアクセスできます。詳細については、「[the section called “extract_jdbc_conf”](#)」を参照してください。

コンソール設定 – AWS Glue ジョブは、関連する AWS Glue 接続を使用して Amazon VPC サブネットへの接続を設定します。セキュリティ資料を直接管理している場合は、AWS コンソールで、NETWORK タイプ [Additional network connection] (追加のネットワーク接続) を指定してルーティングを設定します。AWS Glue 接続 API の詳細については、「[the section called “接続”](#)」を参照してください

Spark プログラムにカスタム認証フローまたは一般的でない認証フローがある場合、セキュリティ資料をハンズオンで管理する必要があるかもしれません。AWS Glue 接続が適切でない場合、セキュリティ資料をシークレットマネージャーで安全にホストし、ジョブで提供される boto3 または AWS SDK を介してアクセスします。

Apache Spark の設定

複雑な移行では、ワークロードに対応するために Spark 設定が変更されることがよくあります。Apache Spark の最新バージョンでは、ランタイム設定を SparkSession で設定できます。AWS Glue 3.0 以降のジョブでは SparkSession が提供され、ランタイム設定の変更ができます。[Apache Spark 設定](#)。Spark のチューニングは複雑で、AWS Glue は、すべての Spark 設定に対するサポートを保証するものではありません。移行に相当する Spark レベルの設定が必要な場合は、サポートにお問い合わせください。

カスタム設定のセット

移行された Spark プログラムは、カスタム設定を選択するように設計されている場合があります。AWS Glue はジョブ引数を使用して、ジョブとジョブの実行レベルで構成を設定できます。ジョブ引数の詳細については、「[the section called “ジョブのパラメータ”](#)」を参照してください。ライブラリを通じて、ジョブのコンテキスト内でジョブの引数にアクセスできます。AWS Glue は、ジョブに設定された引数とジョブ実行時に設定された引数の間で一貫したビューを提供するユーティリティ関数を提供します。Python の「[the section called “getResolvedOptions”](#)」、Scala の「[the section called “GlueArgParser”](#)」を参照してください。

Java コードの移行

「[the section called “サードパーティライブラリ”](#)」で説明されているように、Java や Scala などの JVM 言語で生成されたクラスを依存関係に含めることができます。依存関係には main メソッドが含まれます。AWS Glue Scala ジョブのエントリーポイントとして、依存関係内に main メソッドを使用できます。これにより、Java に main メソッドを記述するか、独自のライブラリ標準にパッケージ化された main メソッドを再使用できます。

依存関係から main メソッドを使用するには、次のいずれかの方法を実行します: デフォルトの GlueApp オブジェクトを提供する編集ペインの内容をクリアします。依存関係内のクラスの完全修

飾名を、キー `--class` を使用したジョブ引数として指定します。これで、ジョブ実行をトリガーできるようになります。

AWS Glue が main メソッドに引き渡す引数の順序や構造は設定できません。既存のコードで、AWS Glue の設定を読み取る必要がある場合、以前のコードとの非互換性を引き起こす可能性があります。getResolvedOptions を使用していると、このメソッドを呼び出すのに適した条件もありません。AWS Glue で生成されたメインメソッドから依存関係を直接呼び出すことをおすすめします。この例については、以下の AWS Glue ETL スクリプトをご覧ください。

```
import com.amazonaws.services.glue.util.GlueArgParser

object GlueApp {
  def main(sysArgs: Array[String]) {
    val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)

    // Invoke static method from JAR. Pass some sample arguments as a String[], one
    // defined inline and one taken from the job arguments, using getResolvedOptions
    com.mycompany.myproject.MyClass.myStaticPublicMethod(Array("string parameter1",
    args("JOB_NAME")))

    // Alternatively, invoke a non-static public method.
    (new com.mycompany.myproject.MyClass).someMethod()
  }
}
```

AWS Glue での Ray ジョブの使用

このセクションでは、AWS Glue for Ray のジョブを使用する方法について説明します。AWS Glue for Ray のスクリプトを記述する方法の詳細については、「[the section called “AWS Glue for Ray”](#)」セクションを参照してください。

トピック

- [AWS Glue for Ray の使用を開始する](#)
- [サポートされている Ray のランタイム環境](#)
- [Ray ジョブのワーカーの考慮](#)
- [Ray ジョブでジョブパラメータを使用する](#)
- [メトリクスによる Ray ジョブのモニタリング](#)

AWS Glue for Ray の使用を開始する

AWS Glue for Ray を使用するには、AWS Glue for Spark で使用しているものと同じ、AWS Glue ジョブとインタラクティブセッションを使用します。AWS Glue ジョブは、同じスクリプトを一定の頻度で繰り返し実行するように設計されていますが、インタラクティブセッションは、プロビジョニングされた同じリソースに対して、コードのスニペットを順番に実行するように設計されています。

AWS Glue ETL と Ray は下部構造が異なるため、スクリプトではさまざまなツール、機能、設定にアクセスすることができます。AWS Glue が管理する新しい計算フレームワークとして、Ray は異なるアーキテクチャを持ちます。また、その機能を説明するポキャブラリーも異なります。詳細については、Ray のドキュメントの「[Architecture Whitepapers](#)」(アーキテクチャホワイトペーパー)を参照してください。

Note

AWS Glue for Ray は、米国東部 (バージニア北部)、米国東部 (オハイオ)、米国西部 (オレゴン)、アジアパシフィック (東京)、欧州 (アイルランド) でご利用いただけます。

AWS Glue Studio コンソールでの Ray ジョブ

AWS Glue Studio コンソールの [ジョブ] ページでは、AWS Glue Studio でジョブを作成する際に、新たに [Ray script editor] を選択できます。コンソールに Ray ジョブを作成するときはこちらを選択します。ジョブとその使用方法についての詳細は、「[AWS Glue Studio でビジュアル ETL ジョブを作成する](#)」を参照してください。

The screenshot shows the 'Jobs' page in AWS Glue Studio. At the top, there is a breadcrumb 'AWS Glue Studio > Jobs'. Below that, the 'Jobs' section is visible. The main content area is titled 'Create job' and includes an 'Info' link and a 'Create' button. There are six radio button options for creating a job:

- Visual with a source and target
Start with a source, ApplyMapping transform, and target.
- Visual with a blank canvas
Author using an interactive visual interface.
- Spark script editor
Write or upload your own Spark code.
- Python Shell script editor
Write or upload your own Python shell script.
- Jupyter Notebook
Write your own code in a Jupyter Notebook for interactive development.
- Ray script editor **New**
Write your own code to run on Ray.

Below these options, there is an 'Options' section with two radio button choices:

- Create a new script with boilerplate code
- Upload and edit an existing script
Choose a local file.

AWS CLI および SDK での Ray ジョブ

AWS CLI での Ray ジョブは、他のジョブと同じ SDK アクションとパラメータを使用します。AWS Glue for Ray では、特定のパラメータに新しい値が導入されています。ジョブの API の詳細については、「[the section called “ジョブ”](#)」を参照してください。

サポートされている Ray のランタイム環境

Spark ジョブでは、GlueVersion によって、AWS Glue for Spark ジョブで利用できる Apache Spark と Python のバージョンが決まります。Python のバージョンは、Spark タイプのジョブでサポートされているバージョンを示します。これは、Ray のランタイム環境の設定方法とは異なります。

Ray ジョブの場合、GlueVersion を 4.0 以降に設定する必要があります。ただし、Ray ジョブで利用できる Ray、Python、その他のライブラリバージョンは、ジョブ定義の Runtime フィールドにより決まります。

Ray2.4 ランタイム環境は、リリースから最低 6 か月間は使用可能です。Ray は進化のスピードが速いため、今後のランタイム環境のリリースを通じてその更新や改良を統合していきます。

有効な値: Ray2.4

ランタイム値	Ray と Python のバージョン
Ray2.4 (AWS Glue 4.0 以降)	Ray 2.4.0 Python 3.9

追加情報

- AWS Glue on Ray のリリースに付属のリリースノートについては、「[the section called “AWS Glue バージョン”](#)」を参照してください。
- ランタイム環境で提供される Python ライブラリについては、「[the section called “Ray ジョブで提供されるモジュール”](#)」を参照してください。

Ray ジョブのワーカーの考慮

AWS Glue は、Ray ジョブを Graviton ベースの新しい EC2 ワーカータイプで実行します。このタイプは Ray ジョブでのみ利用できるものです。これらのワーカーを、Ray の設計の根拠であるワークロードに合わせて適切にプロビジョニングするため、コンピューティングリソースとメモリリソースの比率がほとんどのワーカーとは異なっています。これらのリソースを考慮して、標準データ処理ユニット (DPU) ではなく、メモリ最適化データ処理ユニット (M-DPU) が使用されています。

- 1 基の M-DPU は 4 基の vCPU と 32 GB のメモリに相当します。
- 1 基の DPU は 4 基の vCPU と 16 GB のメモリに相当します。DPU は、AWS Glue with Spark のジョブと対応するワーカーのリソースを考慮して、使用されています。

Ray ジョブでは現在、Z.2X という 1 つのワーカータイプしか使用できません。Z.2X ワーカーは 2 つの M-DPU (8 基の vCPU、64 GB のメモリ) にマップされ、128 GB のディスク容量を有しています。Z.2X マシンには 8 つの Ray ワーカーがあります (vCPU につき 1 つ)。

1 つのアカウントで同時に使用できる M-DPU の数は、サービスクォータの対象となります。AWS Glue のアカウント制限についての詳細は、「[AWS Glue エンドポイントとクォータ](#)」を参照してください。

Ray ジョブで使用できるワーカーノードの数は、ジョブ定義で `--number-of-workers` (NumberOfWorkers) のように指定します。ジョブ API の Ray 値の詳細については、「[the section called “ジョブ”](#)」を参照してください。

また、Ray ジョブが割り当てるワーカーの最小数は、`--min-workers` ジョブパラメータを使用して指定できます。ジョブパラメータについては、「[the section called “リファレンス”](#)」を参照してください。

Ray ジョブでジョブパラメータを使用する

AWS Glue Ray ジョブの引数は、AWS Glue for Spark ジョブの引数を設定するときと同じ方法で設定します。AWS Glue API の詳細については、「[the section called “ジョブ”](#)」を参照してください。AWS Glue Ray ジョブはさまざまな引数を用いて設定できます。これらはこちらのリファレンスに記載されています。また、独自の引数を指定することもできます。

ジョブは、Job Parameters (ジョブパラメータ) 見出しの Job details (ジョブの詳細) タブからコンソールで設定できます。また、ジョブに DefaultArguments、またはジョブ実行に Arguments を設定することで、AWS CLI を使用してジョブを設定することもできます。デフォルトの引数とジョブパラメータは、複数回実行してもジョブに対して有効です。

以下は、特殊なパラメータをセットするために `--arguments` を使用する、ジョブ実行の構文例です。

```
$ aws glue start-job-run --job-name "CSV to CSV" --arguments='--scriptLocation="s3://my_glue/libraries/test_lib.py",--test-environment="true"'
```

引数を設定すると、環境変数を使用して Ray ジョブ内からジョブパラメータにアクセスできます。そうすることで、ジョブを実行するたびにジョブの設定を行うことができます。環境変数の名前は、先頭の `--` を除いたジョブ引数の名前になります。

例えば、上記の例では、変数の名前は `scriptLocation` および `test-environment` になります。次に、標準ライブラリで利用できる方法を使用して引数を取得します (`test_environment = os.environ.get('test-environment')`)。Python で環境変数にアクセスする方法の詳細については、Python のドキュメントの「[os module](#)」を参照してください。

Ray ジョブがログを生成する方法を設定する

デフォルトでは、Ray ジョブは CloudWatch と Amazon S3 に送信されるログとメトリクスを生成します。`--logging_configuration` パラメータを使用してログの生成方法を変更できます。現在のところ、これを使用して Ray ジョブがさまざまなタイプのログを生成しないようにすることができます。このパラメータは、変更したいログや動作に対応するキーを持つ JSON オブジェクトを取ります。次のキーがサポートされています。

- `CLOUDWATCH_METRICS` — ジョブの状態を視覚化するために使用できる CloudWatch メトリクスシリーズを設定します。メトリクスの詳細については、「[the section called “Ray ジョブメトリクス”](#)」を参照してください。
- `CLOUDWATCH_LOGS` — ジョブの実行ステータスに関する Ray アプリケーションレベルの詳細を提供する CloudWatch ログを設定します。ログの詳細については、「[the section called “Ray エラーをトラブルシューティングする”](#)」を参照してください。
- `S3` — AWS Glue が Amazon S3 に何を書き込むかを設定します。主として CloudWatch ログと同様の情報ですが、ログストリームではなくファイルとして書き込まれます。

Ray ログ記録の動作を無効にするには、値 `{"IS_ENABLED": "False"}` を指定します。例えば、CloudWatch メトリクスと CloudWatch ログを無効にするには、以下の設定を指定します。

```
"--logging_configuration": "{\"CLOUDWATCH_METRICS\": {\"IS_ENABLED\": \"False\"},  
  \"CLOUDWATCH_LOGS\": {\"IS_ENABLED\": \"False\"}}"
```

リファレンス

Ray ジョブは、Ray ジョブおよびジョブ実行の、スクリプト環境のセットアップに使用できる、以下の引数名を認識します。

- `--logging_configuration` — Ray ジョブによって作成されるさまざまなログの生成を停止するために使用されます。これらのログは、すべての Ray ジョブでデフォルトで生成されます。形式: 文字列がエスケープされた JSON オブジェクト。詳細については、「[the section called “Ray ジョブがログを生成する方法を設定する”](#)」を参照してください。
- `--min-workers` — Ray ジョブに割り当てられるワーカーノードの最小数。ワーカーノードは、複数のレプリカを実行できます。レプリカは仮想 CPU ごとに 1 つ割り当てられます。形式: 整数。最小値: 0。最大値: ジョブ定義で `--number-of-workers` (`NumberOfWorkers`) に指定できる値。ワーカーノードの考慮に関する詳細は、「[the section called “Ray ジョブのワーカーの考慮”](#)」を参照してください。
- `--object_spilling_config` — AWS Glue for Ray は、Ray のオブジェクトストアで利用できるスペースを拡張する方法として、Amazon S3 の使用をサポートしています。この動作を有効にするには、このパラメータを使用して Ray に object spilling の JSON 設定オブジェクトを提供します。Ray の object spilling 設定の詳細については、Ray のドキュメントの「[Object Spilling](#)」を参照してください。形式: JSON オブジェクト。

AWS Glue for Ray では、一度にディスクへのスピリングまたは Amazon S3 へのスピリングのいずれかのみがサポートされています。この制限内であれば、スピリングに複数の場所を指定できます。Amazon S3 にスピリングする場合は、IAM 権限をこのバケットのジョブに追加することも必要になります。

CLI の構成として JSON オブジェクトを提供する場合は、JSON オブジェクトの文字列をエスケープした文字列として提供する必要があります。例えば、1 つの Amazon S3 パスにスピリングする文字列の値は、`"{\\"type\\": \\"smart_open\\", \\"params\\": {\\"uri\\": \\"s3path\\"}}"` のようになります。AWS Glue Studio で、このパラメータを特別な形式のない JSON オブジェクトとして指定します。

- `--object_store_memory_head` — Ray のヘッドノードの Plasma オブジェクトストアに割り当てられるメモリ。このインスタンスは、クラスター管理のサービスおよびワーカーレプリカを実行します。この値は、ウォームスタート後のインスタンスの空きメモリの割合を表します。このパラメータは、メモリを大量に消費するワークロードの調整に使用します。ほとんどのユースケースにはデフォルト値を使用できます。形式: 正の整数。最小値: 1。最大値: 100。

Plasma の詳細については、Ray のドキュメントで「[The Plasma In-Memory Object Store](#)」(Plasma インメモリオブジェクトストア)を参照してください。

- `--object_store_memory_worker` — Ray のワーカーノードの Plasma オブジェクトストアに割り当てられるメモリ。これらのインスタンスはワーカーレプリカのみを実行します。この値は、ウォームスタート後のインスタンスの空きメモリの割合を表します。このパラメータは、メモリを大量に消費するワークロードを調整するために使用されます。ほとんどのユースケースにはデフォルト値を使用できます。形式: 正の整数。最小値: 1。最大値: 100。

Plasma の詳細については、Ray のドキュメントで「[The Plasma In-Memory Object Store](#)」(Plasma インメモリオブジェクトストア)を参照してください。

- `--pip-install` — インストールされる Python パッケージのセット。この引数を使用して PyPI からパッケージをインストールできます。形式: カンマ区切りリスト。

PyPI パッケージのエントリは、ターゲットパッケージの PyPI 名とバージョンを含む `package==version` の形式です。エントリは、Python のバージョンマッチングを使用してパッケージとバージョンをマッチさせるため、シングリコールではなく、`==` のようにします。バージョンマッチングの演算子は他にもあります。詳細については、Python ウェブサイトの「[PEP 440](#)」を参照してください。`--s3-py-modules` を使用してカスタムのモジュールを提供することもできます。

- `--s3-py-modules` — Python モジュールのディストリビューションをホストする Amazon S3 パスのセットです。形式: カンマで区切られたリスト。

これを使用すれば、独自のモジュールを Ray ジョブに配布できます。`--pip-install` を使用して PyPI からモジュールを提供することもできます。AWS Glue ETL とは異なり、カスタムのモジュールは pip では設定されず、ディストリビューション用に Ray に渡されます。詳細については、「[the section called “Ray ジョブ用の Python モジュールを追加する”](#)」を参照してください。

- `--working-dir` — Amazon S3 でホストされている .zip ファイルへのパスです。このパスには、Ray ジョブを実行しているすべてのノードに配布されるファイルが格納されています。形式: 文字列。詳細については、「[the section called “Ray ジョブへのファイルの提供”](#)」を参照してください。

メトリクスによる Ray ジョブのモニタリング

AWS Glue Studio と Amazon CloudWatch を使用して、Ray ジョブのモニタリングができます。CloudWatch は、分析が可能にする Ray を使用した AWS Glue から未加工のメトリクスを収集

して処理します。これらのメトリクスは AWS Glue Studio コンソールで確認できます。つまり、ジョブの実行をモニタリングできます。

AWS Glue のモニタリング方法についての概要は、「[the section called “CloudWatch メトリクスの使用”](#)」を参照してください。AWS Glue で公開されている CloudWatch メトリクスの使用方法についての概要は、「[the section called “CloudWatch によるモニタリング”](#)」を参照してください。

AWS Glue コンソールで Ray ジョブをモニタリングする

ジョブ実行の詳細ページにある [Run details] セクションには、利用可能なジョブメトリクスを視覚化した、事前構築済みの集計グラフを表示できます。AWS Glue Studio は、ジョブが実行されると、その都度ジョブメトリクスを CloudWatch に送信します。これらを使用することで、クラスターとタスクのプロファイルを作成したり、各ノードに関する詳細情報にアクセスしたりできます。

利用可能なメトリクスグラフの詳細については、「[the section called “Ray ジョブ実行の Amazon CloudWatch メトリクスを表示”](#)」を参照してください。

CloudWatch の Ray ジョブメトリクスの概要

CloudWatch で詳細モニタリングが有効になっている場合、Ray のメトリクスが公開されます。メトリクスは、CloudWatch の Glue/Ray 名前空間に公開されます。

• インスタンスメトリクス

ジョブに割り当てられたインスタンスの CPU、メモリ、およびディスクの使用率に関するメトリクスを公開します。これらのメトリクスは、ExecutorId、ExecutorType および host などの特徴で識別されます。これらのメトリクスは、標準 Linux CloudWatch エージェントのメトリクスのサブセットです。CloudWatch のドキュメントで、メトリクスの名前と特徴に関する情報を確認できます。詳細については、「[CloudWatch エージェントにより収集されるメトリクス](#)」を参照してください。

• Ray クラスターメトリクス

スクリプトを実行する Ray プロセスからこの名前空間にメトリクスを転送し、最も重要なものを提供します。使用可能なメトリクスは Ray のバージョンによって異なる場合があります。ジョブで実行されている Ray のバージョンについての詳細は、「[the section called “AWS Glue バージョン”](#)」を参照してください。

Ray はメトリクスをインスタンスレベルで収集します。また、タスクとクラスターのメトリクスも提供します。Ray の基盤となるメトリクス戦略の詳細については、Ray のドキュメントの「[Metrics](#)」を参照してください。

Note

Glue/Job Metrics/ 名前空間に Ray のメトリクスが公開されることはありません。この名前空間は、AWS Glue ETL ジョブにのみ使用されます。

AWS Glue での Python シェルジョブに関するジョブプロパティの設定

Python シェルジョブを使用して、AWS Glue でシェルとして Python スクリプトを実行できます。Python シェルジョブを使用すると、Python 3.6 または Python 3.9 と互換性のあるスクリプトを実行することができます。

トピック

- [制限事項](#)
- [Python シェルジョブのジョブプロパティの定義](#)
- [サポートされている Python シェルジョブのライブラリ](#)
- [独自の Python ライブラリの提供](#)
- [AWS Glue の Python シェルジョブで AWS CloudFormation を使用する](#)

制限事項

Python シェルジョブには、次の制限があることに注意してください。

- Python シェルを使用してジョブのブックマークを使用することはできません。
- Python 3.9 以降では、.egg ファイルとして Python ライブラリをパッケージ化することはできません。代わりに .whl を使用してください。
- S3 データの一時的なコピーには制限があるため、--extra-files オプションは使用できません。

Python シェルジョブのジョブプロパティの定義

このセクションでは、AWS Glue Studio でのジョブプロパティの定義、または AWS CLI の使用について説明します。

AWS Glue Studio

AWS Glue Studio で Python シェルジョブを定義する場合、次のいくつかのプロパティを指定します。

IAM ロール

ジョブの実行とデータストアへのアクセスに使用されるリソースの認証に使用する AWS Identity and Access Management (IAM) ロールを指定します。AWS Glue でジョブを実行するためのアクセス権限の詳細については、[AWS Glue の Identity and Access Management](#) を参照してください。

タイプ

[Python シェル] を選択して、`pythonshell` という名前のジョブコマンドを使用して Python スクリプトを実行します。

Python バージョン

Python のバージョンを選択します。デフォルトは Python 3.9 です。有効なバージョンは Python 3.6 および Python 3.9 です。

共通分析ライブラリを読み込む (推奨)

Python 3.9 用の共通ライブラリを Python シェルに含めるには、このオプションを選択します。

ライブラリがカスタムであったり、インストール済みのライブラリと競合する場合は、共通ライブラリをインストールしないことを選択することができます。しかしながら、共通ライブラリ以外のライブラリを追加でインストールすることは可能です。

このオプションを選択する場合、`library-set` オプションが `analytics` に設定されます。このオプションの選択を外す場合、`library-set` オプションが `none` に設定されます。

スクリプトのファイル名とスクリプトパス

スクリプトのコードでジョブの手続きロジックを定義します。Amazon Simple Storage Service (Amazon S3) でスクリプト名と場所を指定します。パスのスクリプトディレクトリと同じ名前のファイルが存在していないことを確認します。スクリプトの使用の詳細については、[AWS Glue プログラミングガイド](#) を参照してください。

スクリプト

スクリプトのコードでジョブの手続きロジックを定義します。Python 3.6 または Python 3.9 でスクリプトをコーディングできます。スクリプトは AWS Glue Studio で編集できます。

データ処理単位

このジョブの実行に割り当てられる AWS Glue データ処理ユニット (DPU) の最大数。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表](#)」を参照してください。

値を 0.0625 または 1 に設定できます。デフォルトは 0.0625 です。いずれの場合も、インスタンスのローカルディスクは 20 GB になります。

CLI

また、以下の例のように、AWS CLI を使用して Python シェルジョブを作成することもできます

```
aws glue create-job --name python-job-cli --role Glue_DefaultRole
  --command '{"Name" : "pythonshell", "PythonVersion": "3.9", "ScriptLocation" :
"s3://DOC-EXAMPLE-BUCKET/scriptname.py"}'
  --max-capacity 0.0625
```

Note

--glue-version パラメータは AWS Glue シェルジョブに適用されないため、AWS Glue のバージョンを指定する必要はありません。指定されたバージョンは無視されます。

AWS CLI で作成したジョブはデフォルトで Python 3 になります。有効な Python バージョンは 3 (3.6 に対応) と 3.9 です。Python 3.6 を指定するには、--command パラメータにこのタプルを追加します: "PythonVersion": "3"。

Python 3.9 を指定するには、--command パラメータにこのタプルを追加します: "PythonVersion": "3.9"。

Python シェルジョブによって使用される最大容量を設定するには、--max-capacity パラメータを指定します。Python シェルジョブに --allocated-capacity パラメータは使用できません。

サポートされている Python シェルジョブのライブラリ

Python 3.9 を使用する Python シェルでは、必要に応じて事前にパッケージ済みのライブラリセットを使用するようにライブラリセットを選択できます。ライブラリセットを選択するには、library-set オプションを使用します。有効な値は、analytics および none です。

Python シェルジョブを実行する環境は、次のライブラリをサポートしています。

Python バージョン	Python 3.6	Python 3.9	
ライブラリセット	該当なし	分析	なし
avro		1.11.0	
awscli	116.242	1.23.5	1.23.5
awswrangler		2.15.1	
botocore	1.12.232	1.24.21	1.23.5
boto3	1.9.203	1.21.21	
Elasticsearch		8.2.0	
numpy	1.16.2	1.22.3	
pandas	0.24.2	1.4.2	
psycopg2		2.9.3	
pyathena		2.5.3	
PyGreSQL	5.0.6		
PyMySQL		1.0.2	
pyodbc		4.0.32	
pyorc		0.6.0	
redshift-connector		2.0.907	
リクエスト	2.22.0	2.27.1	
scikit-learn	0.20.3	1.0.2	
scipy	1.2.1	1.8.0	

Python バージョン	Python 3.6	Python 3.9
SQLAlchemy		1.4.36
s3fs		2022.3.0

科学計算用の Python シェルジョブで、NumPy ライブラリを使用できます。詳細については、「[NumPy](#)」を参照してください。以下の例は、Python シェルジョブで使用できる NumPy スクリプトを示しています。このスクリプトは、「Hello world」といくつかの数値計算の結果を出力します。

```
import numpy as np
print("Hello world")

a = np.array([20,30,40,50])
print(a)

b = np.arange( 4 )

print(b)

c = a-b

print(c)

d = b**2

print(d)
```

独自の Python ライブラリの提供

PIP を使用する

Python 3.9 を使用する Python シェルでは、ジョブレベルで追加の Python モジュールや異なるバージョンを指定することもできます。--additional-python-modules オプションでコンマ区切りの Python モジュールのリストを指定することで、新しいモジュールを追加したり、既存のモジュールのバージョンを変更したりできます。Python シェルジョブを使用する場合、Amazon S3 でホストされているカスタム Python モジュールにこのパラメータを指定することはできません。

例えば、scikit-learn モジュールを更新したり新しく追加したりするには、次のキー/値を使用します: "--additional-python-modules", "scikit-learn==0.21.3"

AWS Glue は Python パッケージインストーラ (pip3) を使用して追加のモジュールをインストールします。--additional-python-modules 値の中で、追加の pip3 オプションを渡すことができます。例えば、"scikit-learn==0.21.3 -i https://pypi.python.org/simple/" と指定します。pip3 からの非互換性や制限は、すべて適用されます。

Note

将来的に互換性が失われないように、Python 3.9 用にビルドされたライブラリを使用することをお勧めします。

Egg または Whl ファイルを使用する

1 つ以上の Python ライブラリが .egg または .whl ファイルとして既にパッケージ化されている場合があります。その場合は、次の例のように、「--extra-py-files」フラグの下で AWS Command Line Interface (AWS CLI) を使用して、それらをジョブに指定できます。

```
aws glue create-job --name python-redshift-test-cli --role role --command '{"Name" : "pythonshell", "ScriptLocation" : "s3://MyBucket/python/library/redshift_test.py"}' --connections Connections=connection-name --default-arguments '{"--extra-py-files" : ["s3://DOC-EXAMPLE-BUCKET/EGG-FILE", "s3://DOC-EXAMPLE-BUCKET/WHEEL-FILE"]}'
```

Python ライブラリから .egg ファイルまたは .whl ファイルを作成する方法がわからない場合は、次のステップを使用してください。この例は、macOS、Linux、および Windows Subsystem for Linux (WSL) で適用できます。

Python の .egg ファイルまたは .whl ファイルを作成するには

1. Virtual Private Cloud (VPC) で Amazon Redshift クラスターを作成し、テーブルにデータを追加します。
2. クラスターを作成した VPC-SecurityGroup-Subnet の組み合わせの AWS Glue 接続を作成します。接続に成功するかどうかをテストします。
3. redshift_example という名前のディレクトリを作成して、setup.py という名前のファイルを作成します。次のコードを setup.py に貼り付けます。

```
from setuptools import setup
```

```
setup(
    name="redshift_module",
    version="0.1",
    packages=['redshift_module']
)
```

4. `redshift_example` ディレクトリに `redshift_module` ディレクトリを作成します。`redshift_module` ディレクトリに、ファイル `__init__.py` と `pygresql_redshift_common.py` を作成します。
5. `__init__.py` ファイルは空のままにします。`pygresql_redshift_common.py` 内に次のコードを貼り付けます。`port`、`db_name`、`user`、および `password_for_user` を、Amazon Redshift クラスターに固有の詳細に置き換えます。`table_name` を Amazon Redshift のテーブルの名前と置き換えます。

```
import pg

def get_connection(host):
    rs_conn_string = "host=%s port=%s dbname=%s user=%s password=%s" % (
        host, port, db_name, user, password_for_user)

    rs_conn = pg.connect(dbname=rs_conn_string)
    rs_conn.query("set statement_timeout = 1200000")
    return rs_conn

def query(con):
    statement = "Select * from table_name;"
    res = con.query(statement)
    return res
```

6. `redshift_example` ディレクトリに変更します (現在のディレクトリではない場合)。
7. 次のいずれかを行います。
 - `.egg` ファイルを作成するには、次のコマンドを実行します。

```
python setup.py bdist_egg
```

- `.whl` ファイルを作成するには、次のコマンドを実行します。

```
python setup.py bdist_wheel
```

8. 前述のコマンドに必要な依存関係をインストールします。
9. このコマンドでは、`dist` ディレクトリにファイルを作成します。
 - `egg` ファイルを作成した場合は、`redshift_module-0.1-py2.7.egg` という名前になります。
 - `wheel` ファイルを作成した場合は、`redshift_module-0.1-py2.7-none-any.whl` という名前になります。

このファイルを Amazon S3 にアップロードします。

この例では、アップロードするファイルパスは `s3://DOC-EXAMPLE-BUCKET/EGG-FILE` または `s3://DOC-EXAMPLE-BUCKET/WHEEL-FILE` に配置します。

10. AWS Glue ジョブのスクリプトとして使用する Python ファイルを作成し、このファイルに次のコードを追加します。

```
from redshift_module import pygresql_redshift_common as rs_common

con1 = rs_common.get_connection(redshift_endpoint)
res = rs_common.query(con1)

print "Rows in the table cities are: "

print res
```

11. 前述のファイルを Amazon S3 にアップロードします。この例では、アップロードするファイルパスは `s3://DOC-EXAMPLE-BUCKET/scriptname.py` に配置します。
12. このスクリプトを使用して、Python シェルジョブを作成します。AWS Glue コンソールの [Job properties] (ジョブプロパティ) ページの [Python library path] (Python ライブラリパス) ボックスで `.egg/.whl` ファイルへのパスを指定します。複数の `.egg/.whl` ファイルおよび Python ファイルがある場合は、このボックスにカンマ区切りリストとして指定します。

`.egg` ファイルを変更または名前変更する場合、ファイル名は "python setup.py bdist_egg" コマンドによって生成されたデフォルト名を使用するか、Python モジュールの命名規則に従う必要があります。詳細については、[Style Guide for Python Code](#) を参照してください。

次の例のように、AWS CLI でコマンドを使用してジョブを作成します。

```
aws glue create-job --name python-redshift-test-cli --role Role --command
'{"Name" : "pythonshell", "ScriptLocation" : "s3://DOC-EXAMPLE-BUCKET/
scriptname.py"}'
    --connections Connections="connection-name" --default-arguments '{"--extra-
py-files" : ["s3://DOC-EXAMPLE-BUCKET/EGG-FILE", "s3://DOC-EXAMPLE-BUCKET/WHEEL-
FILE"]}'
```

ジョブが実行されると、スクリプトは、Amazon Redshift クラスターの `[table_name]` テーブルに作成された行をプリントします。

AWS Glue の Python シェルジョブで AWS CloudFormation を使用する

AWS Glue の Python シェルジョブで AWS CloudFormation を使用できます。以下に例を示します。

```
AWSTemplateFormatVersion: 2010-09-09
Resources:
  Python39Job:
    Type: 'AWS::Glue::Job'
    Properties:
      Command:
        Name: pythonshell
        PythonVersion: '3.9'
        ScriptLocation: 's3://bucket/location'
      MaxRetries: 0
      Name: python-39-job
      Role: RoleName
```

Python シェルジョブ用の Amazon CloudWatch Logs グループは、`/aws-glue/python-jobs/output` です。エラーについては、ロググループ `/aws-glue/python-jobs/error` を参照してください。

AWS Glue のモニタリング

モニタリングは、AWS Glue およびその他の AWS ソリューションの信頼性、可用性、およびパフォーマンスを維持する上で重要な部分です。AWS には、AWS Glue を監視したり、問題が発生し

たときに報告したり、必要に応じて自動的にアクションを実行するために使用する監視ツールが用意されています。

以下の自動化されたモニタリングツールを使用して、AWS Glue をモニタリングし、問題が発生したときにレポートできます。

- Amazon CloudWatch Events は、AWS リソースの変更を示すシステムイベントをほぼリアルタイムのストリームとして提供します。CloudWatch Events によって、自動イベント駆動型コンピューティングが有効になります。特定のイベントを監視し、これらのイベントが発生したときに他の AWS サービスで自動アクションをトリガーするルールを記述できます。詳細については、「[Amazon CloudWatch Events ユーザーガイド](#)」を参照してください。
- Amazon CloudWatch Logs によって Amazon EC2 インスタンス、AWS CloudTrail、などからのログファイルのモニタリング、保存、アクセスができます。CloudWatch Logs は、ログファイル内の情報をモニタリングし、特定のしきい値が満たされたときに通知します。高い耐久性を備えたストレージにログデータをアーカイブすることも可能です。詳細については、「[Amazon CloudWatch Logs ユーザーガイド](#)」を参照してください。
- AWS CloudTrail は、AWS アカウントにより、またはそのアカウントに代わって行われた API コールや関連イベントを取得し、指定した Amazon S3 バケットにログファイルを配信します。AWS を呼び出すユーザーとアカウント、呼び出しの送信元 IP アドレス、および呼び出しが発生した時刻を特定できます。詳細については、[AWS CloudTrail ユーザーガイド](#)を参照してください。

さらに、AWS Glue コンソールで以下のインサイトにアクセスし、ジョブのデバッグやプロファイリングに役立てることができます。

- Spark ジョブ – 選択した CloudWatch メトリクスシリーズの視覚化を確認できます。また、新しいジョブは Spark UI にアクセスできます。詳細については、「[the section called “スパークジョブのモニタリング”](#)」を参照してください。
- Ray ジョブ – 選択した CloudWatch メトリクスシリーズの視覚化を確認できます。詳細については、「[the section called “Ray ジョブメトリクス”](#)」を参照してください。

トピック

- [AWSAWS Glue のタグ](#)
- [CloudWatch Events を使用した AWS Glue の自動化](#)
- [AWS Glue リソースのモニタリング](#)
- [AWS Glue による AWS CloudTrail API コールログ記録](#)

AWSAWS Glue のタグ

AWS Glue リソースを管理しやすくするために、オプションでいくつかの AWS Glue リソースタイプに独自のタグを割り当てることができます。タグとは、AWS リソースに割り当てられるラベルです。タグはそれぞれ、1つのキーとオプションの1つの値で設定されており、どちらもお客様側が定義します。AWS Glue のタグを使用して、自分のリソースを整理し識別することができます。タグを使用してコスト計算レポートを作成し、リソースへのアクセスを制限することができます。AWS Identity and Access Management を使用している場合は、AWS アカウント内のどのユーザーがタグを作成、編集、削除する許可を持つかをコントロールできます。タグ関連の API を呼び出すためのアクセス許可に加えて、接続でタグ付け API を呼び出すための `glue:GetConnection` アクセス許可と、データベースでタグ付け API を呼び出すための `glue:GetDatabase` アクセス許可も必要です。詳細については、「[AWS Glue での ABAC](#)」を参照してください。

AWS Glue では、以下のリソースにタグを付けることができます。

- Connection
- データベース
- Crawler
- インタラクティブセッション
- 開発エンドポイント
- ジョブ
- Trigger トリガー)
- ワークフロー
- ブループリント
- 機械学習変換
- データ品質ルールセット
- ストリームスキーマ
- ストリームスキーマレジストリ

Note

ベストプラクティスとして、これらの AWS Glue リソースのタグ付けを可能にするには、常にポリシーに `glue:TagResource` アクションを含めます。

AWS Glue でタグを使用するときは、以下について検討します。

- エンティティ 1 つにつき、最大 50 までのタグがサポートされています。
- AWS Glue では、タグをキーと値のペアのリスト ({"string": "string" ...} 形式) で指定します。
- オブジェクトにタグを作成する場合、タグキーは必須で、タグ値はオプションです。
- タグキーとタグ値の大文字と小文字は区別されます。
- タグキーとタグ値にプレフィックス aws を含めることはできません。そのようなタグに対するオペレーションは許可されていません。
- タグキーの最大長は UTF-8 で 128 Unicode 文字です。タグキーを空または null にすることはできません。
- タグ値の最大長は UTF-8 で 256 Unicode 文字です。タグ値は空または null にすることができます。

AWS Glue 接続のタグ付けのサポート

リソースタグに基づ

き、CreateConnection、UpdateConnection、GetConnection、DeleteConnection のアクションの許可を制限することができます。これにより、データカタログから JDBC 接続情報を取得する必要がある JDBC データソースを持つ AWS Glue ジョブに対して、最小特権のアクセス制御を実装することができます。

使用例

タグ ["connection-category", "dev-test"] との AWS Glue 接続を作成する。

IAM ポリシーの GetConnection アクションのタグ条件を指定します。

```
{
  "Effect": "Allow",
  "Action": [
    "glue:GetConnection"
  ],
  "Resource": "*",
  "Condition": {
    "ForAnyValue:StringEquals": {
      "aws:ResourceTag/tagKey": "dev-test"
    }
  }
}
```

```
}  
}
```

例

次の例では、タグが割り当てられたジョブを作成します。

AWS CLI

```
aws glue create-job --name job-test-tags --role MyJobRole --command  
  Name=glueetl,ScriptLocation=S3://aws-glue-scripts//prod-job1  
  --tags key1=value1,key2=value2
```

AWS CloudFormation JSON

```
{  
  "Description": "AWS Glue Job Test Tags",  
  "Resources": {  
    "MyJobRole": {  
      "Type": "AWS::IAM::Role",  
      "Properties": {  
        "AssumeRolePolicyDocument": {  
          "Version": "2012-10-17",  
          "Statement": [  
            {  
              "Effect": "Allow",  
              "Principal": {  
                "Service": [  
                  "glue.amazonaws.com"  
                ]  
              },  
              "Action": [  
                "sts:AssumeRole"  
              ]  
            }  
          ]  
        },  
        "Path": "/",  
        "Policies": [  
          {  
            "PolicyName": "root",  
            "PolicyDocument": {  
              "Version": "2012-10-17",
```

```
        "Statement": [
          {
            "Effect": "Allow",
            "Action": "*",
            "Resource": "*"
          }
        ]
      }
    ]
  },
  "MyJob": {
    "Type": "AWS::Glue::Job",
    "Properties": {
      "Command": {
        "Name": "glueetl",
        "ScriptLocation": "s3://aws-glue-scripts//prod-job1"
      },
      "DefaultArguments": {
        "--job-bookmark-option": "job-bookmark-enable"
      },
      "ExecutionProperty": {
        "MaxConcurrentRuns": 2
      },
      "MaxRetries": 0,
      "Name": "cf-job1",
      "Role": {
        "Ref": "MyJobRole",
        "Tags": {
          "key1": "value1",
          "key2": "value2"
        }
      }
    }
  }
}
```

AWS CloudFormation YAML

Description: AWS Glue Job Test Tags
Resources:

```
MyJobRole:
  Type: AWS::IAM::Role
  Properties:
    AssumeRolePolicyDocument:
      Version: '2012-10-17'
      Statement:
        - Effect: Allow
          Principal:
            Service:
              - glue.amazonaws.com
          Action:
            - sts:AssumeRole
    Path: "/"
    Policies:
      - PolicyName: root
        PolicyDocument:
          Version: '2012-10-17'
          Statement:
            - Effect: Allow
              Action: "*"
              Resource: "*"

MyJob:
  Type: AWS::Glue::Job
  Properties:
    Command:
      Name: glueetl
      ScriptLocation: s3://aws-glue-scripts//prod-job1
    DefaultArguments:
      "--job-bookmark-option": job-bookmark-enable
    ExecutionProperty:
      MaxConcurrentRuns: 2
    MaxRetries: 0
    Name: cf-job1
    Role:
      Ref: MyJobRole
    Tags:
      key1: value1
      key2: value2
```

詳細については、「[AWSタグ付け戦略](#)」を参照してください。

タグを使用してアクセスをコントロールする方法については、「[AWS Glue での ABAC](#)」を参照してください。

CloudWatch Events を使用した AWS Glue の自動化

Amazon CloudWatch Events を使用して、AWS のサービスを自動化し、アプリケーションの可用性の問題やリソースの変更などのシステムイベントに自動的に対応できます。AWS のサービスからのイベントは、ほぼリアルタイムに CloudWatch Events に送信されます。簡単なルールを記述して、注目するイベントと、イベントがルールに一致した場合に自動的に実行するアクションを指定できます。自動的にトリガーできるオペレーションには、以下が含まれます。

- AWS Lambda 関数の呼び出し
- Amazon EC2 Run Command の呼び出し
- Amazon Kinesis Data Streams へのイベントの中継
- AWS Step Functions ステートマシンのアクティブ化
- Amazon SNS トピックまたは Amazon SQS キューの通知

AWS Glue で CloudWatch Events を使用する例をいくつか次に示します。

- ETL ジョブが成功したときの Lambda 関数のアクティブ化
- ETL ジョブが失敗したときの Amazon SNS トピックへの通知

次の CloudWatch Events は、AWS Glue によって生成されます。

- "detail-type":"Glue Job State Change" のイベントが SUCCEEDED、FAILED、TIMEOUT、STOPPED に対して生成されます。
- ジョブ遅延通知のしきい値を超えると、"detail-type":"Glue Job Run Status" のイベントが RUNNING、STARTING、STOPPING ジョブの実行に対して生成されます。これらのイベントを受信するには、ジョブ遅延通知のしきい値プロパティを設定する必要があります。

ジョブ遅延通知のしきい値を超えると、ジョブの実行ステータスごとに 1 つだけイベントが生成されます。

- "detail-type":"Glue Crawler State Change" のイベントが Started、Succeeded、Failed に対して生成されます。
- "detail-type":"Glue Data Catalog Database State Change" のイベントが、CreateDatabase、DeleteDatabase、CreateTable、DeleteTable、および BatchDeleteTable に対して生成されます。例えば、テーブルが作成または削除されると、CloudWatch Events に通知が送信されます。通知イベントの順序や有無に依存するプログラ

ムは作成できないことに注意してください。通知イベントは順番が違っていたり不足している可能性があります。イベントは、ベストエフォートベースで発生します。通知の詳細情報。

- `typeOfChange` には、API オペレーションの名前が含まれます。
- `databaseName` には影響を受けるデータベースの名前が含まれます。
- `changedTables` には、通知ごとに最大 100 の影響を受けるテーブルの名前が含まれます。テーブル名が長いと、複数の通知が作成されることがあります。
- `"detail-type": "Glue Data Catalog Table State Change"` のイベントが、`UpdateTable`、`CreatePartition`、`BatchCreatePartition`、`UpdatePartition`、`DeletePartition` および `BatchDeletePartition` に対して生成されます。例えば、テーブルまたはパーティションが更新されると、CloudWatch Events に通知が送信されます。通知イベントの順序や有無に依存するプログラムは作成できないことに注意してください。通知イベントは順番が違っていたり不足している可能性があります。イベントは、ベストエフォートベースで発生します。通知の詳細情報。
 - `typeOfChange` には、API オペレーションの名前が含まれます。
 - `databaseName` には影響を受けるリソースを含むデータベースの名前が含まれます。
 - `tableName` には影響を受けるテーブルの名前が含まれます。
 - `changedPartitions` は 1 つの通知で影響を受けるパーティションを最大 100 まで指定します。パーティション名が長いと、複数の通知が作成されることがあります。

たとえば、`Year` と `Month` の 2 つのパーティションキーがある場合、`"2018,01"`、`"2018,02"` は、パーティション `"Year=2018"` and `"Month=01"`、およびパーティション `"Year=2018"` and `"Month=02"` を変更します。

```
{
  "version": "0",
  "id": "abcdef00-1234-5678-9abc-def012345678",
  "detail-type": "Glue Data Catalog Table State Change",
  "source": "aws.glue",
  "account": "123456789012",
  "time": "2017-09-07T18:57:21Z",
  "region": "us-west-2",
  "resources": ["arn:aws:glue:us-west-2:123456789012:database/default/foo"],
  "detail": {
    "changedPartitions": [
      "2018,01",
      "2018,02"
    ],
    "databaseName": "default",
  }
}
```

```
"tableName": "foo",  
"typeOfChange": "BatchCreatePartition"  
}  
}
```

詳細については、「[Amazon CloudWatch Events ユーザーガイド](#)」を参照してください。AWS Glue 固有のイベントについては、「[AWS Glue イベント](#)」を参照してください。

AWS Glue リソースのモニタリング

AWS Glue には、想定外の過剰なプロビジョニングや請求額を増やすことを目的とした悪意のある行為からお客様を保護するためのサービス制限があります。制限はサービスの保護にもなります。AWS サービスクォータコンソールにログインすると、お客様は現在のリソース制限を確認し、(必要に応じて) 増加をリクエストできます。

AWS Glue では、サービスのリソース使用率を Amazon CloudWatch 内のパーセンテージで表示し、使用状況をモニタリングするように CloudWatch アラームを設定できます。Amazon CloudWatch は、Amazon インフラストラクチャで実行されている AWS リソースと顧客アプリケーションをモニタリングします。メトリクスは無料で使用できます。次のメトリクスがサポートされています。

- アカウントあたりのワークフローの数
- アカウントあたりのトリガー数
- アカウントあたりのジョブ数
- アカウントあたりの同時ジョブの実行数
- アカウントあたりのブループrintの数の数
- アカウントあたりのインタラクティブセッションの数

リソースメトリクスの設定と使用

この機能を使用するには、Amazon CloudWatch コンソールに移動してメトリクスを表示し、アラームを設定します。メトリクスは AWS/Glue 名前空間の下にあり、実際のリソース使用量をリソースクォータで割ったパーセンテージです。CloudWatch メトリクスはアカウントに配信され、費用はかかりません。例えば、10 件のワークフローを作成していて、サービスクォータで最大 200 件のワークフローを作成できる場合、使用率は $10/200 = 5\%$ となり、グラフでは 5 のデータポイントがパーセンテージとして表示されます。具体的には:

```
Namespace: AWS/Glue
```

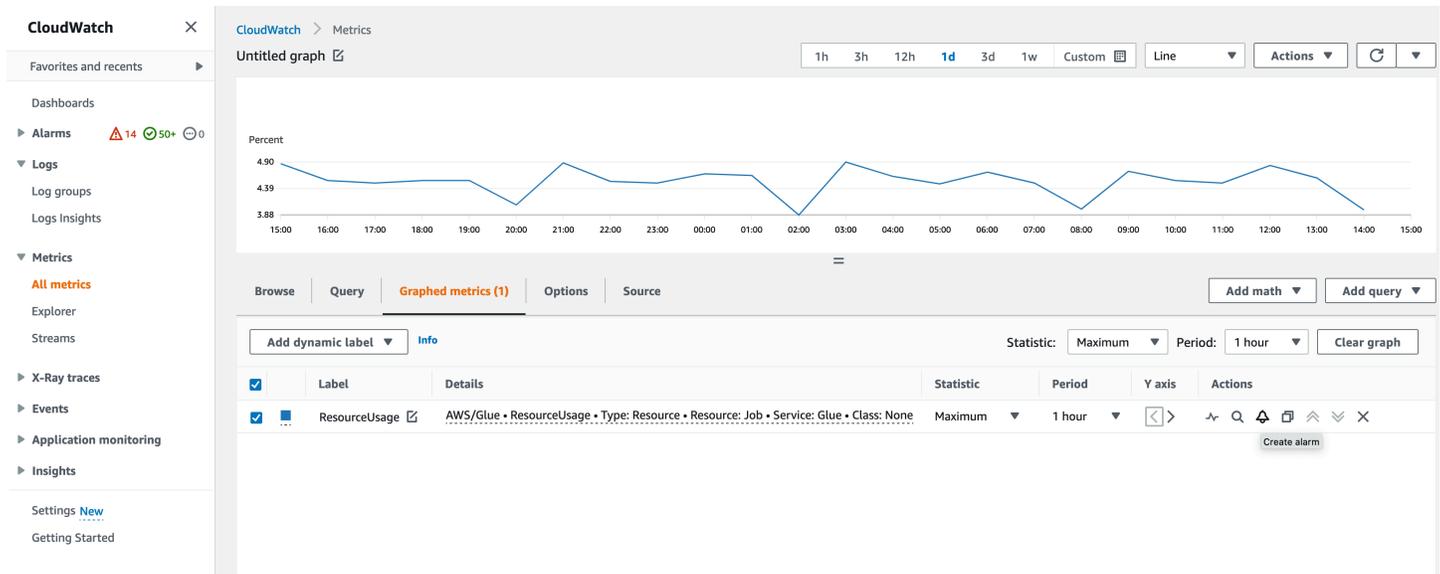
Metric name: ResourceUsage

Type: Resource

Resource: Workflow (or Trigger, Job, JobRun, Blueprint, InteractiveSession)

Service: Glue

Class: None



CloudWatch コンソールでクラスターメトリクスにアラームを作成するには:

1. メトリクスを見つけたら、[グラフ化したメトリクス] に移動します。
2. [アクション] の [アラームを作成] をクリックします。
3. 必要に応じてアラームを設定します。

リソース使用量が変化 (増加や減少) するたびにメトリクスを出力します。ただし、リソースの使用量が変わらない場合でも、メトリクスは 1 時間ごとに出力されるため、CloudWatch グラフは連続的になります。データポイントの欠落を防ぐため、1 時間未満の期間を設定することはお勧めしません。

また、以下の例のように、AWS CloudFormation を使用してアラームを設定することもできます。この例では、ワークフローリソースの使用率が 80% に達すると、アラームがトリガーされて既存の SNS トピックにメッセージが送信されます。サブスクライブすると通知を受け取ることができます。

```
{
  "Type": "AWS::CloudWatch::Alarm",
  "Properties": {
```

```
"AlarmName": "WorkflowUsageAlarm",
"ActionsEnabled": true,
"OKActions": [],
"AlarmActions": [
  "arn:aws:sns:af-south-1:085425700061:Default_CloudWatch_Alarms_Topic"
],
"InsufficientDataActions": [],
"MetricName": "ResourceUsage",
"Namespace": "AWS/Glue",
"Statistic": "Maximum",
"Dimensions": [{
  "Name": "Type",
  "Value": "Resource"
},
{
  "Name": "Resource",
  "Value": "Workflow"
},
{
  "Name": "Service",
  "Value": "Glue"
},
{
  "Name": "Class",
  "Value": "None"
}
],
"Period": 3600,
"EvaluationPeriods": 1,
"DatapointsToAlarm": 1,
"Threshold": 80,
"ComparisonOperator": "GreaterThanThreshold",
"TreatMissingData": "notBreaching"
}
}
```

AWS Glue による AWS CloudTrail API コールのログ記録

AWS Glue は AWS CloudTrail という、AWS Glue のユーザー、ロール、または AWS のサービスが実行したアクションを記録するサービスと統合しています。CloudTrail は、AWS Glue のすべての API コールをイベントとしてキャプチャします。キャプチャされたコールには、AWS Glue コンソールのコールと、AWS Glue API オペレーションへのコードのコールが含まれます。証跡を作成する場

合は、AWS Glue のイベントなど、Amazon S3 バケットへの CloudTrail イベントの継続的な配信を有効にすることができます。証跡を設定しない場合でも、CloudTrail コンソールの イベント履歴 で最新のイベントを表示できます。CloudTrail が収集した情報を使用して、AWS Glue に対して行われた要求、要求が行われた IP アドレス、要求を行った人、要求が行われた日時、および追加の詳細を判別できます。

CloudTrail の詳細については、「[AWS CloudTrail ユーザーガイド](#)」を参照してください。

CloudTrail の AWS Glue 情報

CloudTrailは、アカウントを作成するとAWSアカウントで有効になります。AWS Glue でアクティビティが発生すると、そのアクティビティは Event history イベント履歴で AWS のその他のサービスのイベントと共に CloudTrail イベントに記録されます。最近のイベントは、AWSアカウントで表示、検索、ダウンロードできます。詳細については、「[CloudTrail イベント履歴でのイベントの表示](#)」を参照してください。

AWSのイベントなど、AWS Glueアカウントのイベントの継続的なレコードについては、追跡を作成します。証跡により、CloudTrail はログファイルを Amazon S3 バケットに配信できます。デフォルトでは、コンソールで証跡を作成すると、すべての AWS リージョンに証跡が適用されます。証跡は、AWSパーティションのすべてのリージョンからのイベントをログに記録し、指定した Amazon S3 バケットにログファイルを配信します。さらに、CloudTrail ログで収集したイベントデータをより詳細に分析し、それに基づいて対応するため、他の AWS サービスを構成できます。詳細については、次を参照してください:

- [AWS アカウントの証跡の作成](#)
- 「[CloudTrail がサポートされているサービスと統合](#)」
- [CloudTrail 用 Amazon SNS 通知の構成](#)
- 「[複数のリージョンからCloudTrailログファイルを受け取る](#)」および「[複数のアカウントからCloudTrailログファイルを受け取る](#)」

すべての AWS Glue アクションは CloudTrail によってログに記録され、[AWS Glue API](#) に記録されます。例えば、CreateDatabase、CreateTable、CreateScript の各アクションを呼び出すと、CloudTrail ログファイルにエントリが生成されます。

各イベントまたはログエントリには、誰がリクエストを生成したかという情報が含まれます。同一性情報は次の判断に役立ちます。

- リクエストが、ルートと IAM ユーザー認証情報のどちらを使用して送信されたか。

- リクエストがロールまたはフェデレーションユーザーの一時的なセキュリティ認証情報を使用して行われたかどうか。
- リクエストが別の AWS サービスによって行われたかどうか。

詳細については、[CloudTrail userIdentity 要素](#)を参照してください。

ただし、CloudTrail は呼び出しに関するすべての情報は記録しません。たとえば、接続リクエストで使用される ConnectionProperties などの機密情報は記録せず、次の API によって返される応答の代わりに null を記録します。

BatchGetPartition	GetCrawlers	GetJobs	GetTable
CreateScript	GetCrawlerMetrics	GetJobRun	GetTables
GetCatalogImportStatus	GetDatabase	GetJobRuns	GetTableVersions
GetClassifier	GetDatabases	GetMapping	GetTrigger
GetClassifiers	GetDataflowGraph	GetObjects	GetTriggers
GetConnection	GetDevEndpoint	GetPartition	GetUserDefinedFunction
GetConnections	GetDevEndpoints	GetPartitions	GetUserDefinedFunctions
GetCrawler	GetJob	GetPlan	

AWS Glue ログファイルエントリの理解

追跡は、指定した Amazon S3 バケットにイベントをログファイルとして配信するように設定できます。CloudTrail のログファイルは、単一か複数のログエントリを含みます。イベントは任意ソースからの単一リクエストを表し、リクエストされたアクション、アクションの日時、リクエストパラメータなどの情報を含みます。CloudTrail ログファイルは、パブリック API コールの順序付けられたスタックトレースではないため、特定の順序では表示されません。

以下の例は、DeleteCrawler アクションを示す CloudTrail ログエントリです。

```
{
  "eventVersion": "1.05",
  "userIdentity": {
    "type": "IAMUser",
    "principalId": "AKIAIOSFODNN7EXAMPLE",
    "arn": "arn:aws:iam::123456789012:user/johndoe",
    "accountId": "123456789012",
    "accessKeyId": "AKIAIOSFODNN7EXAMPLE",
    "userName": "johndoe"
  },
  "eventTime": "2017-10-11T22:29:49Z",
```

```
"eventSource": "glue.amazonaws.com",
"eventName": "DeleteCrawler",
"awsRegion": "us-east-1",
"sourceIPAddress": "72.21.198.64",
"userAgent": "aws-cli/1.11.148 Python/3.6.1 Darwin/16.7.0 botocore/1.7.6",
"requestParameters": {
  "name": "tes-alpha"
},
"responseElements": null,
"requestID": "b16f4050-aed3-11e7-b0b3-75564a46954f",
"eventID": "e73dd117-cfd1-47d1-9e2f-d1271cad838c",
"eventType": "AwsApiCall",
"recipientAccountId": "123456789012"
}
```

この例では、CreateConnection アクションを表す CloudTrail ログエントリを示しています。

```
{
  "eventVersion": "1.05",
  "userIdentity": {
    "type": "IAMUser",
    "principalId": "AKIAIOSFODNN7EXAMPLE",
    "arn": "arn:aws:iam::123456789012:user/johndoe",
    "accountId": "123456789012",
    "accessKeyId": "AKIAIOSFODNN7EXAMPLE",
    "userName": "johndoe"
  },
  "eventTime": "2017-10-13T00:19:19Z",
  "eventSource": "glue.amazonaws.com",
  "eventName": "CreateConnection",
  "awsRegion": "us-east-1",
  "sourceIPAddress": "72.21.198.66",
  "userAgent": "aws-cli/1.11.148 Python/3.6.1 Darwin/16.7.0 botocore/1.7.6",
  "requestParameters": {
    "connectionInput": {
      "name": "test-connection-alpha",
      "connectionType": "JDBC",
      "physicalConnectionRequirements": {
        "subnetId": "subnet-323232",
        "availabilityZone": "us-east-1a",
        "securityGroupIdList": [
          "sg-12121212"
        ]
      }
    }
  }
}
```

```
    }  
  }  
},  
"responseElements": null,  
"requestID": "27136ebc-afac-11e7-a7d6-ab217e5c3f19",  
"eventID": "e8b3baeb-c511-4597-880f-c16210c60a4a",  
"eventType": "AwsApiCall",  
"recipientAccountId": "123456789012"  
}
```

AWS Glue ジョブ実行ステータス

AWS Glue 抽出、変換、ロード (ETL) ジョブのステータスは、実行中または停止後に表示できます。ステータスは AWS Glue コンソール、AWS Command Line Interface (AWS CLI)、または AWS Glue API の [GetJobRun アクション](#) で表示できます。

ジョブ実行ステータスに

は、STARTING、RUNNING、STOPPING、STOPPED、SUCCEEDED、FAILED、ERROR、WAITING、TIMEOUT があります。

次の表に、ジョブの異常終了を示すステータスを一覧表示します。

ジョブ実行ステータス	説明
FAILED	ジョブが許容される最大同時実行数を超過しているか、不明な終了コードで終了しました。
ERROR	ワークフロー、スケジュールトリガー、またはイベントトリガーが削除されたジョブを実行しようとした。
TIMEOUT	ジョブの実行時間が、指定されたタイムアウト値を超えました。

WAITING ステータスは、ジョブの実行がリソース待ちであることを示しています。次の表には、さまざまなクラスのジョブの待機動作が記載されています。

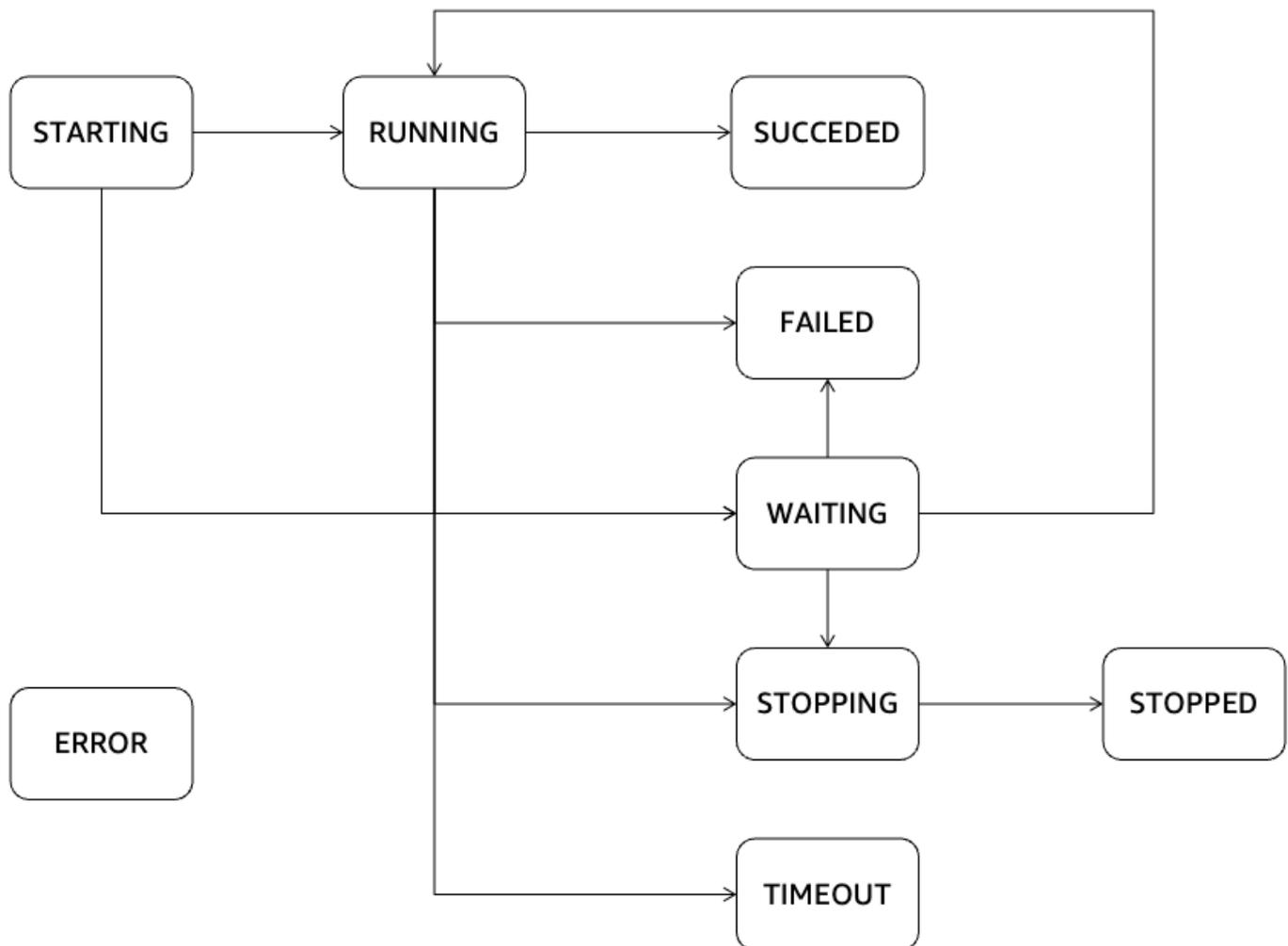
ジョブタイプ	Behavior
Spark ジョブ (標準)	<p>maxRetries 設定に基づいて再試行するように設定されていないジョブは、[待機中] 状態になる可能性があります。サービスが実行を開始するのに十分なリソースを取得できない場合、新しいジョブの実行は [待機中] 状態になります。これは、アカウントのサービスクォータまたはリージョンのキャパシティ制限により発生し、次のいずれかのエラーケースにつながる可能性があります。</p> <ul style="list-style-type: none">• アカウントあたりの同時実行ジョブの最大実行数を超過しました• ジョブごとの最大同時ジョブ実行数を超過した (アカウントレベルのサービスクォータと、MaxConcurrentRuns でジョブについて指定した制限を含みます)• 最大同時コンピューティング (DPU 使用量) を超過した• リソースを利用できません。 <p>AWS Glue サービスクォータの詳細については、「AWS Glue エンドポイントとクォータ」を参照してください。AWS Glue がリソースを待機する時間は、状況によって異なる場合があります。ジョブは、リソースの取得を試みる際に、非終了ステータス間で遷移することがあります。最終的に、リソースを取得できない場合、ジョブは FAILED に移行します。AWSGlue は、最大 15 分間または 10 回の試行のうち、いずれかの要件が満たされるまで再試行します。</p>
Spark ジョブ (Flex)	サービスが実行を開始するのに十分なリソースを取得できない場合、新しいジョブの実行

ジョブタイプ	Behavior
	<p>は WAITING 状態になり、実行の開始が遅れます。実行は最大 20 分間にわたって WAITING 状態になります (タイムアウト時間はサービスが制御)。15 分経過すると、サービスは強制起動を試み、使用可能な容量に応じて実行が開始するか失敗し、適切なエラーメッセージが表示されます。</p>

Python シェルジョブ

Spark を使用する標準ジョブと同じ動作。

次の状態図は、AWS Glue ジョブのライフサイクルを通じて予想される状態遷移の概要を示しています。この情報はすべてのジョブタイプに適用されます。



AWS Glue ストリーミング

AWS Glue のコンポーネントである Streaming を使用すると AWS Glue、ストリーミングデータをほぼリアルタイムで効率的に処理できるため、データの取り込み、処理、機械学習などの重要なタスクを実行できます。Apache Spark Streaming フレームワークを使用して、AWS Glue ストリーミングデータを大規模に処理できるサーバーレスサービスを提供します。は、サーバーレスインフラストラクチャ、自動スケーリング、ビジュアルジョブ開発、ストリーミングジョブ用のインスタントノートブック、その他のパフォーマンス向上など、Apache Spark 上にさまざまな最適化 AWS Glue を提供します。

ストリーミングのユースケース

AWS Glue ストリーミングの一般的なユースケースには、次のようなものがあります。

N ear-real-time データ処理: AWS Glue ストリーミングを使用すると、組織はストリーミングデータをほぼリアルタイムで処理できるため、インサイトを導き出し、最新情報に基づいてタイムリーな意思決定を行うことができます。

不正検出: AWS Glue ストリーミングデータをリアルタイムで分析するために Streaming を利用するため、クレジットカード詐欺、ネットワーク侵入、オンライン詐欺などの不正行為を検出するために役立ちます。受信データを継続的に処理して分析することで、疑わしいパターンや異常を迅速に特定できます。

ソーシャルメディア分析: AWS Glue ストリーミングは、ツイート、投稿、コメントなどのリアルタイムのソーシャルメディアデータを処理できるため、組織は傾向のモニタリング、感情分析、ブランド評価のリアルタイム管理を行うことができます。

Internet of Things (IoT) 分析: AWS Glue ストリーミングは、IoT デバイス、センサー、およびコネクテッドマシンによって生成されたデータの高速ストリームの処理と分析に適しています。これにより、リアルタイムの監視、異常検知、予知保全、およびその他の IoT 分析のユースケースが可能になります。

クリックストリーム分析: AWS Glue ストリーミングは、ウェブサイトまたはモバイルアプリケーションからのリアルタイムのクリックストリームデータを処理および分析できます。これにより、企業はユーザーの行動に関する洞察を得たり、ユーザー体験をパーソナライズしたり、リアルタイムのクリックストリームデータに基づいてマーケティングキャンペーンを最適化したりすることができます。

ログのモニタリングと分析: AWS Glue ストリーミングは、サーバー、アプリケーション、またはネットワークデバイスからのログデータをリアルタイムで継続的に処理および分析できます。これは、異常の検出、問題のトラブルシューティング、システムの状態とパフォーマンスの監視に役立ちます。

レコメンデーションシステム: AWS Glue ストリーミングは、ユーザーアクティビティデータをリアルタイムで処理し、レコメンデーションモデルを動的に更新できます。これにより、ユーザーの行動や好みに基づいた、パーソナライズされたリアルタイムのレコメンデーションが可能になります。

これらは、AWS Glue ストリーミングを適用できるさまざまなユースケースの例です。AWS エコシステムやマネージドサービスとの統合により、クラウドでのリアルタイムのストリーム処理と分析に便利な選択肢となります。

AWS Glue ストリーミングを使用する利点は何ですか？

AWS Glue ストリーミングを使用する利点は次のとおりです。

- **サーバーレス**: AWS Glue ストリーミングはサーバーレスであるため、インフラストラクチャを管理する必要はありません。これにより、運用上のオーバーヘッドが軽減され、ユーザーは、インフラストラクチャ管理ではなくデータ処理および分析タスクに集中できます。
- **自動スケーリング**: AWS Glue ストリーミングは、ワークロードに基づいて処理容量を動的に調整する自動スケーリング機能を提供します。データ量の変動に合わせて自動的にスケールアウトまたはスケールインを行い、最適なパフォーマンスとリソース使用率を確保します。
- **ビジュアル開発**: ストリーミングジョブの開発は複雑になる可能性があります。AWS Glue ストリーミングは、ビジュアルオーサリングツールである AWS Glue Studio を提供することで、この課題に対処します。AWS Glue Studio は、ストリーミングワークフローの作成プロセスを簡素化し、デベロッパーがストリーミングアプリケーションを視覚的に設計および管理できるようにし、学習曲線を短縮し、生産性を向上させます。
- **コスト効率の高い**: サーバーレスサービスである AWS Glue ストリーミングは、インフラストラクチャのプロビジョニングと保守が不要になり、コスト効率が向上します。ユーザーへの請求は、ストリーミングジョブの実行中に消費されたリソースに基づいて行われるため、実際の使用量に基づくコストの最適化とスケーリングが可能になります。
- **複雑なワークロードを処理します**: AWS Glue ストリーミングは、複雑なストリーミングワークロードを処理するように設計されています。大量のリアルタイムデータの処理と分析、高度な変換のサポート、他の AWS サービスとの統合が可能で、高度なストリーミングデータパイプラインと分析ワークフローが可能になります。

- **ロックインなし**：AWS Glue ストリーミングは柔軟性を提供し、ベンダーのロックインを回避します。ユーザーは、より広範な AWS エコシステムの一部として AWS Glue ストリーミングを活用し、他の AWS サービスとシームレスに統合できます。これにより、特定のテクノロジーやプラットフォームに縛られることなく、既存のデータソース、アプリケーション、サービスと簡単に統合できます。

AWS Glue ストリーミングを使用するタイミング

ストリーミングのユースケースに関して言えば、多くの選択肢があります。以下のシナリオでは AWS Glue ストリーミングをお勧めします。

1. バッチ処理に AWS Glue または Spark を既に使用している場合は、AWS Glue ストリーミングが最適な選択肢です。新しい言語やフレームワークを学習しなくても、ストリーミングジョブの構築にシームレスに移行できます。既存の知識とインフラストラクチャを活用して、AWS Glue Streaming はジョブ開発プロセスを簡素化し、データ処理機能をリアルタイムストリーミングシナリオに簡単に拡張できます。
2. バッチ、ストリーミング、イベント駆動型のワークロードを処理するために統合サービスまたは製品が必要な場合は、AWS Glue ストリーミングがソリューションとなります。AWS Glue Streaming を使用すると、データ処理のニーズを 1 つのフレームワークに統合できるため、複数のシステムを管理する複雑さがなくなります。これにより、さまざまなワークロードタイプ間の一貫性と互換性を確保しながら、多様なデータワークフローを効率的に開発および保守できます。
3. AWS Glue ストリーミングは、ストリームやリレーショナルデータベース間の結合など、非常に大きなストリーミングデータボリュームや複雑な変換を含むシナリオに適しています。大量のデータストリームを効率的に処理して分析できるため、要求の厳しいワークロードにも簡単に取り組むことができます。高速データ取り込みでも複雑なデータ操作でも、AWS Glue ストリーミングのスケーラビリティと高度な処理機能により、最適なパフォーマンスと正確な結果が得られます。
4. ストリーミングジョブの構築に視覚的なアプローチを希望する場合は、ストリーミングアプリケーションを視覚的に設計および管理できる AWS Glue Studio AWS Glue を提供し、開発プロセスを簡素化します。この直感的なインターフェイスにより、開発者はビジュアルインターフェイスを使用してストリーミングワークフローを作成、設定、監視できるため、習得時間が短縮され、生産性が向上します。
5. AWS Glue ストリーミングは、厳格な SLAs (サービスレベルアグリーメント) が 10 秒を超える near-real-time ユースケースに最適です。

6. Apache Iceberg、Apache Hudi、または Delta Lake を使用してトランザクションデータレイクを構築する場合、AWS Glue Streaming はこれらのオープンテーブル形式をネイティブにサポートします。このシームレスな統合により、これらのトランザクションデータレイクからストリーミングデータを直接処理できるようになり、データ整合性、完全性、互換性が確保されます。
7. さまざまなデータターゲットのストリーミングデータを取り込む必要がある場合：AWS Glue Streaming は、Amazon Redshift、Amazon RDS、Amazon Aurora、Oracle、SQL Server などのさまざまなデータターゲットにネイティブターゲットを提供します。

サポートされているデータソース

AWS Glue ストリーミングでは、次のデータソースがサポートされています。

- Amazon Kinesis
- Amazon MSK (Managed Streaming for Apache Kafka)
- セルフマネージド Apache Kafka

サポートされるデータターゲット

AWS Glue ストリーミングは、次のようなさまざまなデータターゲットをサポートします。

- Data Catalog でサポートされる AWS Glue データターゲット
- Amazon S3
- Amazon Redshift
- MySQL
- PostgreSQL
- Oracle
- Microsoft SQL Server
- Snowflake
- JDBC を使用して接続できるすべてのデータベース
- Apache Iceberg、Delta、および Apache Hudi
- AWS Glue Marketplace コネクタ

チュートリアル: AWS Glue Studio を使用して最初のストリーミングワークロードを構築する

このチュートリアルでは、AWS Glue Studio を使用してストリーミングジョブを作成する方法について説明します。AWS GlueStudio は、AWS Glue ジョブを作成するためのビジュアルインターフェイスです。

Amazon Kinesis Data Streams、Apache Kafka、および Amazon Managed Streaming for Apache Kafka (Amazon MSK) で、ストリーミングソースからのデータを消費し、連続的に実行するストリーミング抽出、変換、ロード (ETL) ジョブを作成できます。

前提条件

このチュートリアルを実行するには、AWS Glue、Amazon Kinesis、Amazon S3、Amazon Athena、AWS CloudFormation、AWS Lambda、Amazon Cognito を使用する AWS コンソール権限を持つユーザーが必要です。

Amazon Kinesis からストリーミングデータを消費する

トピック

- [Kinesis Data Generator でモックデータを生成する](#)
- [AWS Glue Studio で AWS Glue Streaming ジョブを作成する](#)
- [変換を実行し、変換した結果を Amazon S3 に保存する](#)

Kinesis Data Generator でモックデータを生成する

Kinesis Data Generator (KDG) を使用して JSON 形式のサンプルデータを合成的に生成できます。完全な手順と詳細については、[ツールのドキュメント](#)を参照してください。

1. 開始するに

は、

をクリックして、ご使用の AWS 環境で AWS CloudFormation テンプレートを実行します。

Note

Kinesis Data Generator の Amazon Cognito ユーザーなどのリソースが AWS アカウントに既に存在しているために、CloudFormation テンプレートの実行が失敗する場合があります。

ます。これは、別のチュートリアルやブログで既に設定していたことが原因である可能性があります。これに対処するには、新しい AWS アカウントでテンプレートの実行をやり直すか、別の AWS リージョンを試してみてください。これらの方法により、既存のリソースと競合することなくチュートリアルを実行できるようになります。

このテンプレートにより、Kinesis データストリームと Kinesis Data Generator アカウントがプロビジョニングされます。また、データを保持する Amazon S3 バケットと、このチュートリアルに必要な権限を持つ Glue サービスロールが作成されます。

2. KDG が認証に使用する [ユーザー名] と [パスワード] を入力します。後で使用するために、ユーザー名とパスワードをメモしておきます。
3. 最後のステップまで [次へ] を選択していきます。IAM リソースの作成を承認します。パスワードが最小要件を満たしていないなど、画面上部にエラーがないか確認し、テンプレートをデプロイします。
4. スタックの [出力] タブに移動します。テンプレートがデプロイされると、生成されたプロパティ `KinesisDataGeneratorUrl` が表示されます。その URL をクリックします。
5. メモしておいた [ユーザー名] と [パスワード] を入力します。
6. 使用しているリージョンを選択し、Kinesis ストリーム `GlueStreamTest-{AWS::AccountId}` を選択します。
7. 以下のテンプレートを入力します。

```
{
  "ventilatorid": {{random.number(100)}},
  "eventtime": "{{date.now("YYYY-MM-DD HH:mm:ss")}}",
  "serialnumber": "{{random.uuid}}",
  "pressurecontrol": {{random.number(
    {
      "min":5,
      "max":30
    }
  )}},
  "o2stats": {{random.number(
    {
      "min":92,
      "max":98
    }
  )}},
  "minutevolume": {{random.number(
```

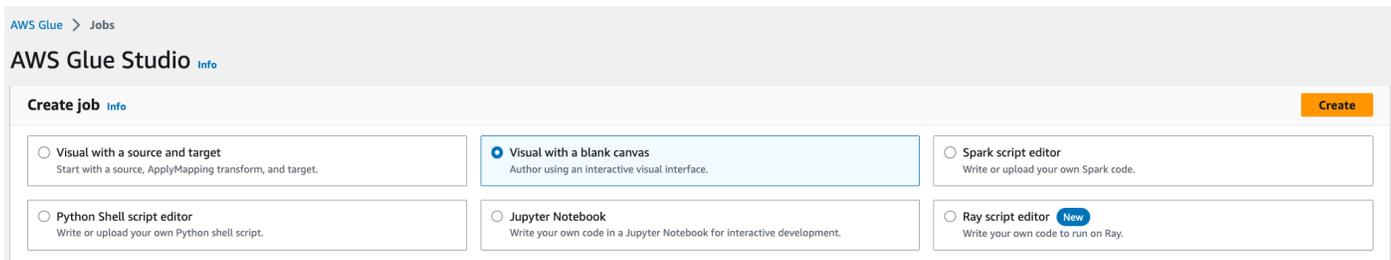
```
{
  "min":5,
  "max":8
}
}},
"manufacturer": "{random.arrayElement(
  ["3M", "GE","Vyaire", "Getinge"]
)}"
}
```

[テストテンプレート] を使用してモックデータを表示し、[データを送信する] を使用してモックデータを Kinesis に取り込むことができるようになりました。

8. [データを送信する] をクリックして、5~10K のレコードを Kinesis に生成します。

AWS Glue Studio で AWS Glue Streaming ジョブを作成する

1. 同じリージョンのコンソールで AWS Glue に移動します。
2. 左側のナビゲーションバーの [データ統合と ETL] にある [ETL ジョブ] を選択します。
3. [空白のキャンバスのビジュアル] を使用して AWS Glue ジョブを作成します。



4. [ジョブの詳細] タブに移動します。
5. AWS Glue のジョブ名には DemoStreamingJob を入力します。
6. [IAM ロール] では、CloudFormation テンプレートによってプロビジョニングされたロール glue-tutorial-role-\${AWS::AccountId} を選択します。
7. [Glue バージョン] では、Glue 3.0 を選択します。その他のオプションはすべて、デフォルト設定のままにします。

Basic properties [Info](#)**Name****Description - optional**

Descriptions can be up to 2048 characters long.

IAM Role

Role assumed by the job with permission to access your data stores. Ensure that this role has permission to your Amazon S3 sources, targets, temporary directory, scripts, and any libraries used by the job.

  **Type**

The type of ETL job. This is set automatically based on the types of data sources you have selected.

Glue version [Info](#) **Language** **Worker type**

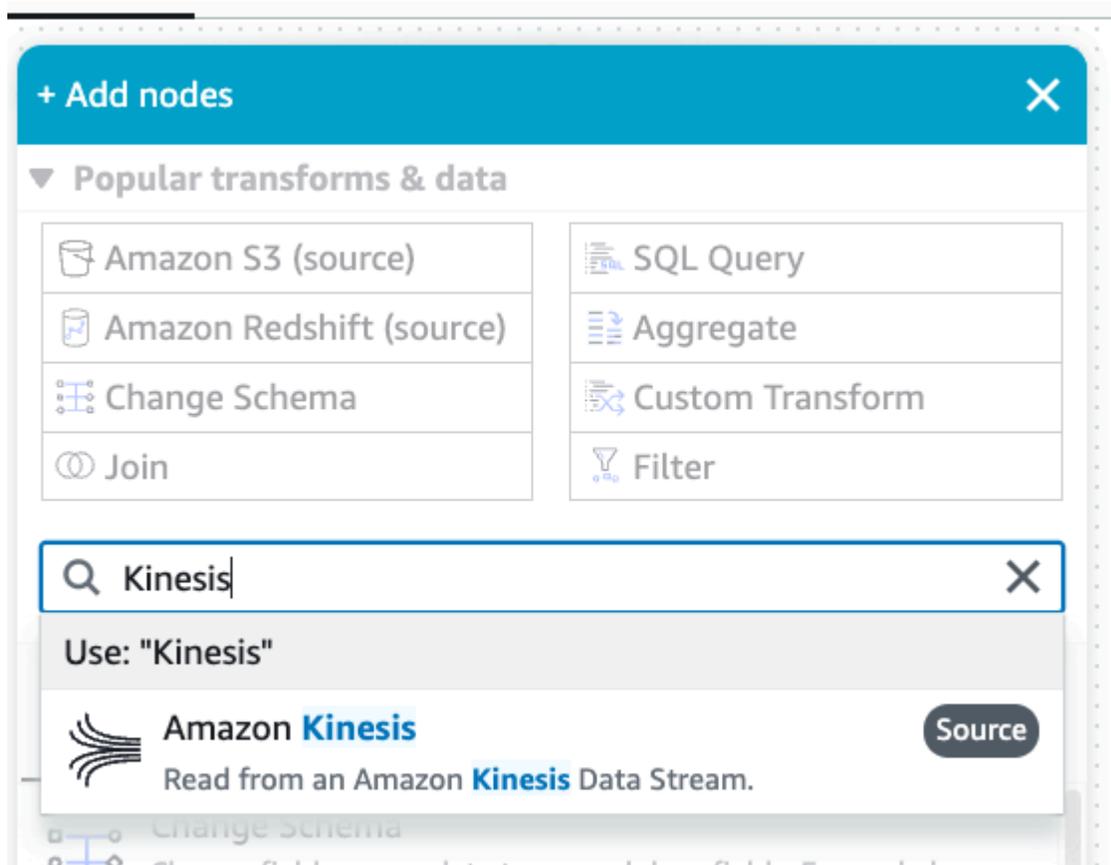
Set the type of predefined worker that is allowed when a job runs.

 **Automatically scale the number of workers**

- AWS Glue will optimize costs and resource usage by dynamically scaling the number of workers up and down throughout the job run. Requires Glue 3.0 or later.

8. [ビジュアル] タブに移動します。

9. プラスアイコンをクリックします。検索バーに「Kinesis」と入力します。[Amazon Kinesis] データソースを選択します。



10[データソースのプロパティ - Kinesis ストリーム] タブの [Amazon Kinesis ソース] で [ストリームの詳細] を選択します。

11[データストリームの場所] で [自分のアカウントにストリームを配置する] を選択します。

12.使用しているリージョンを選択します。

13.GlueStreamTest-{AWS::AccountId} ストリームを選択します。

14.他の設定はすべて、デフォルトのままにします。

Data source properties - Kinesis Stream | Output schema | Data preview 

Name
Amazon Kinesis

Amazon Kinesis Source [Info](#)

Stream details
 Data Catalog table

Location of data stream
 Stream is located in my account
 Stream is located in another account

Region
US East (Ohio) us-east-2

Stream name [Info](#)
GlueStreamTest- 

Data format
JSON

Starting position
Select the position where the job will start reading from the input stream.
Earliest
Start reading from the oldest available record in the stream.

Window size [Info](#)
Enter the time in seconds spent between batch calls.
100

15[データのプレビュー] タブに移動します。

16[データプレビューセッションを開始] をクリックすると、KDG が生成したモックデータがプレビューされます。AWS Glue Streaming ジョブのために前もって作成した Glue サービスロールを選択します。

プレビューデータが表示されるまで 30~60 秒かかります。[表示するデータがありません] と表示された場合は、歯車アイコンをクリックして、[サンプリングする行数] を 100 に変更します。

サンプルデータは以下のように表示されます。

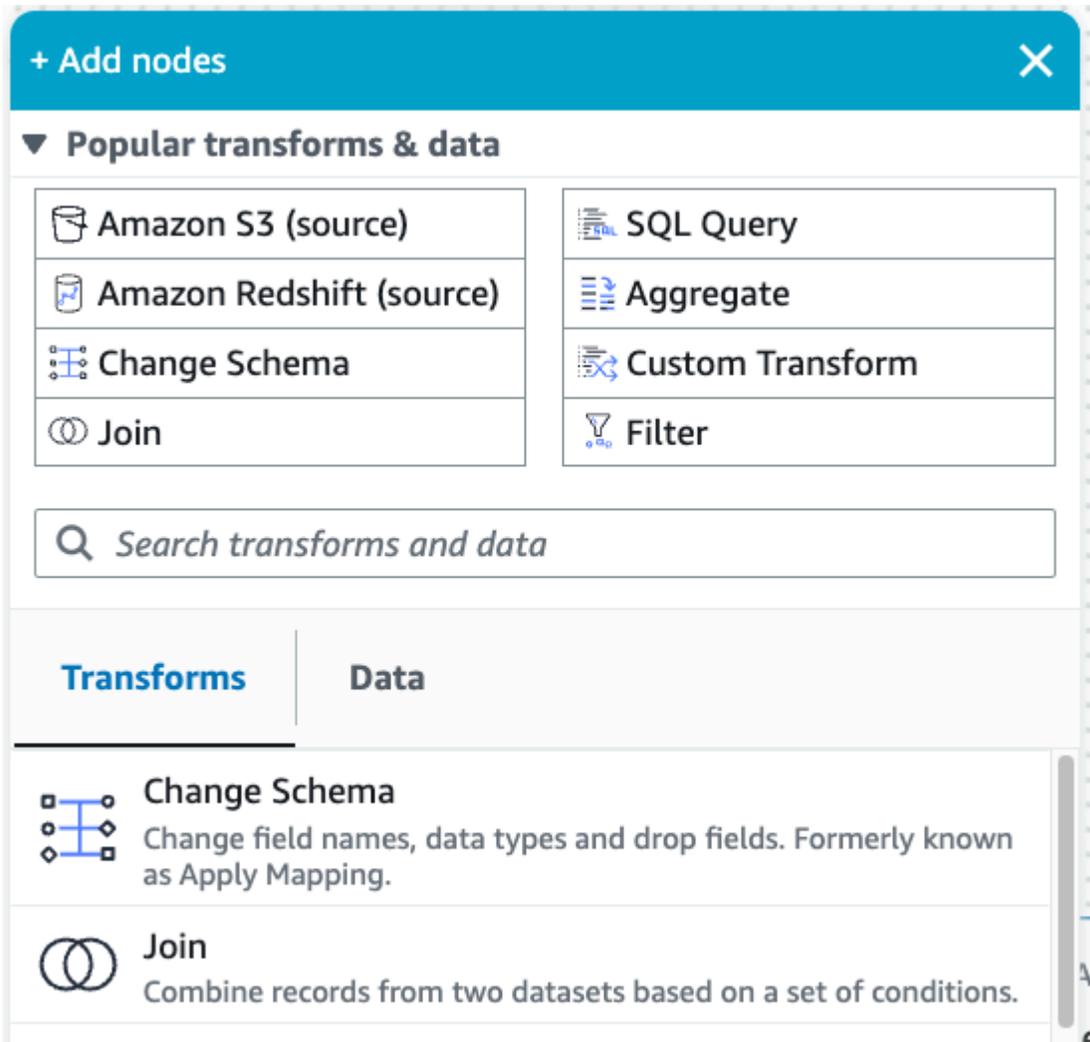
Data source properties - Kinesis Stream		Output schema	Data preview				
Data preview (100) Info		Previewing 7 of 7 fields					
<input type="text" value="Filter sample dataset"/>							
eventtime	manufacturer	minutevolume	o2stats	pressurecontrol	serialnumber	ventilatorid	
2023-06-26 14:25:37	Vyaire	5	95	7	9e79ae66-33a7-48e5-ab78-a61271199d5d	92	
2023-06-26 14:25:37	3M	5	98	17	cfb845ca-b513-4c27-9543-74dd222fc537	10	
2023-06-26 14:25:37	GE	8	98	23	90ba966c-6676-4567-a584-e267e714e57d	37	
2023-06-26 14:25:37	Vyaire	8	92	16	77f78f41-be24-47dc-b25c-05428bd76a0b	56	
2023-06-26 14:25:37	Getinge	6	92	23	ddf7b9e1-d0f7-4381-8aea-06a934583f5c	28	
2023-06-26 14:25:37	Getinge	5	92	6	c3ca9991-9b97-43e7-a866-59acbc6c5b17	84	
2023-06-26 14:25:37	3M	8	98	21	93c49e41-868b-4b5b-b725-06b4b1fb0a09	68	
2023-06-26 14:25:37	Vyaire	8	92	18	e46abe8d-b02f-43e6-91bf-c4700719f846	10	
2023-06-26 14:25:37	Vyaire	8	93	16	b3946e38-6292-4afd-8695-ada5cc09d0dd	15	
2023-06-26 14:25:37	GE	8	93	10	e3f7390d-1e68-4def-9dae-5c98b1d85d9d	3	
2023-06-26 14:25:37	Vyaire	8	98	17	a3917233-fe7f-4105-8728-779bd7ab1379	8	
2023-06-26 14:25:37	Getinge	8	98	16	06a8e8ff-cae4-4438-9714-33324f1524c9	93	
2023-06-26 14:25:37	Getinge	6	96	14	7af06237-bddf-4615-b9ac-05d05d484ba0	13	
2023-06-26 14:25:37	3M	8	93	8	bf9985f6-04b8-442b-b7f9-24b1db6b5a37	81	
2023-06-26 14:25:37	Getinge	6	97	28	e67f4220-3070-4951-b4e0-c86b7489de10	19	
2023-06-26 14:25:37	3M	6	92	15	77954206-535e-4ef8-a1fe-0da5ece049a6	31	
2023-06-26 14:25:37	Vyaire	7	94	25	81303a43-6206-46cb-851f-fc3986491bf9	32	

推測されたスキーマは [出力スキーマ] タブでも確認できます。

Data source properties - Kinesis Stream		Output schema	Data preview
Schema Info			
Key			Data type
eventtime			string
manufacturer			string
minutevolume			long
o2stats			long
pressurecontrol			long
serialnumber			string
ventilatorid			long

変換を実行し、変換した結果を Amazon S3 に保存する

1. ソースノードを選択した状態で、左上のプラスアイコンをクリックして [変換] ステップを追加します。
2. [スキーマ変更] ステップを選択します。



3. このステップでは、フィールドの名前を変更し、フィールドのデータ型を変換することができます。o2stats 列の名前を OxygenSaturation に変更し、すべての long データ型を int に変換します。

Transform
Output schema
Data preview

Name

Change Schema

Node parents
Choose which nodes will provide inputs for this one.

Choose one or more parent node

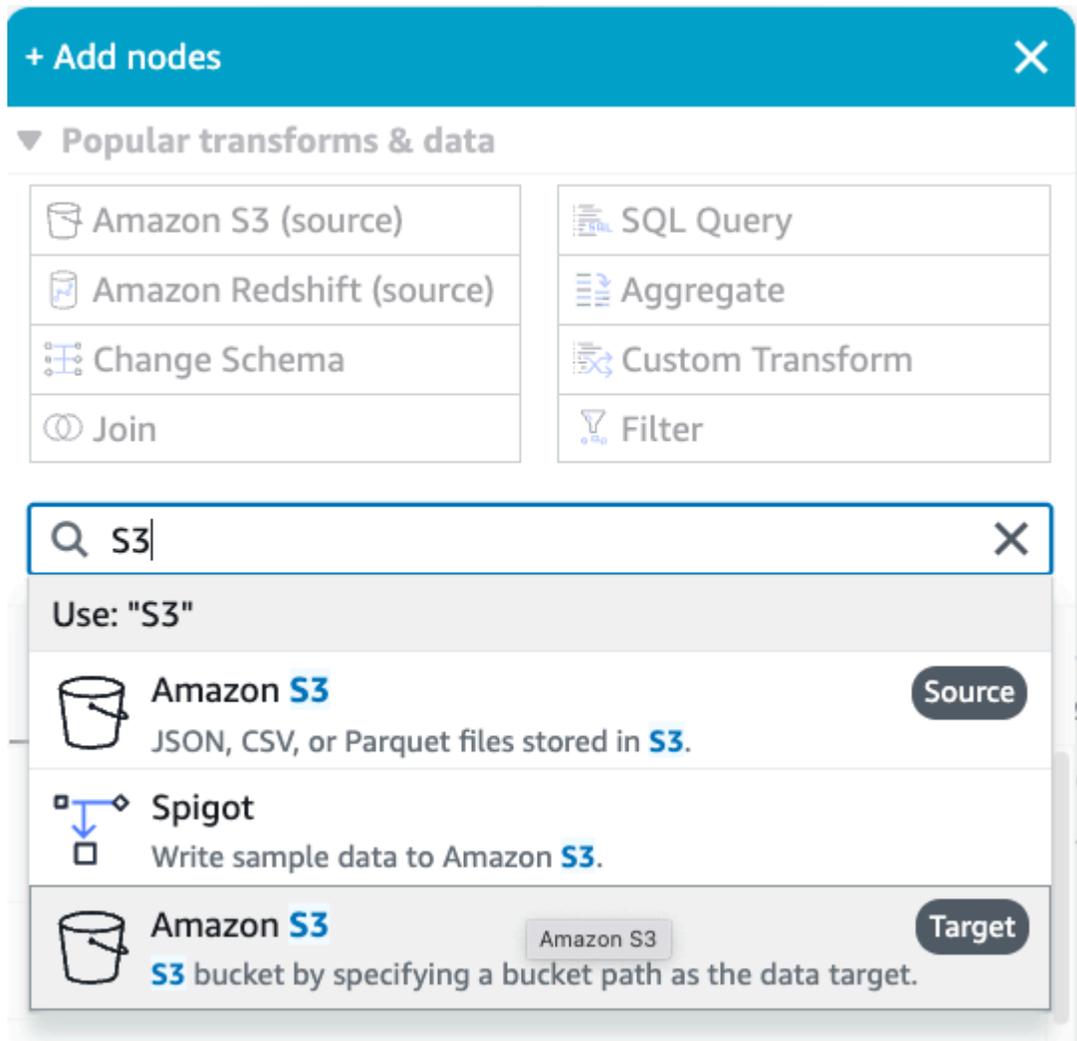
Amazon Kinesis ✕

Kinesis - DataSource

Change Schema (Apply mapping)

Source key	Target key	Data type	Drop
eventtime	<input type="text" value="eventtime"/>	string ▼	<input type="checkbox"/>
manufacturer	<input type="text" value="manufacturer"/>	string ▼	<input type="checkbox"/>
minutevolume	<input type="text" value="minutevolume"/>	int ▼	<input type="checkbox"/>
o2stats	<input type="text" value="OxygenSaturation"/>	int ▼	<input type="checkbox"/>
pressurecontrol	<input type="text" value="pressurecontrol"/>	int ▼	<input type="checkbox"/>
serialnumber	<input type="text" value="serialnumber"/>	string ▼	<input type="checkbox"/>
ventilatorid	<input type="text" value="ventilatorid"/>	int ▼	<input type="checkbox"/>

4. プラスアイコンをクリックして [Amazon S3] ターゲットを追加します。検索ボックスに「S3」と入力し、[Amazon S3 - ターゲット] 変換ステップを選択します。



5. ターゲットファイル形式として [Parquet] を選択します。
6. 圧縮タイプとして [Snappy] を選択します。
7. CloudFormation テンプレート streaming-tutorial-s3-target-{AWS::AccountId} によって作成された [S3 ターゲットの場所] を入力します。
8. [データカタログにテーブルを作成し、それ以降の実行時にスキーマを更新して新しいパーティションを追加する] を選択します。
9. Amazon S3 ターゲットテーブルのスキーマを保存するターゲットの [データベース] と [テーブル名] を入力します。

Name

Amazon S3

Node parents

Choose which nodes will provide inputs for this one.

Choose one or more parent node

Change Schema ✕
ApplyMapping - Transform

Format

Parquet

Compression Type

Snappy

S3 Target Location

Choose an S3 location in the format s3://bucket/prefix/object/ with a trailing slash (/).

s3://



View

Browse S3

Data Catalog update options [Info](#)

Choose how you want to update the Data Catalog table's schema and partitions. These options will only apply if the Data Catalog table is an S3 backed source.

- Do not update the Data Catalog
- Create a table in the Data Catalog and on subsequent runs, update the schema and add new partitions
- Create a table in the Data Catalog and on subsequent runs, keep existing schema and add new partitions

Database

Choose the database from the AWS Glue Data Catalog.

demo

**▶ Use runtime parameters**

Table name

Enter a table name for the AWS Glue Data Catalog.

demo_stream_transform_result

10[スクリプト] タブをクリックすると、生成されたコードが表示されます。

11.右上の [保存] をクリックして ETL コードを保存し、[実行] をクリックして AWS Glue Streaming ジョブを開始します。

[実行ステータス] は [実行] タブで確認できます。ジョブを 3~5 分間実行してから、ジョブを停止します。

Visual	Script	Job details	Runs	Data quality <small>New</small>	Schedules	Version Control
Job runs (1/1) Info						
<input type="text" value="Filter job runs by property"/>						
Run status	Retry	Start time	End time	Duration		
● ● Running	0	06/26/2023 15:58:05	-	35 s		

12 Amazon Athena で作成された新しいテーブルを確認します。

● Query 9

```
1 select * from "demo_stream_transform_result"
```

SQL Ln 1, Col 45

Run again
Explain
Cancel
Clear
Create

● Reuse query results up to 60 minutes ago

Query results | Query stats

● Completed Time in queue: 137 ms Run time: 894 ms Data scanned: 91.59 KB

Results (2,200)
Copy
Download results

#	eventtime	manufacturer	minutevolume	oxygensaturation	pressurecontrol	serialnumber	ventilatorid	ingest_year	ingest_month	ingest_day
13	2023-06-26 16:03:24	Vyair	6	98	10	8e438321-3bee-423f-9bcd-c693ee475868	91	2023	06	26
17	2023-06-26 16:03:24	3M	5	98	17	a7bcb332-6c52-489e-9a55-c923f3f650d2	64	2023	06	26
19	2023-06-26 16:03:24	Getinge	7	98	24	871a5ed3-4912-4b51-8428-5cb3e1d0034a	30	2023	06	26
27	2023-06-26 16:04:24	Vyair	8	98	8	5e4eeeba-29bb-4add-9013-2307c640b09e	94	2023	06	26
29	2023-06-26 16:04:24	3M	7	98	26	69443bbd-f347-419a-97d0-912cb88b36eb	3	2023	06	26
31	2023-06-26 16:04:24	3M	7	98	16	9d6242e6-7f57-48a4-bbb6-3e1b954454be	8	2023	06	26

チュートリアル: AWS Glue Studio ノートブックを使用して最初のストリーミングワークロードを構築する

このチュートリアルでは、AWS Glue Studio ノートブックを活用して ETL ジョブをインタラクティブに構築および改良し、ほぼリアルタイムのデータ処理を行う方法について説明します。このガイドでは、AWS Glue が初めての方も、スキルセットの向上を目指している方も、AWS Glue インタラクティブセッションノートブックの可能性を最大限に引き出すことができるよう、順を追って説明していきます。

AWS Glue Streaming では、Amazon Kinesis Data Streams、Apache Kafka、Amazon Managed Streaming for Apache Kafka (Amazon MSK) などのストリーミングソースからのデータを消費し、連続的に実行するストリーミング抽出、変換、ロード (ETL) ジョブを作成できます。

前提条件

このチュートリアルを実行するには、AWS Glue、Amazon Kinesis、Amazon S3、Amazon Athena、AWS CloudFormation、AWS Lambda、Amazon Cognito を使用する AWS コンソール権限を持つユーザーが必要です。

Amazon Kinesis からストリーミングデータを消費する

トピック

- [Kinesis Data Generator でモックデータを生成する](#)
- [AWS Glue Studio で AWS Glue Streaming ジョブを作成する](#)
- [クリーンアップ](#)
- [結論](#)

Kinesis Data Generator でモックデータを生成する

Note

前の「[チュートリアル: AWS Glue Studio を使用して最初のストリーミングワークロードを構築する](#)」を既に完了している場合は、Kinesis Data Generator がアカウントに既にインストールされているため、以下のステップ 1~8 をスキップしてセクション「[AWS Glue Studio で AWS Glue Streaming ジョブを作成する](#)」に進んでください。

Kinesis Data Generator (KDG) を使用して JSON 形式のサンプルデータを合成的に生成できます。完全な手順と詳細については、[ツールのドキュメント](#)を参照してください。

1. 開始するに

は、

 Launch Stack 

をクリックして、ご使用の AWS 環境で AWS CloudFormation テンプレートを実行します。

Note

Kinesis Data Generator の Amazon Cognito ユーザーなどのリソースが AWS アカウントに既に存在しているために、CloudFormation テンプレートの実行が失敗する場合があります。これは、別のチュートリアルやブログで既に設定していたことが原因である可能性があります。これに対処するには、新しい AWS アカウントでテンプレートの実行をやり直すか、別の AWS リージョンを試してみてください。これらの方法により、既存のリソースと競合することなくチュートリアルを実行できるようになります。

このテンプレートにより、Kinesis データストリームと Kinesis Data Generator アカウントがプロビジョニングされます。

2. KDG が認証に使用する [ユーザー名] と [パスワード] を入力します。後で使用するために、ユーザー名とパスワードをメモしておきます。
3. 最後のステップまで [次へ] を選択していきます。IAM リソースの作成を承認します。パスワードが最小要件を満たしていないなど、画面上部にエラーがないか確認し、テンプレートをデプロイします。
4. スタックの [出力] タブに移動します。テンプレートがデプロイされると、生成されたプロパティ `KinesisDataGeneratorUrl` が表示されます。その URL をクリックします。
5. メモしておいた [ユーザー名] と [パスワード] を入力します。
6. 使用しているリージョンを選択し、Kinesis ストリーム `GlueStreamTest-{AWS::AccountId}` を選択します。
7. 以下のテンプレートを入力します。

```
{
  "ventilatorid": {{random.number(100)}},
  "eventtime": "{{date.now("YYYY-MM-DD HH:mm:ss")}}",
  "serialnumber": "{{random.uuid}}",
  "pressurecontrol": {{random.number(
    {
      "min":5,
      "max":30
    }
  )}},
  "o2stats": {{random.number(
    {
      "min":92,
```

```
        "max":98
      }
    )}},
    "minutevolume": {{random.number(
      {
        "min":5,
        "max":8
      }
    )}},
    "manufacturer": "{{random.arrayElement(
      ["3M", "GE","Vyair", "Getinge"]
    )}}"
```

[テストテンプレート] を使用してモックデータを表示し、[データを送信する] を使用してモックデータを Kinesis に取り込むことができるようになりました。

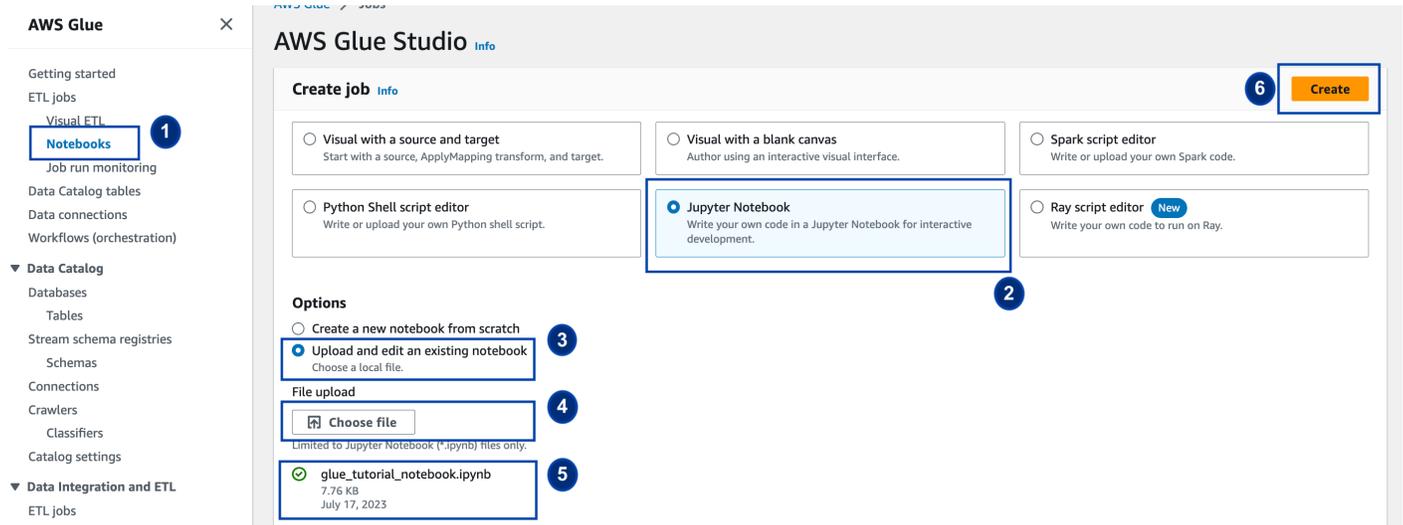
8. [データを送信する] をクリックして、5~10K のレコードを Kinesis に生成します。

AWS Glue Studio で AWS Glue Streaming ジョブを作成する

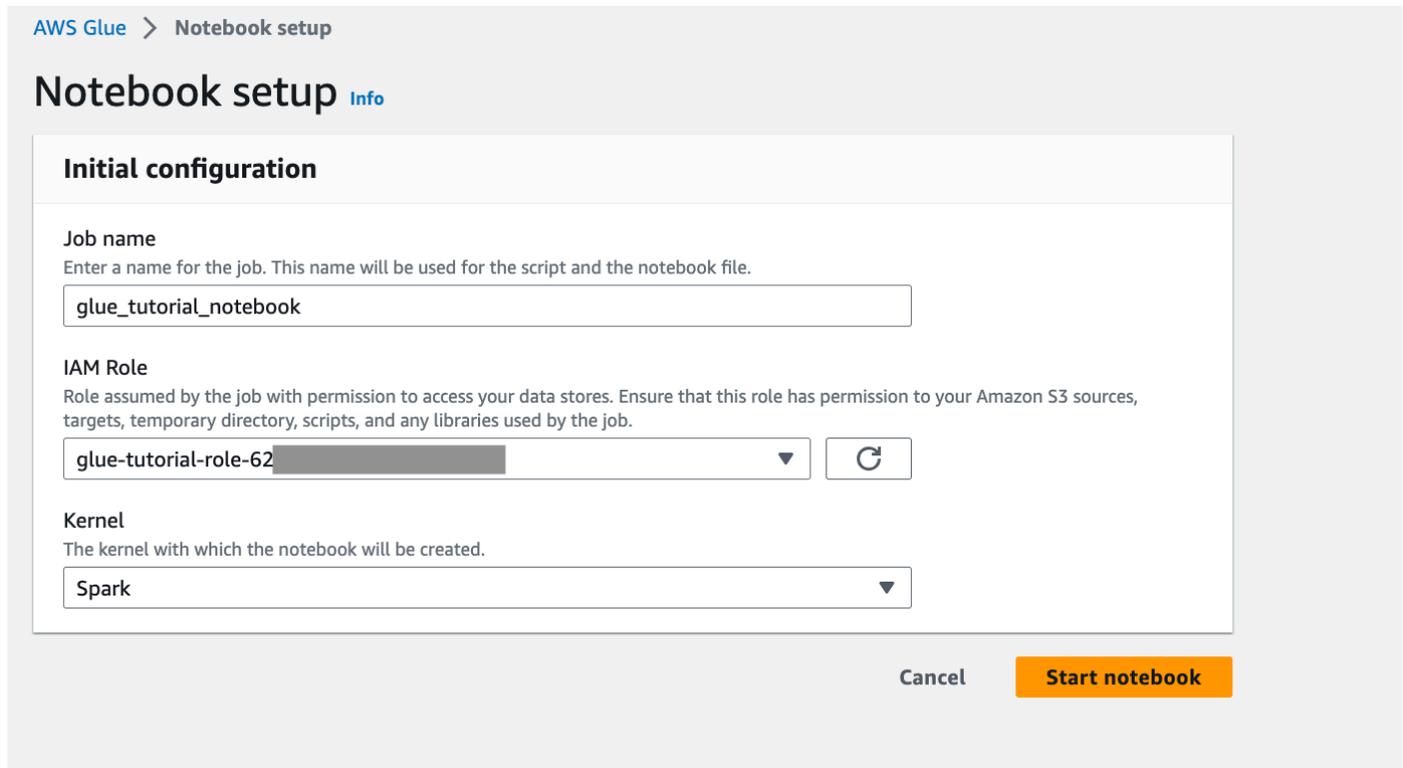
AWS Glue Studio は、データ統合パイプラインの設計、オーケストレーション、監視のプロセスを簡素化するビジュアルインターフェイスです。これにより、ユーザーは大量のコードを記述しなくてもデータ変換パイプラインを構築できます。ビジュアルジョブオーサリングエクスペリエンスとは別に、AWS Glue Studio には、このチュートリアルで使用する AWS Glue インタラクティブセッションによって支援された Jupyter Notebook も含まれています。

AWS Glue Streaming インタラクティブセッションジョブを設定する

1. 提供されている [ノートブックファイル](#) をダウンロードして、ローカルディレクトリに保存します。
2. AWS Glue コンソールを開き、左側のペインで [ノートブック] > [Jupyter Notebook] > [既存のノートブックをアップロードして編集する] の順にクリックします。前のステップのノートブックをアップロードし、[作成] をクリックします。



3. ジョブに対して名前とロールを指定し、デフォルトの Spark カーネルを選択します。次に [ノートブックの開始] をクリックします。[IAM ロール] では、CloudFormation テンプレートによってプロビジョニングされたロールを選択します。これは CloudFormation の [出力] タブで確認できます。



ノートブックには、このチュートリアルを続けるために必要なすべての手順が記載されています。ノートブックに記載されている手順を実行することも、このチュートリアルに従ってジョブの開発を続けることもできます。

ノートブックのセルを実行する

1. (オプション) 最初のコードセルの `%help` は、使用可能なすべてのノートブックマジックをリストします。ここではこのセルはスキップしても構いませんが、自由に試してみてください。
2. 次のコードブロック `%streaming` から始めます。このマジックにより、ジョブタイプがストリーミングに設定され、AWS Glue Streaming ETL ジョブを開発、デバッグ、デプロイできるようになります。
3. 次のセルを実行して AWS Glue インタラクティブセッションを作成します。出力セルには、セッションの作成を確認するメッセージが表示されます。

Run this cell to set up and start your interactive session.

```
[1]: %glue_version 3.0

import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from awsglue import DynamicFrame
from datetime import datetime
from pyspark.sql.types import StructType, StructField, StringType, LongType
from pyspark.sql.functions import lit,col,from_json
import boto3

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)

Setting Glue version to: 3.0
Authenticating with environment variables and user-defined glue_role_arn: arn:aws:iam::6[redacted]:role/glue-tutorial-role
Trying to create a Glue session for the kernel.
Worker Type: G.1X
Number of Workers: 5
Session ID: [redacted]af
Job Type: gluestreaming
Applying the following default arguments:
--glue_kernel_version 0.37.3
--enable-glue-datacatalog true
Waiting for session 4[redacted] to get into ready status...
Session 48 [redacted] has been created.
```

4. 次のセルでは変数を定義します。値をジョブに適した値に置き換え、セルを実行します。例:

```
output_database_name="default"
output_table_name="test_stream_001"

account_id = boto3.client("sts").get_caller_identity()["Account"]
region_name=boto3.client('s3').meta.region_name
stream_arn_name = "arn:aws:kinesis:{{region_name}}:{{account_id}}:stream/GlueStreamTest-{{account_id}}"
s3_bucket_name = "streaming-tutorial-s3-target-{{account_id}}"

output_location = "s3://{{s3_bucket_name}}/streaming_output/"
checkpoint_location = "s3://{{s3_bucket_name}}/checkpoint_location/"
```

5. データは既に Kinesis Data Streams にストリーミングされているため、次のセルではストリームからの結果を消費します。次のセルを実行します。print ステートメントがないため、このセルからの期待される出力はありません。

6. 次のセルでは、サンプルセットを取得して受信ストリームを調べ、そのスキーマと実際のデータを出力します。例:

Sample and print the incoming records

the sampling is for debugging purpose. You may comment off the entire code cell below, before deploying the actual code

```
[4]: options = {
  --- "pollingTimeInMs": "20000",
  --- "windowSize": "5 seconds"
}
sampled_dynamic_frame = glueContext.getSampleStreamingDynamicFrame(data_frame, options, None)

count_of_sampled_records = sampled_dynamic_frame.count()

print(count_of_sampled_records)

sampled_dynamic_frame.printSchema()

sampled_dynamic_frame.toDF().show(10, False)
```

```
100
root
```

```
|-- eventtime: string
|-- manufacturer: string
|-- minutevolume: long
|-- o2stats: long
|-- pressurecontrol: long
|-- serialnumber: string
|-- ventilatorid: long
```

eventtime	manufacturer	minutevolume	o2stats	pressurecontrol	serialnumber	ventilatorid
2023-07-18 10:20:11	3M	6	92	24	a3e860ba-24b9-41c4-bc10-91c6b35e1406	6
2023-07-18 10:20:11	Vyaire	6	95	6	96101dca-3e88-457f-b390-e3291df48a81	26
2023-07-18 10:20:12	Getinge	8	96	24	18f3d448-1dee-4c80-835b-1a0daa818915	22
2023-07-18 10:20:12	Getinge	7	98	30	25f425cd-b978-4953-9a03-4d607a639364	91
2023-07-18 10:20:12	GE	5	93	25	2cd7cdc2-f5f5-4ff2-ae32-45e5a8922d53	93

7. 次に、実際のデータ変換ロジックを定義します。セルは、マイクロバッチごとにトリガーされる processBatch メソッドで構成されます。セルを実行します。概説すると、受信ストリームに対して次のことを行います。
- 入力列のサブセットを選択します。
 - 列の名前を変更します (o2stats を oxygen_stats に変更します)。
 - 新しい列 (serial_identifier、ingest_year、ingest_month、ingest_day) を取得します。
 - 結果を Amazon S3 バケットに保存し、パーティション化された AWS Glue カタログテーブルも作成します。
8. 最後のセルでは、10 秒ごとにプロセスバッチをトリガーします。セルを実行し、Amazon S3 バケットと AWS Glue カタログテーブルに入力されるまで約 30 秒待ちます。
9. 最後に、Amazon Athena クエリエディタを使用して、保存されたデータを参照します。名前が変更された列と新しいパーティションが表示されます。

1 `select * from test_stream_001 limit 10`

SQL Ln 1, Col 39

Run again Explain Cancel Clear Create

Reuse query results up to 60 minutes ago

Query results Query stats

Completed Time in queue: 164 ms Run time: 1.22 sec Data scanned: 11.76 KB

Results (10) Copy Download results

Search rows

time	manufacturer	oxygen_stats	serialnumber	ventilatorid	serial_identifier	ingest_year	ingest_month	ingest_day
7-18 14:08:12	GE	96	a28895a3-0d57-4d0e-9d5e-86fdc92a5ba8	54	a28895a3	2023	7	18
7-18 14:08:12	Getinge	93	1e7b6e7e-e248-4cc7-971c-7cc7f4bb53e9	94	1e7b6e7e	2023	7	18
7-18 14:08:12	GE	97	52f8b540-4baa-4b90-bc65-986d668e8174	42	52f8b540	2023	7	18
7-18 14:08:12	Vyaire	93	e4ebdf4a-ca96-4465-ba03-681b438d9589	14	e4ebdf4a	2023	7	18
7-18 14:08:12	GE	92	52ba9e2b-748f-4226-9ac0-3767ce900233	33	52ba9e2b	2023	7	18
7-18 14:08:12	Getinge	96	74922910-ddcd-4e03-899b-acdf7487bb6c	8	74922910	2023	7	18

ノートブックには、このチュートリアルを続けるために必要なすべての手順が記載されています。ノートブックに記載されている手順を実行することも、このチュートリアルに従ってジョブの開発を続けることもできます。

AWS Glue ジョブを保存して実行する

インタラクティブセッションノートブックを使用したアプリケーションの開発とテストが完了したら、ノートブックインターフェイスの上部にある [保存] をクリックします。保存後は、アプリケーションをジョブとして実行することもできます。

glue_tutorial_notebook

Stop notebook Download Notebook Actions Save Run

Notebook Script Job details Runs Data quality New Schedules Version Control

Code Download

Glue PySpark

AWS Glue Streaming Tutorials - Working with Studio Notebook

クリーンアップ

アカウントに追加料金が発生しないように、手順の一部として開始したストリーミングジョブを停止してください。これを行うには、ノートブックを停止します。これでセッションが終了します。Amazon S3 バケットを空にして、前にプロビジョニングした AWS CloudFormation スタックを削除します。

結論

このチュートリアルでは、AWS Glue Studio ノートブックを使用して次のことを行う方法を説明しました。

- ノートブックを使用してストリーミング ETL ジョブを作成する
- 受信データストリームをプレビューする
- AWS Glue ジョブを公開することなくコーディングして問題を修正する
- エンドツーエンド作業コードを確認し、デバッグを除去して、ノートブックからステートメントやセルを出力する
- コードを AWS Glue ジョブとして公開する

このチュートリアルの目的は、AWS Glue Streaming とインタラクティブセッションを使用する実践的な使用体験を提供することです。AWS Glue Streaming の個々のユースケースの参考としてご利用いただくことをお勧めします。詳細については、「[AWS Glue のインタラクティブセッションの開始方法](#)」を参照してください。

AWS Glue Streaming の概念

以下のセクションでは、AWS Glue Streaming の概念について説明します。

トピック

- [AWS Glue Streaming ジョブの構造](#)
- [Kafka 接続](#)
- [Kinesis 接続](#)
- [AWS Glue ストリーミングオプション](#)

AWS Glue Streaming ジョブの構造

AWS Glue Streaming ジョブは Spark ストリーミングパラダイムで動作し、Spark フレームワークの構造化されたストリーミングを活用します。ストリーミングジョブは、特定の時間間隔で常にストリーミングデータソースをポーリングし、レコードをマイクロバッチとして取得します。以下のセクションでは、AWS Glue Streaming ジョブのさまざまな部分について説明します。

```
def processBatch(data_frame, batchId):
  2
  if data_frame.count() > 0:
    AmazonKinesis_node1696872487972 = DynamicFrame.fromDF(
      glueContext.add_ingestion_time_columns(data_frame, "hour"),
      glueContext,
      "from_data_frame",
    )
    # Script generated for node Change Schema
    ChangeSchema_node1696872679326 = ApplyMapping.apply(
      frame=AmazonKinesis_node1696872487972,
      mappings=[
        ("eventtime", "string", "eventtime", "string"),
        ("manufacturer", "string", "manufacturer", "string"),
        ("minutevolume", "long", "minutevolume", "int"),
        ("o2stats", "long", "OxygenSaturation", "int"),
        ("pressurecontrol", "long", "pressurecontrol", "int"),
        ("serialnumber", "string", "serialnumber", "string"),
        ("ventilatorid", "long", "ventilatorid", "long"),
        ("ingest_year", "string", "ingest_year", "string"),
        ("ingest_month", "string", "ingest_month", "string"),
        ("ingest_day", "string", "ingest_day", "string"),
        ("ingest_hour", "string", "ingest_hour", "string"),
      ],
      transformation_ctx="ChangeSchema_node1696872679326",
    )
    # Script generated for node Amazon S3
    AmazonS3_node1696872743449_path = (
      "s3://streaming-tutorial-s3-target-
    )
    AmazonS3_node1696872743449 = glueContext.getSink(
      path=AmazonS3_node1696872743449_path,
      connection_type="s3",
      updateBehavior="UPDATE_IN_DATABASE",
      partitionKeys=["ingest_year", "ingest_month", "ingest_day", "ingest_hour"],
      compression="snappy",
      enableUpdateCatalog=True,
      transformation_ctx="AmazonS3_node1696872743449",
    )
    AmazonS3_node1696872743449.setCatalogInfo(
      catalogDatabase="demo", catalogTableName="demo_stream_transform_result"
    )
    AmazonS3_node1696872743449.setFormat("glueparquet")
    AmazonS3_node1696872743449.writeFrame(ChangeSchema_node1696872679326)
    glueContext.forEachBatch(
      frame=dataFrame_AmazonKinesis_node1696872487972,
      batch_function=processBatch,
      options={
        "windowSize": "100 seconds",
        "checkpointLocation": args["TempDir"] + "/" + args["JOB_NAME"] + "/checkpoint/",
      },
    ),
  },
  job.commit()
```

1 ← Entry Point

forEachBatch

forEachBatch メソッドは、AWS Glue Streaming ジョブ実行のエントリーポイントです。AWS Glue Streaming ジョブは、forEachBatch メソッドを使用して、イテレーターのように機能するデータをポーリングします。このデータはストリーミングジョブのライフサイクル中ずっとアクティブなままです。新しいデータがないかストリーミングソースを定期的にポーリングして、最新のデータをマイクロバッチで処理します。

```
glueContext.forEachBatch(
  frame=dataFrame_AmazonKinesis_node1696872487972,
  batch_function=processBatch,
  options={
```

```
        "windowSize": "100 seconds",
        "checkpointLocation": args["TempDir"] + "/" + args["JOB_NAME"] + "/"
checkpoint/",
    },
)
```

forEachBatch の frame プロパティを設定してストリーミングソースを指定します。この例では、ジョブの作成中に空白のキャンバスに作成したソースノードに、ジョブのデフォルトの DataFrame が入力されます。batch_function プロパティを、マイクロバッチ操作ごとに呼び出す function として設定します。受信データのバッチ変換を処理する関数を定義する必要があります。

ソース

processBatch 関数の最初のステップで、プログラムは、forEachBatch のフレームプロパティとして定義した DataFrame のレコード数を検証します。プログラムは、空でない DataFrame に取り込みタイムスタンプを追加します。data_frame.count()>0 句は、最新のマイクロバッチが空ではなく、さらに処理できる状態にあるかどうかを判断します。

```
def processBatch(data_frame, batchId):
    if data_frame.count() >0:
        AmazonKinesis_node1696872487972 = DynamicFrame.fromDF(
            glueContext.add_ingestion_time_columns(data_frame, "hour"),
            glueContext,
            "from_data_frame",
        )
```

マッピング

プログラムの次のセクションでは、マッピングを適用します。Spark DataFrame で Mapping.apply メソッドを使用すると、データ要素に関する変換ルールを定義できます。通常、名前を変更したり、データ型を変更したり、ソースデータ列でカスタム関数を適用してターゲット列にマッピングしたりすることができます。

```
#Script generated for node ChangeSchema
ChangeSchema_node16986872679326 = ApplyMapping.apply(
    frame = AmazonKinesis_node1696872487972,
    mappings = [
```

```

        ("eventtime", "string", "eventtime", "string"),
        ("manufacturer", "string", "manufacturer", "string"),
        ("minutevolume", "long", "minutevolume", "int"),
        ("o2stats", "long", "OxygenSaturation", "int"),
        ("pressurecontrol", "long", "pressurecontrol", "int"),
        ("serialnumber", "string", "serialnumber", "string"),
        ("ventilatorid", "long", "ventilatorid", "long"),
        ("ingest_year", "string", "ingest_year", "string"),
        ("ingest_month", "string", "ingest_month", "string"),
        ("ingest_day", "string", "ingest_day", "string"),
        ("ingest_hour", "string", "ingest_hour", "string"),
    ],
    transformation_ctx="ChangeSchema_node16986872679326",
)
)

```

シンク

このセクションでは、ストリーミングソースから受信したデータセットがターゲットの場所に保存されます。この例では、Amazon S3 の場所にデータを書き込みます。AmazonS3_node_path プロパティの詳細は、キャンバスからのジョブ作成時に使用した設定に基づいて事前入力されます。ユースケースに基づいて `updateBehavior` を設定し、データカタログテーブルを更新しないか、それ以降の実行時にデータカタログを作成してデータカタログスキーマを更新するか、カタログテーブルを作成してそれ以降の実行時にスキーマ定義を更新しないかを決定できます。

`partitionKeys` プロパティは、ストレージパーティションオプションを定義します。デフォルトの動作では、ソースセクションで使用できるようになった `ingestion_time_columns` ごとにデータをパーティション化します。`compression` プロパティでは、ターゲットへの書き込み時に適用される圧縮アルゴリズムを設定できます。圧縮方式として Snappy、LZO、または GZIP を設定するオプションがあります。`enableUpdateCatalog` プロパティは、AWS Glue カタログテーブルを更新する必要があるかどうかを制御します。このプロパティで使用できるオプションは `True` または `False` です。

```

#Script generated for node Amazon S3
AmazonS3_node1696872743449 = glueContext.getSink(
    path = AmazonS3_node1696872743449_path,
    connection_type = "s3",
    updateBehavior = "UPDATE_IN_DATABASE",
    partitionKeys = ["ingest_year", "ingest_month", "ingest_day", "ingest_hour"],

```

```
compression = "snappy",
enableUpdateCatalog = True,
transformation_ctx = "AmazonS3_node1696872743449",
)
```

AWS Glue カタログシンク

ジョブのこのセクションでは、AWS Glue カタログテーブルの更新動作を制御します。設計中の AWS Glue ジョブに関連する AWS Glue カタログデータベース名とテーブル名ごとに catalogDatabase プロパティと catalogTableName プロパティを設定します。setFormat プロパティを使用してターゲットデータのファイル形式を定義できます。この例では、データを parquet 形式で保存します。

このチュートリアルを参照して AWS Glue Streaming ジョブを設定して実行すると、Amazon Kinesis Data Streams で生成されたストリーミングデータは、Snappy 圧縮方式、parquet 形式で、Amazon S3 の場所に保存されます。ストリーミングジョブが正常に実行されると、Amazon Athena でデータをクエリできるようになります。

```
AmazonS3_node1696872743449 = setCatalogInfo(
    catalogDatabase = "demo", catalogTableName = "demo_stream_transform_result"
)
AmazonS3_node1696872743449.setFormat("glueparquet")
AmazonS3_node1696872743449.writeFormat("ChangeSchema_node16986872679326")
)
```

Kafka 接続

Kafka クラスターまたは Amazon Managed Streaming for Apache Kafka への接続を指定します。

Data Catalog テーブルに格納されている情報を使用するか、データストリームに直接アクセスするための情報を指定することにより、Kafka データストリームへ読み込むまたは書き込むことができます。Kafka から Spark に情報を読み取り DataFrame、それを Glue に変換できます。AWS DynamicFrameKafka には JSON DynamicFrames 形式で書き込むことができます。データストリームに直接アクセスする場合は、これらのオプションを使用して、データストリームへのアクセス方法に関する情報を提供します。

getCatalogSource または create_data_frame_from_catalog を使用して Kafka ストリーミングソースからレコードを消費するか、getCatalogSink または write_dynamic_frame_from_catalog を使用して Kafka にレコードを書き込み、ジョブに Data Catalog データベースおよびテーブル名の情報があり、それを使用して Kafka ストリーミングソースから読み込むためのいくつかの基本パラメータを取得できる場合。getSource、getCatalogSink、getSourceWithFormat、getSinkWithFormat、createDataFrame を使用する場合、ここで説明する接続オプションを使用し、これらの基本パラメータを指定する必要があります。

GlueContext クラスで指定されたメソッドの次の引数を使用し、Kafka の接続オプションを指定できます。

- Scala
 - connectionOptions: getSource、createDataFrameFromOptions、getSink で使用
 - additionalOptions: getCatalogSource、getCatalogSink で使用
 - options: getSourceWithFormat、getSinkWithFormat で使用
- Python
 - connection_options:
create_data_frame_from_options、write_dynamic_frame_from_options で使用
 - additional_options:
create_data_frame_from_catalog、write_dynamic_frame_from_catalog で使用
 - options: getSource、getSink で使用

ストリーミング ETL ジョブに関する注意事項と制限事項については、「[the section called “ストリーミング ETL に関する注意と制限”](#)」を参照してください。

Kafka の設定

インターネット経由で利用可能な Kafka AWS ストリームに接続するための前提条件はありません。

AWS Glue Kafka 接続を作成して、接続認証情報を管理できます。詳細については、「[the section called “Kafka データストリームの接続の作成”](#)」を参照してください。AWS Glue ジョブ設定で、**ConnectionName #####**、メソッド呼び出しでパラメーターに **ConnectionName #####**。connectionName

場合によっては、追加の前提条件を設定する必要があります。

- IAM 認証で Amazon Managed Streaming for Apache Kafka を使用する場合は、適切な IAM 設定が必要になります。
- Amazon VPC で Amazon Managed Streaming for Apache Kafka を使用する場合は、適切な Amazon VPC 設定が必要になります。Amazon VPC 接続情報を提供する AWS Glue 接続を作成する必要があります。AWS Glue 接続を追加ネットワーク接続として含めるには、ジョブ設定が必要です。

ストリーミング ETL ジョブの前提条件の詳細については、「[the section called “ストリーミング ETL ジョブ”](#)」を参照してください。

例: Kafka ストリームからの読み込み

[the section called “forEachBatch”](#) と組み合わせて使用します。

Kafka ストリーミングソースの例 :

```
kafka_options =
  { "connectionName": "ConfluentKafka",
    "topicName": "kafka-auth-topic",
    "startingOffsets": "earliest",
    "inferSchema": "true",
    "classification": "json"
  }
data_frame_datasource0 =
  glueContext.create_data_frame.from_options(connection_type="kafka",
  connection_options=kafka_options)
```

例: Kafka ストリームへの書き込み

Kafka への書き込み例:

getSink メソッドを使った例:

```
data_frame_datasource0 =
  glueContext.getSink(
    connectionType="kafka",
    connectionOptions={
      JsonOptions("""{
        "connectionName": "ConfluentKafka",
        "classification": "json",
```

```

    "topic": "kafka-auth-topic",
    "typeOfData": "kafka"}
    """}},
transformationContext="dataframe_ApacheKafka_node1711729173428")
.getDataFrame()

```

`write_dynamic_frame.from_options` メソッドを使った例:

```

kafka_options =
  { "connectionName": "ConfluentKafka",
    "topicName": "kafka-auth-topic",
    "classification": "json"
  }
data_frame_datasource0 =
  glueContext.write_dynamic_frame.from_options(connection_type="kafka",
  connection_options=kafka_options)

```

Kafka 接続オプションのリファレンス

読み込むとき、`"connectionType": "kafka"` に次の接続オプションを使用します。

- `"bootstrap.servers"` (必須) ブートストラップサーバーの URL のリスト (例: `b-1.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094`)。このオプションは API 呼び出しで指定するか、データカタログ内のテーブルメタデータで定義する必要があります。
- `"security.protocol"` (必須) ブローカーと通信するために使用されるプロトコル。使用できる値は、`"SSL"` または `"PLAINTEXT"` です。
- `"topicName"` (必須) サブスクライブするトピックのカンマ区切りリスト。"`topicName`"、"`assign`"、または "`subscribePattern`" の中から、いずれか 1 つのみを指定する必要があります。
- `"assign"`: (必須) 消費する特定の `TopicPartitions` を指定する JSON 文字列。"`topicName`"、"`assign`"、または "`subscribePattern`" の中から、いずれか 1 つのみを指定する必要があります。

例: `'{"topicA":[0,1],"topicB":[2,4]}'`

- `"subscribePattern"`: (必須) サブスクライブする先のトピックリストを識別する Java の正規表現文字列。"`topicName`"、"`assign`"、または "`subscribePattern`" の中から、いずれか 1 つのみを指定する必要があります。

例: `'topic.*'`

- "classification" (必須) レコード内のデータで使用されるファイル形式。データカタログを通じて提供されていない限り、必須です。
- "delimiter" (オプション) classification が CSV の場合に使用される値の区切り文字。デフォルトは「,」です。
- "startingOffsets": (オプション) Kafka トピック内で、データの読み取りを開始する位置 使用できる値は、"earliest" または "latest" です。デフォルト値は "latest" です。
- "startingTimestamp": (オプション、AWS Glue バージョン 4.0 以降でのみサポート) データを読み取る Kafka トピック内のレコードのタイムスタンプ。指定できる値は UTC 形式 (yyyy-mm-ddTHH:MM:SSZ のパターン) のタイムスタンプ文字列です (Z は UTC タイムゾーンのオフセットを +/- で表します。例: 「2023-04-04T08:00:00-04:00」)。

注: AWS Glue ストリーミングスクリプトの接続オプションリストには、「StartingOffsets」または「StartingTimestamp」のいずれかしか表示できません。これらのプロパティの両方を含めると、ジョブが失敗します。

- "endingOffsets": (オプション) バッチクエリの終了位置。設定が可能な値は、"latest" または、各 TopicPartition の終了オフセットを指定する JSON 文字列のいずれかです。

JSON 文字列の場合、{"topicA":{"0":23,"1":-1},"topicB":{"0":-1}} の形式を使用します。オフセットとして値 -1 を設定する場合、"latest" の意味になります。

- "pollTimeoutMs": (オプション) Spark ジョブエグゼキュータで、Kafka からデータをポーリングする際のタイムアウト値 (ミリ秒単位)。デフォルト値は、512 です。
- "numRetries": (オプション) Kafka オフセットのフェッチが失敗したと判断される前の再試行回数。デフォルト値は、3 です。
- "retryIntervalMs": (オプション) Kafka オフセットのフェッチを開始するまでの待機時間 (ミリ秒)。デフォルト値は、10 です。
- "maxOffsetsPerTrigger": (オプション) 処理されるオフセットの最大数を、トリガー間隔ごとのレート上限で指定する値。指定されたオフセットの合計数は、異なるボリュームの topicPartitions 間で均等に分割されます。デフォルト値は「null」です。この場合、コンシューマーは既知の最新のオフセットまで、すべてのオフセットを読み取ります。
- "minPartitions": (オプション) Kafka から読み取ることを想定する、最小のパーティション数。デフォルト値は「null」です。これは、Spark パーティションの数が Kafka パーティションの数に等しいことを意味します。
- "includeHeaders": (オプション) Kafka ヘッダーを含めるかどうかを決定します。このオプションが「true」に設定されている場合、データ出力には、「glue_streaming_kafka_headers」という

名前で `Array[Struct(key: String, value: String)]` 型の列が追加されます。デフォルト値は「false」です。このオプションは、AWS Glue バージョン 3.0 以降でのみ使用可能です。

- "schema": (inferSchema に false を設定した場合は必須) ペイロードの処理に使用するスキーマ。分類が avro である場合、提供されるスキーマには Avro スキーマ形式を使用する必要があります。分類が avro 以外の場合、提供されるスキーマは DDL スキーマ形式である必要があります。

以下に、スキーマの例を示します。

Example in DDL schema format

```
'column1' INT, 'column2' STRING , 'column3' FLOAT
```

Example in Avro schema format

```
{
  "type": "array",
  "items":
  {
    "type": "record",
    "name": "test",
    "fields":
    [
      {
        "name": "_id",
        "type": "string"
      },
      {
        "name": "index",
        "type":
        [
          "int",
          "string",
          "float"
        ]
      }
    ]
  }
}
}
```

- "inferSchema": (オプション) デフォルト値は「false」です。「true」に設定すると、実行時に、スキーマが foreachbatch 内のペイロードから検出されます。

- "avroSchema": (非推奨) Avro 形式を使用する場合に、Avro データのスキーマを指定するために使用されるパラメータです。このパラメータは非推奨となりました。schema パラメータを使用します。
- "addRecordTimestamp": (オプション) このオプションを「true」に設定すると、トピックが対応するレコードを受信した時刻を表示する「__src_timestamp」という列が、データ出力に追加で表示されます。デフォルト値は、「false」です。このオプションは AWS Glue バージョン 4.0 以降でサポートされています。
- "emitConsumerLagMetrics": (オプション) オプションを 'true' に設定すると、バッチごとに、トピックが受信した最も古いレコードからレコードが届くまでの期間のメトリクスが出力されます。AWS Glue CloudWatchメトリクスの名前は「glue.driver.stream」です。maxConsumerLagInMs」。デフォルト値は、「false」です。このオプションは AWS Glue バージョン 4.0 以降でサポートされています。

書き込むとき、"connectionType": "kafka" に次の接続オプションを使用します。

- "connectionName"(必須) Kafka クラスターへの接続に使用される AWS Glue 接続の名前 (Kafka ソースと同様)。
- "topic" (必須) トピック列が存在する場合、トピック設定オプションが設定されていない限り、指定された行を Kafka に書き込む際にトピック列の値がトピックとして使用されます。つまり、topic 設定オプションはトピック列をオーバーライドします。
- "partition" (オプション) 有効なパーティション番号が指定された場合、レコードの送信時にその partition が使用されます。

パーティションが指定されずに key が存在する場合、キーのハッシュを使用してパーティションが選択されます。

key と partition のどちらも存在しない場合、そのパーティションに少なくとも batch.size バイトが生成された際に変更内容のステイッキーパーティション分割に基づいてパーティションが選択されます。

- "key" (オプション) partition が null の場合、パーティション分割に使用されます。
- "classification" (オプション) レコードのデータに使用されるファイル形式。JSON、CSV、Avro のみをサポートしています。

Avro 形式を使用すると、シリアル化するためにカスタムの AvroSchema を提供できますが、シリアル化解除する場合にもソースに提供する必要がありますことに注意してください。それ以外の場合、デフォルトではシリアル化に Apache を使用します。AvroSchema

さらに、[Kafka プロデューサーの設定パラメーター](#)を更新することにより、必要に応じて Kafka シンクを微調整できます。接続オプションには許可リストがないことに注意してください。すべてのキー値のペアはそのままシンクに保持されます。

ただし、有効にならないオプションの小さな拒否リストがあります。詳細については、「[Apache 固有の設定](#)」を参照してください。

Kinesis 接続

Amazon Kinesis Data Streams から読み取る、または Amazon Kinesis Data Streams に書き込むには、データカタログテーブルに格納されている情報を使用するか、データストリームに直接アクセスするための情報を指定します。Kinesis から Spark に情報を読み取り DataFrame、それを Glue に変換することができます。AWS DynamicFrameKinesis には JSON DynamicFrames 形式で書き込むことができます。データストリームに直接アクセスする場合は、これらのオプションを使用して、データストリームへのアクセス方法に関する情報を提供します。

`getCatalogSource` または `create_data_frame_from_catalog` を使用して Kinesis ストリーミングソースからレコードを消費する場合、ジョブはデータカタログデータベースとテーブル名の情報を持っており、それを使用して Kinesis ストリーミングソースから読み込むためのいくつかの基本パラメータを取得することができます。`getSource`、`getSourceWithFormat`、`createDataFrameFromOptions`、または `create_data_frame_from_options` を使用している場合、ここで説明する接続オプションを使用して、これらの基本パラメータを指定する必要があります。

Kinesis の接続オプションは、`GlueContext` クラス内の指定されたメソッドに対して以下の引数で指定することができます。

- Scala
 - `connectionOptions`: `getSource`、`createDataFrameFromOptions`、`getSink` で使用
 - `additionalOptions`: `getCatalogSource`、`getCatalogSink` で使用
 - `options`: `getSourceWithFormat`、`getSinkWithFormat` で使用
- Python
 - `connection_options`:
`create_data_frame_from_options`、`write_dynamic_frame_from_options` で使用
 - `additional_options`:
`create_data_frame_from_catalog`、`write_dynamic_frame_from_catalog` で使用
 - `options`: `getSource`、`getSink` で使用

ストリーミング ETL ジョブに関する注意事項と制限事項については、「[the section called “ストリーミング ETL に関する注意と制限”](#)」を参照してください。

Kinesis の設定

AWS Glue Spark ジョブで Kinesis データストリームに接続するには、いくつかの前提条件が必要です。

- 読み取りの場合、AWS Glue ジョブには Kinesis データストリームへの読み取りアクセスレベル IAM 権限が必要です。
- 書き込みを行う場合、AWS Glue ジョブには Kinesis データストリームへの書き込みアクセスレベルの IAM 権限が必要です。

場合によっては、追加の前提条件を設定する必要があります。

- AWS Glue ジョブが (通常は他のデータセットに接続するための) 追加のネットワーク接続で設定されていて、その接続のいずれかが Amazon VPC ネットワークオプションを提供している場合、ジョブは Amazon VPC 経由で通信するように指示されます。この場合、Amazon VPC を介して通信するように Kinesis データストリームを設定する必要があります。そのためには、Amazon VPC と Kinesis データストリームの間インターフェイス VPC エンドポイントを作成します。詳細については、「[インターフェイス VPC エンドポイントと Amazon Kinesis Data Streams の使用](#)」を参照してください。
- 別のアカウントで Amazon Kinesis Data Streams を指定する場合は、クロスアカウントアクセスを許可するようにロールとポリシーを設定する必要があります。詳細については、「[Example: Read From a Kinesis Stream in a Different Account](#)」を参照してください。

ストリーミング ETL ジョブの前提条件の詳細については、「[the section called “ストリーミング ETL ジョブ”](#)」を参照してください。

Kinesis からの読み込み

例: Kinesis ストリームからの読み込み

[the section called “forEachBatch”](#) と組み合わせて使用します。

Amazon Kinesis ストリーミングソースの例:

```
kinesis_options =
```

```
{ "streamARN": "arn:aws:kinesis:us-east-2:777788889999:stream/fromOptionsStream",
  "startingPosition": "TRIM_HORIZON",
  "inferSchema": "true",
  "classification": "json"
}
data_frame_datasource0 =
glueContext.create_data_frame.from_options(connection_type="kinesis",
connection_options=kinesis_options)
```

Kinesis への書き込み

例: Kinesis ストリームへの書き込み

[the section called “forEachBatch”](#) と組み合わせて使用します。JSON DynamicFrame 形式でストリームに書き込まれます。何度か再試行してもジョブの書き込みが実行できないと、ジョブは失敗します。デフォルトでは、DynamicFrame 各レコードは Kinesis ストリームに個別に送信されます。この動作は、`aggregationEnabled` と関連するパラメータを使用して設定できます。

ストリーミングジョブから Amazon Kinesis への書き込みの例:

Python

```
glueContext.write_dynamic_frame.from_options(
    frame=frameToWrite
    connection_type="kinesis",
    connection_options={
        "partitionKey": "part1",
        "streamARN": "arn:aws:kinesis:us-east-1:111122223333:stream/streamName",
    }
)
```

Scala

```
glueContext.getSinkWithFormat(
    connectionType="kinesis",
    options=JsonOptions("""{
        "streamARN": "arn:aws:kinesis:us-
east-1:111122223333:stream/streamName",
        "partitionKey": "part1"
    }"""),
)
.writeDynamicFrame(frameToWrite)
```

Kinesis 接続パラメータ

Amazon Kinesis Data Streams への接続を指定します。

Kinesis ストリーミングデータソースには、次の接続オプションを使用します。

- "streamARN" (必須) 読み取り/書き込みに使用されます。Kinesis データストリームの ARN。
- "classification" (読み取りに必須) 読み取りで使用。レコード内のデータで使用されるファイル形式。データカタログを通じて提供されていない限り、必須です。
- "streamName" - (オプション) 読み取りで使用。読み取り対象/読み取り元の Kinesis データストリームの名前。endpointUrl で使用。
- "endpointUrl" - (オプション) 読み取りで使用。デフォルト: "https://kinesis.us-east-1.amazonaws.com"。Kinesis AWS ストリームのエンドポイント。特別なリージョンに接続する場合を除き、これを変更する必要はありません。
- "partitionKey" - (オプション) 書き込みに使用。レコードを作成する際に使用される Kinesis パーティションキー。
- "delimiter" (オプション) 読み取りに使用。classification が CSV の場合に使用される値の区切り文字。デフォルトは「,」です。
- "startingPosition": (オプション) 読み込みに使用。Kinesis データストリーム内の、データの読み取り開始位置。指定できる値は "latest"、"trim_horizon"、"earliest"、または UTC 形式のタイムスタンプ文字列であり、この文字列のパターンは yyyy-mm-ddTHH:MM:SSZ です (Z は UTC タイムゾーンのオフセットを +/- で表します。例: 「2023-04-04T08:00:00-04:00」)。デフォルト値は、"latest" です。注: の UTC "startingPosition" 形式のタイムスタンプ文字列は、AWS Glue バージョン 4.0 以降でのみサポートされます。
- "failOnDataLoss": (オプション) アクティブなシャードがないか、有効期限が切れている場合、ジョブは失敗します。デフォルト値は、"false" です。
- "awsSTSRoleARN": (オプション) 読み取り/書き込みに使用。() を使用して引き受けるロールの Amazon リソースネーム AWS Security Token Service (ARN AWS STS)。このロールには、Kinesis データストリームのレコードの説明操作または読み取り操作の権限が必要です。このパラメーターは、別のアカウントのデータストリームにアクセスするときに使用する必要があります。"awsSTSSessionName" と組み合わせて使用します。
- "awsSTSSessionName": (オプション) 読み取り/書き込みに使用。AWS STS を使って、ロールを担うセッションの識別子。このパラメーターは、別のアカウントのデータストリームにアクセスするときに使用する必要があります。"awsSTSRoleARN" と組み合わせて使用します。

- "awsSTSEndpoint": (オプション) 引き受けたロールで Kinesis AWS STS に接続するとき使用するエンドポイント。これにより、AWS STS デフォルトのグローバルエンドポイントでは不可能な VPC 内のリージョナルエンドポイントを使用できます。
- "maxFetchTimeInMs": (オプション) 読み込みに使用。ジョブエグゼキューターが Kinesis データストリームから現在のバッチのレコードを読み取るために費やした最大時間は、ミリ秒 (ms) 単位で指定されます。この時間内に複数の GetRecords API コールを行うことができます。デフォルト値は、1000 です。
- "maxFetchRecordsPerShard": (オプション) 読み込みに使用。1 マイクロバッチあたりに Kinesis データストリームでシャードごとにフェッチするレコードの最大数。メモ: ストリーミングジョブが既に Kinesis (同じ get-records 呼び出しで) から余分なレコードを読み取っている場合、クライアントはこの制限を超えることができます。maxFetchRecordsPerShard が厳密である必要がある場合、maxRecordPerRead の倍数にする必要があります。デフォルト値は、100000 です。
- "maxRecordPerRead": (オプション) 読み込みに使用。getRecords オペレーションごとに、Kinesis データストリームからフェッチするレコードの最大数。デフォルト値は、10000 です。
- "addIdleTimeBetweenReads": (オプション) 読み込みに使用。2 つの連続する getRecords オペレーション間の遅延時間を追加します。デフォルト値は、"False" です。このオプションは、Glue バージョン 2.0 以降でのみ設定可能です。
- "idleTimeBetweenReadsInMs": (オプション) 読み込みに使用。2 つの連続する getRecords オペレーション間での、最短の遅延時間 (ミリ秒単位で指定)。デフォルト値は、1000 です。このオプションは、Glue バージョン 2.0 以降でのみ設定可能です。
- "describeShardInterval": (オプション) 読み込みに使用。スクリプトが呼び出す 2 つの ListShards API 間での、再シャードリングを考慮すべき最小時間。詳細については、Amazon Kinesis Data Streams デベロッパーガイドの「[リシャードリングのための戦略](#)」を参照してください。デフォルト値は、1s です。
- "numRetries": (オプション) 読み込みに使用。Kinesis Data Streams API リクエストを再試行する最大の回数。デフォルト値は、3 です。
- "retryIntervalMs": (オプション) 読み込みに使用。Kinesis Data Streams API 呼び出しを再試行するまでのクールオフ期間 (ミリ秒単位で指定)。デフォルト値は、1000 です。
- "maxRetryIntervalMs": (オプション) 読み込みに使用。再試行で 2 つの Kinesis Data Streams API を呼び出す間の最大クールオフ期間 (ミリ秒単位で指定)。デフォルト値は、10000 です。

- "avoidEmptyBatches": (オプション) 読み込みに使用。バッチ処理を開始する前に、Kinesis データストリームで未読のデータをチェックすることで、空のマイクロバッチジョブを作成しないようにします。デフォルト値は、"False"です。
- "schema": (inferSchema に false を設定した場合は必須) 読み取りに使用。ペイロードの処理に使用するスキーマ。分類が avro である場合、提供されるスキーマには Avro スキーマ形式を使用する必要があります。分類が avro 以外の場合、提供されるスキーマは DDL スキーマ形式である必要があります。

以下に、スキーマの例を示します。

Example in DDL schema format

```
`column1` INT, `column2` STRING , `column3` FLOAT
```

Example in Avro schema format

```
{
  "type": "array",
  "items":
  {
    "type": "record",
    "name": "test",
    "fields":
    [
      {
        "name": "_id",
        "type": "string"
      },
      {
        "name": "index",
        "type":
        [
          "int",
          "string",
          "float"
        ]
      }
    ]
  }
}
```

- "inferSchema": (オプション) 読み込みに使用。デフォルト値は、「false」です。「true」に設定すると、実行時に、スキーマが foreachbatch 内のペイロードから検出されます。
- "avroSchema": (非推奨) 読み取りに使用。Avro 形式を使用する場合に、Avro データのスキーマを指定するために使用されるパラメータです。このパラメータは非推奨となりました。schema パラメータを使用します。
- "addRecordTimestamp": (オプション) 読み込みに使用。このオプションが「true」に設定されている場合、データ出力には、対応するレコードがストリームによって受信された時刻を表示する「__src_timestamp」という名前が付けられた追加の列が含まれます。デフォルト値は、「false」です。このオプションは AWS Glue バージョン 4.0 以降でサポートされています。
- "emitConsumerLagMetrics": (オプション) 読み込みに使用。このオプションを 'true' に設定すると、バッチごとに、ストリームが受信した最も古いレコードからそのレコードが到着するまでの間のメトリクスが出力されます。AWS Glue CloudWatchメトリクスの名前は「glue.driver.stream」です。maxConsumerLagInMs」。デフォルト値は、「false」です。このオプションは AWS Glue バージョン 4.0 以降でサポートされています。
- "fanoutConsumerARN": (オプション) 読み込みに使用。streamARN で指定されたストリームの Kinesis ストリームコンシューマの ARN。Kinesis 接続の拡張ファンアウトモードを有効にするために使用されます。拡張ファンアウトが使用された Kinesis ストリームの使用に関する詳細については、「[the section called “Kinesis ストリーミングジョブでの拡張ファンアウトの使用”](#)」を参照してください。
- "recordMaxBufferedTime" - (オプション) 書き込みに使用。デフォルト: 1000 (ミリ秒)。レコードが書き込まれるのを待っている間にバッファリングされる最大時間。
- "aggregationEnabled" - (オプション) 書き込みに使用。デフォルト: true。Kinesis に送信する前にレコードを集約するかどうかを指定します。
- "aggregationMaxSize" - (オプション) 書き込みに使用。デフォルト: 51200 (バイト) レコードがこの制限よりも大きい場合、そのレコードはアグリゲータをバイパスします。注: Kinesis では、レコードサイズに 50 KB の制限が適用されます。これを 50 KB を超えて設定すると、サイズ超過のレコードは Kinesis によって拒否されます。
- "aggregationMaxCount" - (オプション) 書き込みに使用。デフォルト: 4294967295。集計されたレコードにパックされる項目の最大数。
- "producerRateLimit" - (オプション) 書き込みに使用。デフォルト: 150 (%)。1 つのプロデューサー (ジョブなど) からの送信されるシャード単位のスループットをバックエンド制限のパーセンテージとして制限できます。
- "collectionMaxCount" - (オプション) 書き込みに使用。デフォルト: 500。1 PutRecords 回のリクエストにまとめるアイテムの最大数。

- "collectionMaxSize" - (オプション) 書き込みに使用。デフォルト: 5242880 (バイト)。1 PutRecords 回のリクエストで送信するデータの最大量。

AWS Glue ストリーミングオプション

Kafka クラスターまたは Amazon Managed Streaming for Apache Kafka への接続を指定します。

Data Catalog テーブルに格納されている情報を使用するか、データストリームに直接アクセスするための情報を指定することにより、Kafka データストリームへ読み込むまたは書き込むことができます。Kafka から Spark に情報を読み取り DataFrame、それを Glue に変換できます。AWS DynamicFrameKafka には JSON DynamicFrames 形式で書き込むことができます。データストリームに直接アクセスする場合は、これらのオプションを使用して、データストリームへのアクセス方法に関する情報を提供します。

getCatalogSource または create_data_frame_from_catalog を使用して Kafka ストリーミングソースからレコードを消費するか、getCatalogSink または write_dynamic_frame_from_catalog を使用して Kafka にレコードを書き込み、ジョブに Data Catalog データベースおよびテーブル名の情報があり、それを使用して Kafka ストリーミングソースから読み込むためのいくつかの基本パラメータを取得できる場合。getSource、getCatalogSink、getSourceWithFormat、getSinkWithFormat、createDataFrame を使用する場合、ここで説明する接続オプションを使用し、これらの基本パラメータを指定する必要があります。

GlueContext クラスで指定されたメソッドの次の引数を使用し、Kafka の接続オプションを指定できます。

- Scala
 - connectionOptions: getSource、createDataFrameFromOptions、getSink で使用
 - additionalOptions: getCatalogSource、getCatalogSink で使用
 - options: getSourceWithFormat、getSinkWithFormat で使用
- Python
 - connection_options:
create_data_frame_from_options、write_dynamic_frame_from_options で使用
 - additional_options:
create_data_frame_from_catalog、write_dynamic_frame_from_catalog で使用
 - options: getSource、getSink で使用

ストリーミング ETL ジョブに関する注意事項と制限事項については、「[the section called “ストリーミング ETL に関する注意と制限”](#)」を参照してください。

AWS Glue Streaming 自動スケーリング

以下のセクションでは、AWS Glue Streaming 自動スケーリングについて説明します。

AWS Glue Studio の Auto-Scaling を有効にする

AWS Glue Studio の [Job details] (ジョブ詳細) タブで、タイプに [Spark] または [Spark Streaming] を選択し、[Glue version] (Glue バージョン) には **[Glue 3.0]** もしくは **[Glue 4.0]** を選択します。次に、[Worker type] (ワーカータイプ) の下にチェックボックスが表示されます。

- [Automatically scale the number of workers] (ワーカー数を自動的にスケーリングする) を選択します。
- ワーカーの最大数をセットして、ジョブ実行に投入できるワーカーの最大数を定義します。

Visual | Script | **Job details** | Runs | Data quality | Schedules

Version Control

Type
The type of ETL job. This is set automatically based on the types of data sources you have selected.

Spark

Glue version [Info](#)

Glue 4.0 - Supports spark 3.3, Scala 2, Python 3 ▼

Language

Python 3 ▼

Worker type
Set the type of predefined worker that is allowed when a job runs.

G 1X
(4vCPU and 16GB RAM) ▼

Automatically scale the number of workers
AWS Glue will optimize costs and resource usage by dynamically scaling the number of workers up and down throughout the job run. Requires Glue 3.0 or later.

Maximum number of workers
The number of workers you want AWS Glue to allocate to this job.

10

AWS CLI または SDK を使用した Auto Scaling の有効化

Auto Scaling を有効にするには、ジョブ実行の AWS CLI から、次の設定で `start-job-run` を実行します。

```
{
  "JobName": "<your job name>",
  "Arguments": {
    "--enable-auto-scaling": "true"
  },
  "WorkerType": "G.2X", // G.1X and G.2X are allowed for Auto Scaling Jobs
}
```

```
"NumberOfWorkers": 20, // represents Maximum number of workers
...other job run configurations...
}
```

ETL ジョブ実行が終了したら、`get-job-run` を呼び出して、実行されたジョブの実際のリソース使用量を DPU 秒単位で確認することもできます。注: 新しいフィールド [DPUSeconds] は、Auto Scaling が有効になっている AWS Glue 3.0 以降のパッチジョブに対してのみ表示されます。このフィールドは、ストリーミングジョブではサポートされません。

```
$ aws glue get-job-run --job-name your-job-name --run-id jr_xx --endpoint https://
glue.us-east-1.amazonaws.com --region us-east-1
{
  "JobRun": {
    ...
    "GlueVersion": "3.0",
    "DPUSeconds": 386.0
  }
}
```

同じ設定で [AWS Glue SDK](#) を使用して、Auto Scaling でジョブ実行を設定することもできます。

仕組み

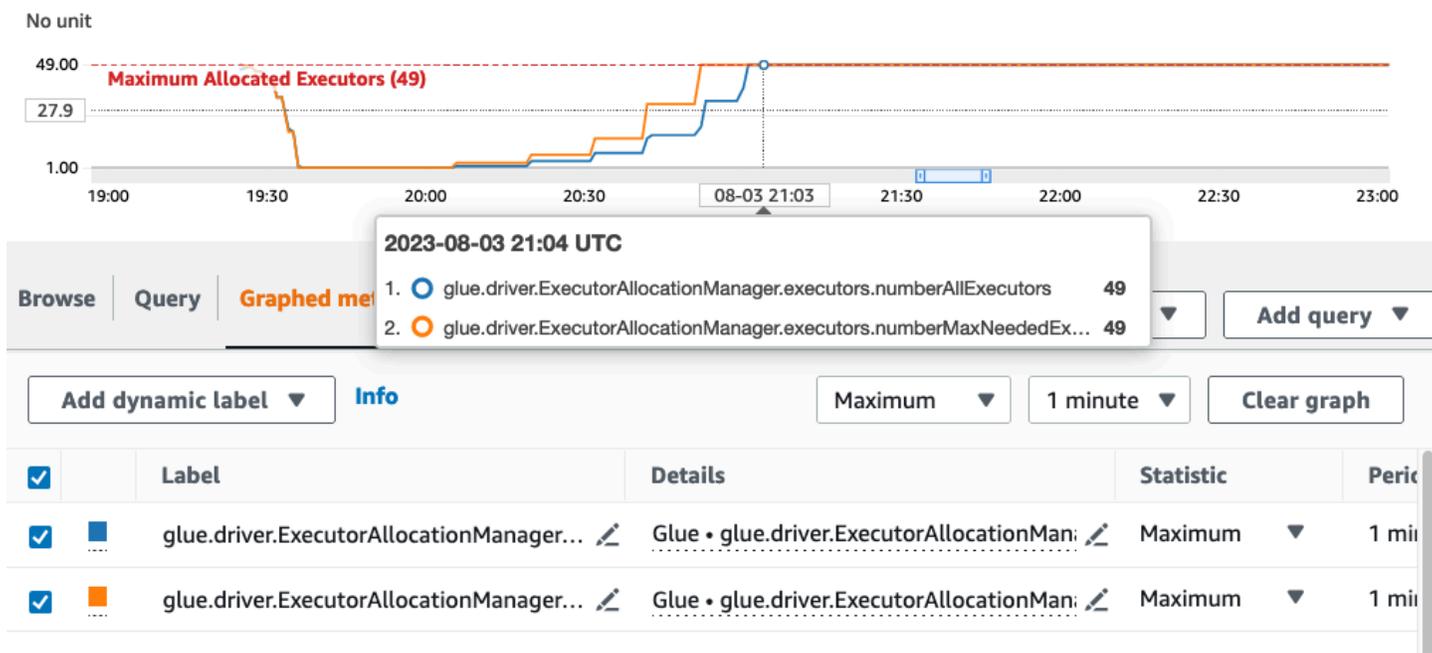
マイクロバッチ全体のスケーリング

次の例を使用して、自動スケーリングの仕組みを説明します。

- 50 DPU から始まる AWS Glue ジョブがあります。
- 自動スケーリングが有効になっています。

この例では、AWS Glue は、いくつかのマイクロバッチの「batchProcessingTimeInMs」メトリクスを確認し、設定されたウィンドウサイズ内でジョブが完了するかどうかを判断します。ジョブがすぐに完了する場合は、その完了するまでの時間によっては AWS Glue がスケールダウンすることもあります。「numberAllExecutors」でプロットされたこのメトリクスを Amazon CloudWatch でモニタリングして、自動スケーリングがどのように機能するかを確認できます。

エグゼキュターの数は、マイクロバッチが完了するたびに指数関数的にスケールアップまたはスケールダウンします。Amazon CloudWatch モニタリングログからわかるように、AWS Glue は、必要なエグゼキュターの数 (オレンジ色の線) を確認し、それに合わせてエグゼキュター (青い線) を自動的にスケーリングしています。



AWS Glue がエグゼキュターの数スケールダウンし、データ量が増加してマイクロバッチ処理時間が長くなることに気付くと、AWS Glue が指定された上限である 50 DPU にスケールアップします。

マイクロバッチ内でのスケーリング

上記の例では、システムが、完了したいいくつかのマイクロバッチをモニタリングして、スケールアップするかスケールダウンするかを決定しています。ウィンドウが長くなると、いくつかのマイクロバッチを待つよりも、自動スケーリングがマイクロバッチ内で速く応答する必要があります。このような場合は、追加の設定 `--auto-scale-within-microbatch` を `true` に指定して使用できます。以下に示すように、これを AWS Glue Studio で AWS Glue ジョブプロパティに追加できます。

Job parameters [Info](#)

Key Value - optional

You can add 49 more parameters.

AWS Glue ストリーミングのメンテナンスウィンドウ

AWS Glue は定期的にメンテナンスアクティビティを実行します。これらのメンテナンスウィンドウ中に、AWS Glue はストリーミングジョブを再起動する必要があります。メンテナンスウィンドウを指定することで、ジョブを再起動するタイミングを制御できます。このセクションでは、メンテナンスウィンドウをセットアップできる場所と、考慮すべき特定の動作の概要を説明します。

トピック

- [メンテナンスウィンドウの設定](#)
- [メンテナンスウィンドウの動作](#)
- [ジョブのモニタリング](#)
- [データ損失処理](#)

メンテナンスウィンドウの設定

AWS Glue Studio または APIs を使用してメンテナンスウィンドウを設定できます。

AWS Glue Studio でのメンテナンスウィンドウの設定

AWS Glue ストリーミングジョブのジョブ詳細ページでメンテナンスウィンドウを指定できます。GMT で日時を指定できます。AWS Glue は指定された時間枠内にジョブを再起動します。

Maintenance window

Restart on

at hours (GMT)

For maintenance reasons, AWS Glue will restart streaming jobs within 3 hours of the specified maintenance window. You have the option to designate the start time in GMT for this maintenance. For more information, refer to documentation.

API でのメンテナンスウィンドウの設定

代わりに、ジョブの作成 API でメンテナンスウィンドウを設定できます。API を使用してメンテナンスウィンドウを設定する例を次に示します。

```
aws glue create-job --name jobName --role roleArnForTheJob --command
Name=gluestreaming,ScriptLocation=s3-path-to-the-script --maintenance-window="Sun:10"
```

コマンドの例は次のとおりです。

```
aws glue create-job --name testMaintenance --role arn:aws:iam::012345678901:role/
Glue_DefaultRole --command Name=gluestreaming,ScriptLocation=s3://glue-example-test/
example.py --maintenance-window="Sun:10"
```

メンテナンスウィンドウの動作

AWS Glue は、ジョブを再起動するタイミングを決定するために、一連のステップを実行します。

1. 新しいストリーミングジョブが開始されると、AWS Glue まずジョブ実行に関連付けられたタイムアウトがあるかどうかをチェックします。タイムアウトを使用すると、ジョブの終了時刻を設定できます。タイムアウトが7日未満の場合、ジョブは再起動されません。
2. タイムアウトが7日を超える場合、はメンテナンスウィンドウがジョブに設定されている AWS Glue かどうかを確認します。そのウィンドウが取得され、そのウィンドウがジョブ実行に割り当てられると、は指定されたメンテナンスウィンドウから3時間以内にジョブ AWS Glue を再開します。例えば、月曜日の GMT 午前 10 時にメンテナンスウィンドウを設定すると、ジョブは GMT 午前 10 時から GMT 午後 1 時までの間に再開されます。
3. メンテナンスウィンドウが設定されていない場合、は再起動時間をジョブ実行開始時刻の7日前 AWS Glue に自動的に設定します。例えば、ジョブを 7/1/2024 午前 12 時に開始し、メンテナンスウィンドウを指定しなかった場合、ジョブは 7/8/2024 午前 12 時に再起動するように設定されます。

Note

既にストリーミングジョブを実行している場合、この変更は 2024 年 7 月 1 日以降に影響します。6 月 30 日までメンテナンスウィンドウを設定する時間があります。7 月 1 日以降、開始したストリーミングジョブは、このドキュメントに従って再起動されます。追加のサポートが必要な場合は、AWS サポートにお問い合わせください。

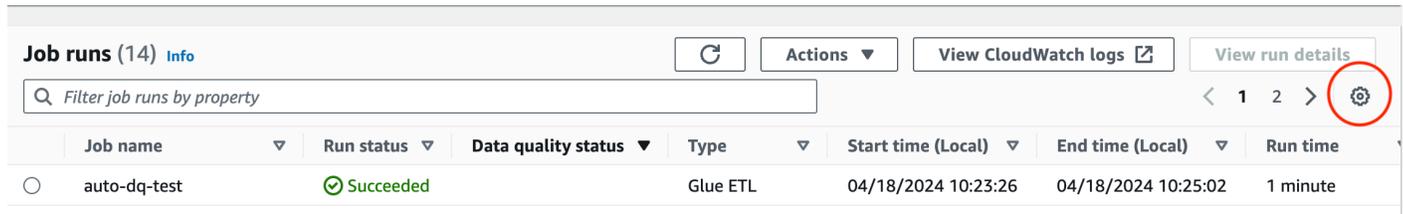
4. 特に進行中のマイクロバッチが処理されない場合、ジョブを再起動できない AWS Glue ことがあります。このような場合、ジョブは中断されません。これらのインスタンス AWS Glue では、は 14 日後にジョブを再起動します。この場合、メンテナンスウィンドウは適用されません。

ジョブのモニタリング

AWS Glue Studio Monitoring ページでジョブをモニタリングできます。

ストリーミングジョブの次回の再起動予定時刻を確認するには、モニタリングページのジョブ実行テーブルに列を表示します。

1. テーブルの右上にある歯車アイコンをクリックします。



The screenshot shows the 'Job runs (14) Info' section of the AWS Glue Studio Monitoring page. It includes a search bar, a refresh button, and buttons for 'Actions', 'View CloudWatch logs', and 'View run details'. Below these is a table with columns for Job name, Run status, Data quality status, Type, Start time (Local), End time (Local), and Run time. The first row shows a job named 'auto-dq-test' with a status of 'Succeeded'. A red circle highlights the gear icon in the top right corner of the table area.

Job name	Run status	Data quality status	Type	Start time (Local)	End time (Local)	Run time
○ auto-dq-test	✔ Succeeded		Glue ETL	04/18/2024 10:23:26	04/18/2024 10:25:02	1 minute

2. 下にスクロールして、予想される再起動時間列をオンにします。UTC およびローカル時間オプションの両方を使用できます。

worker type	<input type="checkbox"/>
DPU hours	<input checked="" type="checkbox"/>
Last modified (Local)	<input type="checkbox"/>
Worker utilization	<input type="checkbox"/>
Data skewness	<input type="checkbox"/>
Start time (UTC)	<input type="checkbox"/>
End time (UTC)	<input type="checkbox"/>
Last modified (UTC)	<input type="checkbox"/>
Data quality	<input checked="" type="checkbox"/>
Expected restart time (UTC)	<input checked="" type="checkbox"/>
Expected restart time (Local)	<input type="checkbox"/>

Cancel
Confirm

3. その後、テーブルの列を表示できます。

Job runs (14) [Info](#) 🔄 Actions ▼ View CloudWatch logs 📄 View run details

🔍 Filter job runs by property < 1 2 > ⚙️

	Job name ▼	Run status ▼	Type ▼	Start time (Local) ▼	End time (Local) ▼	Expected restart time (Local) ▼
<input type="radio"/>	auto-dq-test	✔️ Succeeded	Glue ETL	04/18/2024 10:23:26	04/18/2024 10:25:02	-
<input type="radio"/>	StreamingTest	🔄 Running	Glue Streaming	04/16/2024 16:32:49	-	04/23/2024 02:00:00
<input type="radio"/>	StreamingProd	🔄 Running	Glue Streaming	04/16/2024 13:45:10	-	04/25/2024 05:00:00

元のジョブのステータスは「EXPIRED」になり、新しいジョブインスタンスのステータスは「RUNNING」になります。再起動された新しいジョブ実行には、最初のジョブ実行 ID の連結としてジョブ実行 ID と、再起動数を表すプレフィックス「restart_」が付けられます。例えば、最初のジョブ実行 ID が の場合jr_1234、再起動されたジョブ実行には最初の再起動jr1234_restart_1の ID が含まれます。2 回目の再起動は 2 jr1234_restart_2 回目の再起動と同じです。

再試行は、再起動による影響を受けません。自動再試行により実行が失敗し、新しい実行が開始された場合、再起動のカウンターは 1 から再び開始されます。例えば、 で実行が失敗した場合jr_1234_attempt_3_restart_5、自動再試行は ID で新しい実行を開始します。この試行が 7 日後に再開jr_id1_attempt_4されると、新しい実行 ID は になりますjr_id1_attempt_4_restart_1。

データ損失処理

メンテナンスの再起動中、AWS Glue ストリーミングは、前のジョブ実行と再起動されたジョブ実行の間のデータの整合性と一貫性を確保するプロセスに従います。AWS Glue は、ジョブの再起動間のデータの整合性と一貫性を保証するものではないことに注意してください。ストリーミングジョブ内で重複するデータを処理するために、アーキテクチャ上の考慮事項をお勧めします。

1. メンテナンスの再起動条件の検出: AWS Glue Streaming は、7 日後にメンテナンスウィンドウに達した場合や、14 日後にハード再起動が必要な場合など、メンテナンスの再起動をトリガーする必要があるタイミングを示す条件をモニタリングします。
2. 正常な終了の呼び出し: メンテナンスの再起動条件が満たされると、AWS Glue Streaming は現在実行中のジョブに対して正常な終了プロセスを開始します。このプロセスには次のステップが含まれます。
 - a. 新しいデータの取り込みの停止: ストリーミングジョブは、入力ソース (Kafka トピック、Kinesis ストリーム、ファイルなど) からの新しいデータの消費を停止します。
 - b. 保留中のデータの処理: ジョブは、内部バッファまたはキューに既に存在するデータを引き続き処理します。
 - c. オフセットとチェックポイントのコミット: ジョブは、最新のオフセットまたはチェックポイントを外部システム (Kafka、Kinesis、Amazon S3 など) にコミットして、再起動されたジョブが前のジョブが中断した場所から確実に取得できるようにします。
3. ジョブの再起動: 正常な終了プロセスが完了すると、AWS Glue Streaming は保持された状態とチェックポイントを使用してジョブを再起動します。再起動されたジョブは、最後にコミットされたオフセットまたはチェックポイントから処理を取得し、データが失われたり重複したりしないようにします。

4. データ処理の再開: 再開されたジョブは、前のジョブが中断した時点からデータ処理を再開します。最後にコミットされたオフセットまたはチェックポイントから入力ソースからの新しいデータの取り込みを継続し、定義された ETL ロジックに従ってデータを処理します。

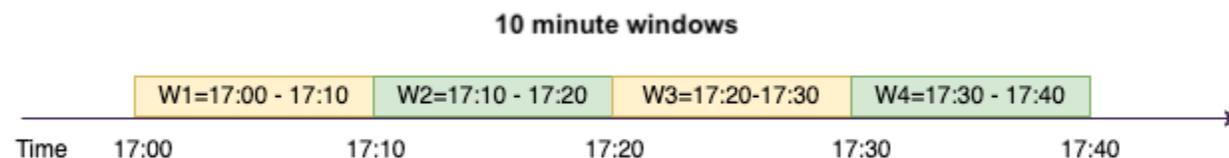
高度な AWS Glue のストリーミングの概念

現代のデータ駆動型のアプリケーションでは、データの重要性は時間の経過と共に低下し、予測データから事後的データへとその価値が移行します。そのため、お客様はデータをリアルタイムで処理して、より迅速な意思決定を行いたいと考えています。IoT センサーなどからのリアルタイムのデータフィードを処理する場合、取り込み中のネットワーク遅延やその他のソース関連の障害により、データの到着が順序付けされていないことや、処理に遅延が発生することがあります。AWS Glue プラットフォームの一部である AWS Glue Streaming は、これらの機能を基盤として、Apache Spark 構造化ストリーミングを利用したスケーラブルでサーバーレスなストリーミング ETL を提供し、ユーザーがリアルタイムでデータを処理できるようにします。

このトピックでは、高度なストリーミングの概念と AWS Glue Streaming の機能について説明します。

ストリームを処理する際の時間に関する考慮事項

ストリームを処理する際の時間には、4 つの概念があります。



- [イベント時間] – イベントが発生した時間。ほとんどの場合、このフィールドはソースでイベントデータ自体に埋め込まれています。
- Event-time-window – 2 つのイベント時間の間の時間枠。上記の図に示すように、W1 は 17:00 event-time-window から 17:10 までのです。各 event-time-window は複数のイベントのグループです。
- [トリガー時間] – トリガー時間により、データの処理と結果の更新が行われる頻度が制御されます。これは、マイクロバッチ処理が開始された時間です。
- [取り込み時間] – ストリームデータがストリーミングサービスに取り込まれた時間。イベント時間がイベント自体に組み込まれていない場合、この時間をウィンドウ処理に使用できる場合があります。

ウィンドウ処理

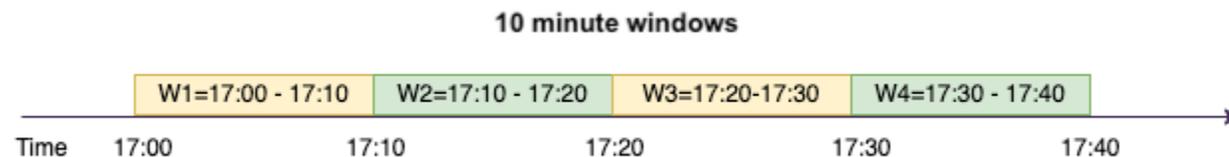
ウィンドウイングは、複数のイベントをグループ化して集計する手法です event-time-window。以下の例で、ウィンドウ処理の利点と、どのような場合に使用するかを説明します。

ビジネスユースケースに応じて、Spark では 3 種類のタイムウィンドウがサポートされます。

- タンブリングウィンドウ — 集計 event-time-windows する一連の重複しない固定サイズ。
- スライディングウィンドウ - 「固定サイズ」であるという点ではタンブリングウィンドウに似ていますが、スライドの継続時間がウィンドウ自体の継続時間よりも短い限り、ウィンドウが重なり合ったりスライドしたりできます。
- セッションウィンドウ - 入力データイベントで開始され、一定の間隔または非アクティブ時間内に入力が受信されている限り、長くなり続けます。セッションウィンドウは、入力に応じて、ウィンドウの長さが静的サイズまたは動的サイズになります。

タンブリングウィンドウ

タンブリングウィンドウは、重複しない一連の固定サイズで、これ event-time-windows に対して集計されます。実際の例を使って理解してみましょう。



ABC Auto 社は、スポーツカーの新しいブランドのマーケティングキャンペーンを行いたいと考えています。彼らは、スポーツカーのファンが最も多い都市を選びたいと考えています。この目標を達成するために、ウェブサイトでは車を紹介する 15 秒間の短い広告を掲載しています。すべての「クリック」と対応する「都市」が記録され、にストリーミングされます Amazon Kinesis Data Streams。10 分のウィンドウでクリック数をカウントし、都市別にグループ化して、需要が最も多い都市を確認したいと考えています。この集計の出力を次に示します。

window_start_time	window_end_time	city	total_clicks
2023-07-10 17:00:00	2023-07-10 17:10:00	ダラス	75
2023-07-10 17:00:00	2023-07-10 17:10:00	シカゴ	10

window_start_time	window_end_time	city	total_clicks
2023-07-10 17:20:00	2023-07-10 17:30:00	ダラス	20
2023-07-10 17:20:00	2023-07-10 17:30:00	シカゴ	50

上記で説明したように、これらはトリガー時間間隔とは異なり event-time-windows ます。例えば、トリガー時間が 1 分ごとであっても、出力結果では 10 分の重なり合わない集計ウィンドウのみが表示されます。最適化のために、トリガー間隔を に合わせることをお勧めします event-time-window。

上の表で、17:00 ~ 17:10 のウィンドウでダラスでは 75 回のクリックがあり、シカゴでは 10 回のクリックがありました。また、どの都市についても 17:10 ~ 17:20 のウィンドウのデータがないため、このウィンドウは省略されています。

これで、ダウンストリーム分析アプリケーションでこのデータをさらに分析して、マーケティングキャンペーンを実施するのに最も適した都市を特定できます。

AWS Glue でのタンブリングウィンドウの使用

1. Amazon Kinesis Data Streams DataFrame を作成し、そこから読み取ります。例：

```
parsed_df = kinesis_raw_df \  
    .selectExpr('CAST(data AS STRING)') \  
    .select(from_json("data", ticker_schema).alias("data")) \  
    .select('data.event_time', 'data.ticker', 'data.trade', 'data.volume',  
    'data.price')
```

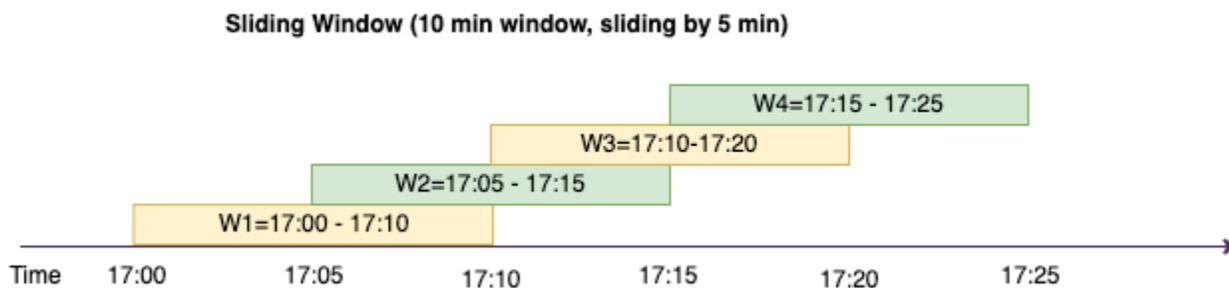
2. タンブリングウィンドウでデータを処理します。以下の例では、入力フィールド「event_time」に基づいて 10 分のタンブリングウィンドウでデータがグループ化され、その出力が Amazon S3 データレイクに書き込まれます。

```
grouped_df = parsed_df \  
    .groupBy(window("event_time", "10 minutes"), "city") \  
    .agg(sum("clicks").alias("total_clicks"))  
  
summary_df = grouped_df \  
    .withColumn("window_start_time", col("window.start")) \  
    .withColumn("window_end_time", col("window.end")) \  
    .withColumn("year", year("window_start_time")) \  
    .withColumn("month", month("window_start_time"))
```

```
.withColumn("month", month("window_start_time")) \  
.withColumn("day", dayofmonth("window_start_time")) \  
.withColumn("hour", hour("window_start_time")) \  
.withColumn("minute", minute("window_start_time")) \  
.drop("window")  
  
write_result = summary_df \  
    .writeStream \  
    .format("parquet") \  
    .trigger(processingTime="10 seconds") \  
    .option("checkpointLocation", "s3a://bucket-stock-stream/stock-  
stream-catalog-job/checkpoint/") \  
    .option("path", "s3a://bucket-stock-stream/stock-stream-catalog-  
job/summary_output/") \  
    .partitionBy("year", "month", "day") \  
    .start()
```

スライディングウィンドウ

スライディングウィンドウは「固定サイズ」であるという点ではタンブリングウィンドウに似ていますが、スライドの継続時間がウィンドウ自体の継続時間よりも短い限り、ウィンドウが重なり合ったりスライドしたりできます。スライディングの性質上、1つの入力を複数のウィンドウにバインドできます。



理解を深めるために、クレジットカードの不正使用の可能性を検知したい銀行の例を考えてみましょう。ストリーミングアプリケーションでは、クレジットカードの取引の流れを継続的にモニタリングできます。これらの取引が10分間のウィンドウに集計され、5分ごとにウィンドウが先にスライドされて、最も古い5分間のデータが削除され、最新の5分間の新しいデータが追加されます。各ウィンドウ内で、取引が国ごとにグループ化され、疑わしいパターンがないかチェックできます。例えば、米国での取引の直後にオーストラリアで発生した取引などです。わかりやすくするために、取引総額が100 USDを超える場合は、その取引を不正使用として分類してみましょう。このようなパ

ターンが検出された場合は、不正使用の可能性があることが示され、カードが凍結される可能性があります。

クレジットカード処理システムは、国と共に各カード ID の取引イベントのストリームを Kinesis に送信しています。AWS Glue ジョブは分析を実行し、次の集計出力を生成します。

window_start_time	window_end_time	card_last_four	country	total_amount
2023-07-10 17:00:00	2023-07-10 17:10:00	6544	米国	85
2023-07-10 17:00:00	2023-07-10 17:10:00	6544	オーストラリア	10
2023-07-10 17:05:45	2023-07-10 17:15:45	6544	米国	50
2023-07-10 17:10:45	2023-07-10 17:20:45	6544	米国	50
2023-07-10 17:10:45	2023-07-10 17:20:45	6544	オーストラリア	150

上記の集計に基づき、取引金額で合計された、5 分ごとにスライドする 10 分のウィンドウが表示されます。外れ値 (オーストラリアでの 150 USD の取引) がある 17:10 ~ 17:20 のウィンドウで異常が検出されます。AWS Glue は、この異常を検出し、boto3 を使用して問題のキーを含むアラームイベントを SNS トピックにプッシュすることができます。さらに、Lambda 関数はこのトピックをサブスクライブしてアクションを実行できます。

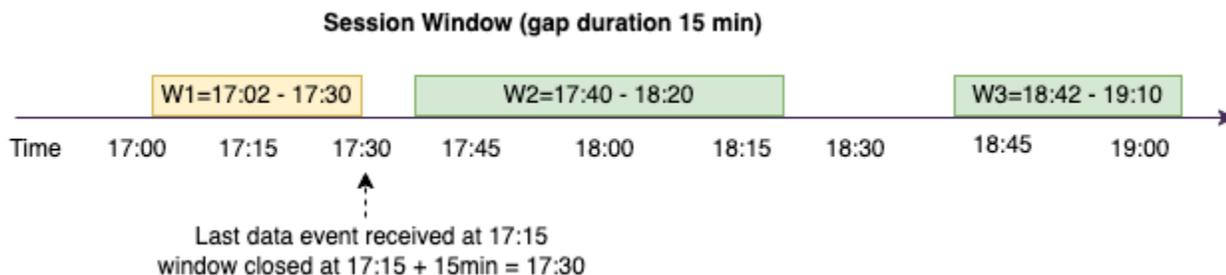
スライディングウィンドウでのデータの処理

以下のように、group-by 句とウィンドウ関数を使用してスライディングウィンドウを実装します。

```
grouped_df = parsed_df \  
    .groupBy(window(col("event_time"), "10 minute", "5 min"), "country",  
    "card_last_four") \  
    .agg(sum("tx_amount").alias("total_amount"))
```

セッションウィンドウ

固定サイズである上記の 2 つのウィンドウとは異なり、セッションウィンドウは、入力に応じて、ウィンドウの長さが静的サイズまたは動的サイズになります。セッションウィンドウは入力データイベントで開始され、一定の間隔または非アクティブ時間内に入力が受信されている限り、長くなり続けます。



例を見てみましょう。ABC hotel 社は、1 週間のうちで最も混雑する時間帯を調べて、宿泊客により良い取引を提供したいと考えています。ゲストがチェックインするとすぐにセッションウィンドウが開始され、Spark はその集約を含む状態を維持します event-time-window。ゲストがチェックインするたびに、イベントが生成され、に送信されます Amazon Kinesis Data Streams。15 分間チェックインがない場合、を閉鎖 event-time-window できるという決定がホテルによって下されます。次の は、新しいチェックインがあると再び開始 event-time-window されます。出力は次のとおりです。

window_start_time	window_end_time	city	total_checkins
2023-07-10 17:02:00	2023-07-10 17:30:00	ダラス	50
2023-07-10 17:02:00	2023-07-10 17:30:00	シカゴ	25
2023-07-10 17:40:00	2023-07-10 18:20:00	ダラス	75
2023-07-10 18:50:45	2023-07-10 19:15:45	ダラス	20

最初のチェックインは、event_time=17:02 に行われました。集計は 17:02 に開始 event-time-window されます。この集計は、15 分以内にイベントが受信される限り継続されます。上の例では、17:15 に最後のイベントが受信され、その後の 15 分間はイベントがありませんでした。その結果、Spark

は 17:15+15min = 17:30 event-time-window にこれを閉じ、17:02 ~ 17:30 に設定します。新しいチェックインデータイベントを受け取った event-time-window 17:47 に新しい が開始されました。

セッションウィンドウでのデータの処理

group-by 句とウィンドウ関数を使用してスライディングウィンドウを実装します。

```
grouped_df = parsed_df \  
    .groupBy(session_window(col("event_time"), "10 minute"), "city") \  
    .agg(count("check_in").alias("total_checkins"))
```

出力モード

出力モードとは、無制限テーブルからの結果が外部シンクに書き込まれるときのモードです。3つのモードを使用できます。次の例では、各マイクロバッチでデータの行がストリーミングおよび処理されるときに、単語の出現回数をカウントしています。

- 完了モード – 現在の で単語数が更新されていなくても、マイクロバッチ処理のたびに結果テーブル全体がシンクに書き込まれます event-time-window。
- 付加モード – デフォルトのモードで、最後のトリガー以降に結果テーブルに追加された新しい単語や行のみがシンクに書き込まれます。このモードは、map、flatMap、filter などのクエリのステートレスストリーミングに適しています。
- 更新モード – 最後のトリガー以降に更新または追加された結果テーブルの単語や行のみがシンクに書き込まれます。

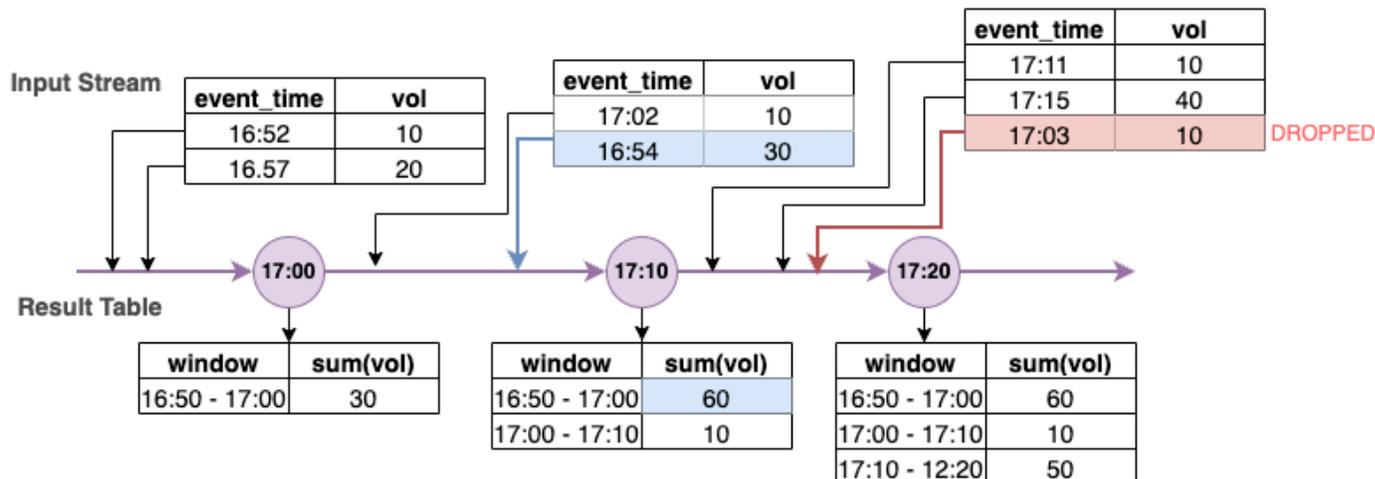
Note

セッションウィンドウでは「更新」の出力モードはサポートされていません。

遅延データとウォーターマークの処理

リアルタイムデータを操作する場合、ネットワークレイテンシーやアップストリーム障害が原因でデータの到着が遅れる可能性があり、失われた に対して集約を再度実行するメカニズムが必要です event-time-window。ただし、そのためにはステートを維持する必要があります。同時に、ステートのサイズを制限するために、古いデータをクリーンアップする必要があります。Spark バージョン 2.1 では、ステートを維持し、遅延データのしきい値をユーザーが指定することができる、ウォーターマークと呼ばれる機能のサポートが追加されました。

前述の株価ティッカーの例について、遅延データの許容しきい値が 10 分以内であると考えてみましょう。単純化のために、タンプリングウィンドウ、ティッカーは AMZ、取引を BUY と仮定します。



上の図では、10 分間のタンプリングウィンドウで合計株数を計算しています。17:00、17:10、および 17:20 にトリガーされています。タイムラインの矢印の上には入力データストリームがあり、下には無制限の結果のテーブルがあります。

最初の 10 分間のタンプリングウィンドウでは、event_time に基づいて集計され、total_volume は 30 と計算されました。2 番目のでは event-time-window、spark は event_time=17:02 の最初のデータイベントを取得しました。これは Spark でこれまでに確認された event_time の最大値であるため、ウォーターマークのしきい値がこの 10 分前に設定されます (つまり、watermark_event_time=16:52)。event_time が 16:52 以降のデータイベントは時間が制限された集計で考慮され、それより前のデータイベントはドロップされます。これにより、Spark はさらに 10 分間中間状態を維持して、遅延データを受け入れることができます。実時間で 17:08 ごろに Spark が event_time=16:54 のイベントを受信し、これはしきい値内でした。したがって、Spark は「16:50 ~ 17:00」を再計算 event-time-window し、合計ボリュームが 30 から 60 に更新されました。

しかし、トリガー時間 17:20 で、Spark が event_time=17:15 のイベントを受信したときに watermark_event_time が 17:05 に設定されます。そのため、event_time=17:03 の遅延データイベントは「遅すぎる」と見なされ、無視されました。

$$\text{Watermark Boundary} = \text{Max(Event Time)} - \text{Watermark Threshold}$$

AWS Glue でのウォーターマークの使用

Spark は、ウォーターマークの境界が越えられるまで、外部シンクへのデータの出力または書き込みを行いません。AWS Glue でウォーターマークを実装するには、以下の例を参照してください。

```
grouped_df = parsed_df \  
    .withWatermark("event_time", "10 minutes") \  
    .groupBy(window("event_time", "5 minutes"), "ticker") \  
    .agg(sum("volume").alias("total_volume"))
```

AWS Glue Streaming ジョブのモニタリング

ストリーミングジョブのモニタリングは、ETL パイプラインを構築するうえで極めて重要です。Spark UI を使用する以外に、Amazon CloudWatch を使用してメトリクスをモニタリングすることもできます。以下は、AWS Glue フレームワークによって出力されるストリーミングメトリクスのリストです。すべての AWS Glue メトリクスが含まれる完全なリストについては、「[Amazon CloudWatch メトリクスを使用した AWS Glue のモニタリング](#)」を参照してください。

AWS Glue では、構造化されたストリーミングフレームワークを使用して入力イベントを処理します。Spark API をコード内で直接使用することも、これらのメトリクスを発行する GlueContext で提供された ForEachBatch を利用することもできます。これらのメトリクスを理解するには、まず windowSize を理解する必要があります。

windowSize: windowSize は、指定するマイクロバッチの間隔です。ウィンドウサイズを 60 秒に指定すると、AWS Glue Streaming ジョブは 60 秒 (それまでに前のバッチが完了していない場合はそれ以上) 待ってから、ストリーミングソースからバッチのデータを読み取り、ForEachBatch で提供される変換を適用します。これはトリガー間隔とも呼ばれます。

メトリクスを詳しく見直して、ヘルスとパフォーマンスの特徴を理解しましょう。

Note

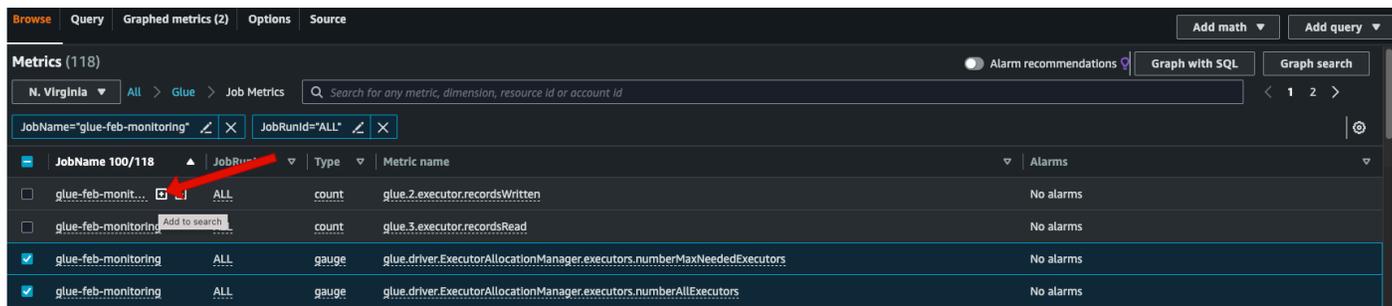
メトリクスは 30 秒ごとに出力されます。windowSize が 30 秒未満の場合、報告されるメトリクスは集計されたものとなります。例えば、windowSize が 10 秒で、マイクロバッチあたり 20 件のレコードを安定して処理しているとします。このシナリオでは、numRecords の出力メトリクス値は 60 になります。

メトリクスは、使用できるデータがない場合は出力されません。また、コンシューマーラグメトリクスの場合は、そのメトリクスを取得する機能を有効にする必要があります。

メトリクスの視覚化

ビジュアルメトリクスをプロットするには:

1. Amazon CloudWatch コンソールの [メトリクス] に移動し、[参照] タブを選択します。次に、「カスタム名前空間」の [Glue] を選択します。



2. [ジョブメトリクス] を選択すると、すべてのジョブのメトリクスが表示されます。
3. JobName=glue-feb-monitoring に基づいてメトリクスをフィルタリングし、次に JobRunId=ALL に基づいてフィルタリングします。下図のように「+」記号をクリックすると、検索フィルターに追加できます。
4. 対象のメトリクスのチェックボックスを選択します。以下の図では、numberAllExecutors と numberMaxNeededExecutors を選択しています。



5. これらのメトリクスを選択したら、[グラフ化したメトリクス] タブに移動して統計を適用できます。
6. メトリクスは 1 分ごとに出力されるため、batchProcessingTimeInMs および maxConsumerLagInMs には 1 分より長い「average」を適用できます。numRecords には、1 分ごとの「sum」を適用できます。
7. [オプション] タブでは、水平方向の windowSize の注釈をグラフに追加できます。

Browse Query Graphed metrics (1) Options Source

Widget type

Line Stacked area Number Gauge Bar Pie

Legend position

Hidden Bottom Right

Live data

Display most recent data point, even when not yet fully aggregated.

Left Y axis

Label milliseconds

Limits Min Auto Max Auto

Show units

Right Y axis

Label Add custom

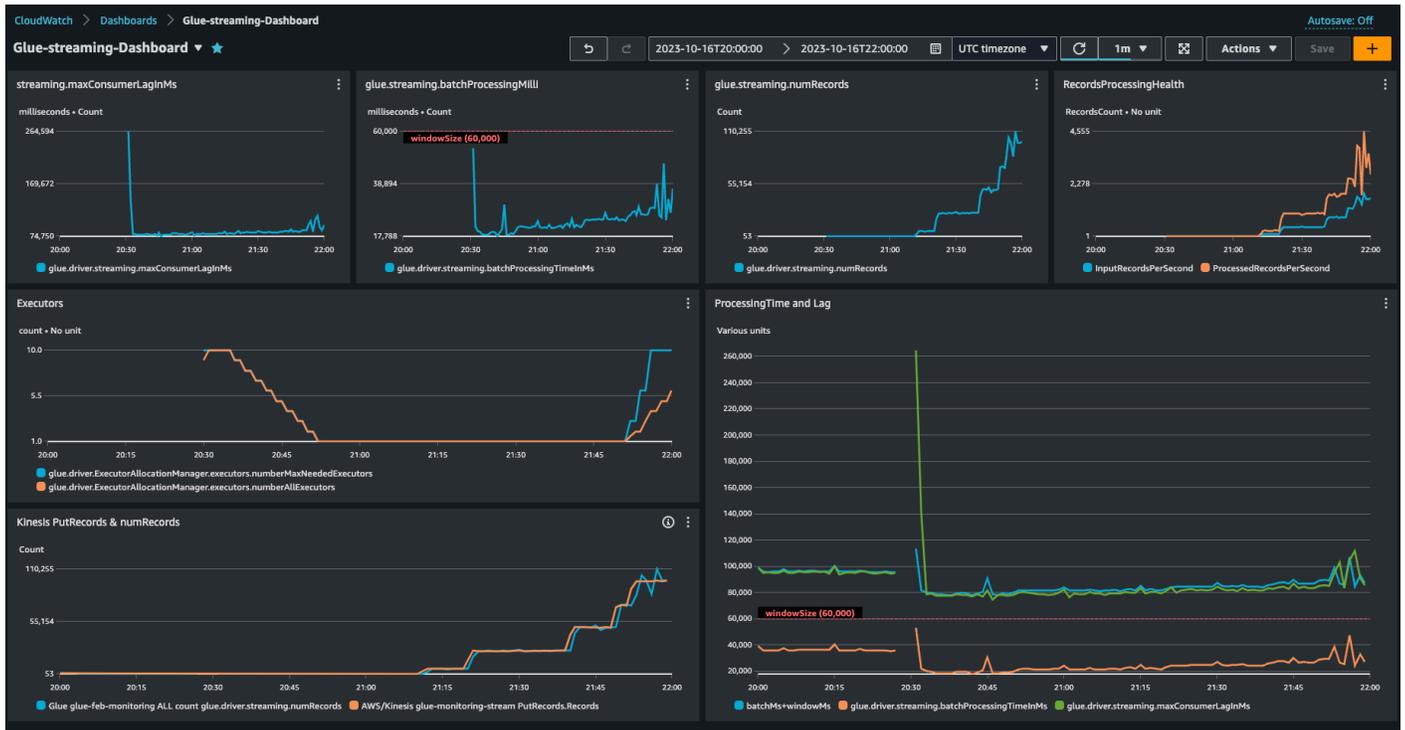
Limits Min Auto Max Auto

Show units

Horizontal annotations / thresholds - New ⓘ

	Label	Value	Fill	Axis	Actions
<input checked="" type="checkbox"/>	windowSize	60000	None	< >	×

8. メトリクスを選択したら、ダッシュボードを作成して追加します。ダッシュボードの例を以下に示します。



メトリクスの詳細

このセクションでは、各メトリクスと、それらが相互にどう関連しているかについて説明します。

レコード数 (メトリクス: streaming.numRecords)

このメトリクスは、処理中のレコードの数を示します。



このストリーミングメトリクスにより、処理中のレコードの数をウィンドウ内で可視化できます。処理中のレコードの数だけでなく、入カトラフィックの動作を把握するのにも役立ちます。

- インジケータ #1 は、トラフィックが安定しており、バーストが発生していない例を示しています。一般に、これは定期的にデータを収集してストリーミングソースに送信する IoT センサーのようなアプリケーションです。
- インジケータ #2 は、他の負荷が安定しているのに、トラフィックで突然バーストが発生した例を示しています。これは、ブラックフライデーのようなマーケティングイベントがあり、クリック数が急増した場合に、クリックストリームアプリケーションで発生することがあります。
- インジケータ #3 は、予測不可能なトラフィックの例を示しています。予想できないトラフィックは問題があることを指しません。これは入力データの性質にすぎません。IoT センサーの例に話を戻すと、何百ものセンサーが気象変化イベントをストリーミングソースに送信していると考えることができます。気象の変化は予測できないため、データも予測できません。トラフィックパターンを理解することは、エグゼキューターのサイズを決定するうえで重要です。入力が非常に多い場合は、自動スケーリングの使用を検討してください (これについては後で詳しく説明します)。



このメトリクスを Kinesis PutRecords メトリクスと組み合わせて、取り込まれるイベントの数と読み取られるレコードの数がほぼ同じになるようにすることができます。これは特に、ラグについて把握したい場合に便利です。取り込み率が高くなると、AWS Glue によって読み取られる numRecords も増加します。

バッチ処理時間 (メトリクス: streaming.batchProcessingTimeInMs)

バッチ処理時間メトリクスは、クラスターがプロビジョニング不足かプロビジョニング過剰かを判断するのに役立ちます。

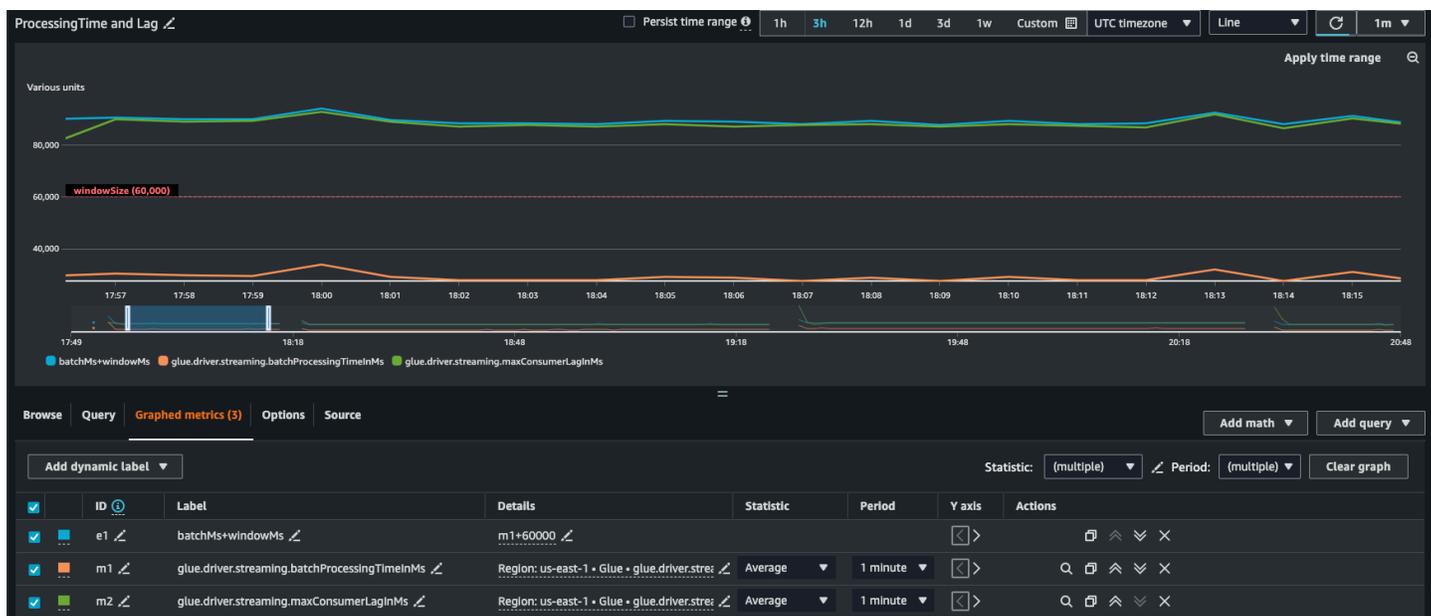


このメトリクスは、レコードの各マイクロバッチの処理にかかったミリ秒数を示します。ここでの主な目的は、この時間をモニタリングして windowSize の間隔を下回るようにすることです。次のウィンドウ間隔で回復するのであれば、batchProcessingTimeInMs が一時的にオーバーしても問題ありません。インジケータ #1 は、ジョブの処理にかかった時間がほぼ安定していることを示しています。ただし、入力レコードの数が増えている場合は、インジケータ #2 で示されているように、ジョブの処理にかかる時間も長くなります。numRecords が増えていない

のに処理時間が長くなっている場合は、エグゼキューターのジョブ処理を詳しく調べる必要があります。batchProcessingTimeInMs が 10 分より長く 120% を超えることがないように、しきい値とアラームを設定することをお勧めします。アラームの設定の詳細については、「[Amazon CloudWatch でのアラームの使用](#)」を参照してください。

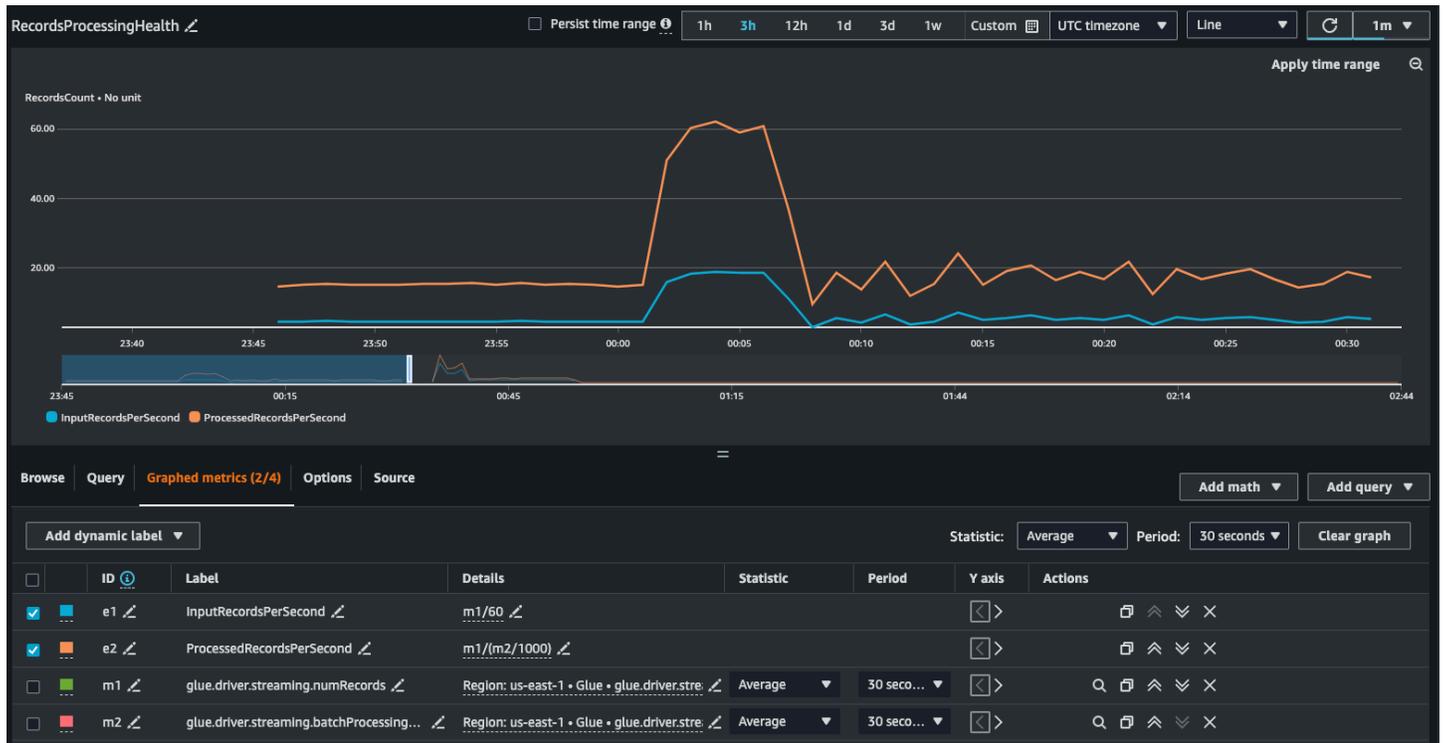
コンシューマーラグ (メトリクス: streaming.maxConsumerLagInMs)

コンシューマーラグメトリクスは、イベントの処理にラグがあるかどうかを把握するのに役立ちます。適切な windowSize を設定していても、ラグが大きすぎると、使用している処理 SLA が満たされない可能性があります。このメトリクスは、emitConsumerLagMetrics 接続オプションを使用して明示的に有効にする必要があります。詳細については、「[KinesisStreamingSourceOptions](#)」を参照してください。



派生メトリクス

より深い洞察を得るために、派生メトリクスを作成して、Amazon CloudWatch のストリーミングジョブについて理解を深めることができます。



派生メトリクスが含まれるグラフを作成して、DPU をさらに使用する必要があるかどうかを判断できます。自動スケーリングはこれを自動的に行うのに役立ちますが、派生メトリクスを使用することで、自動スケーリングが効果的に機能しているかどうかを判断できます。

- InputRecordsPerSecond は、入力レコードを取得する速度を示します。これは次のように算出されます: $\text{入力レコードの数 (glue.driver.streaming.numRecords)} / \text{WindowSize}$ 。
- ProcessingRecordsPerSecond は、レコードを処理する速度を示します。これは次のように算出されます: $\text{入力レコードの数 (glue.driver.streaming.numRecords)} / \text{batchProcessingTimeInMs}$ 。

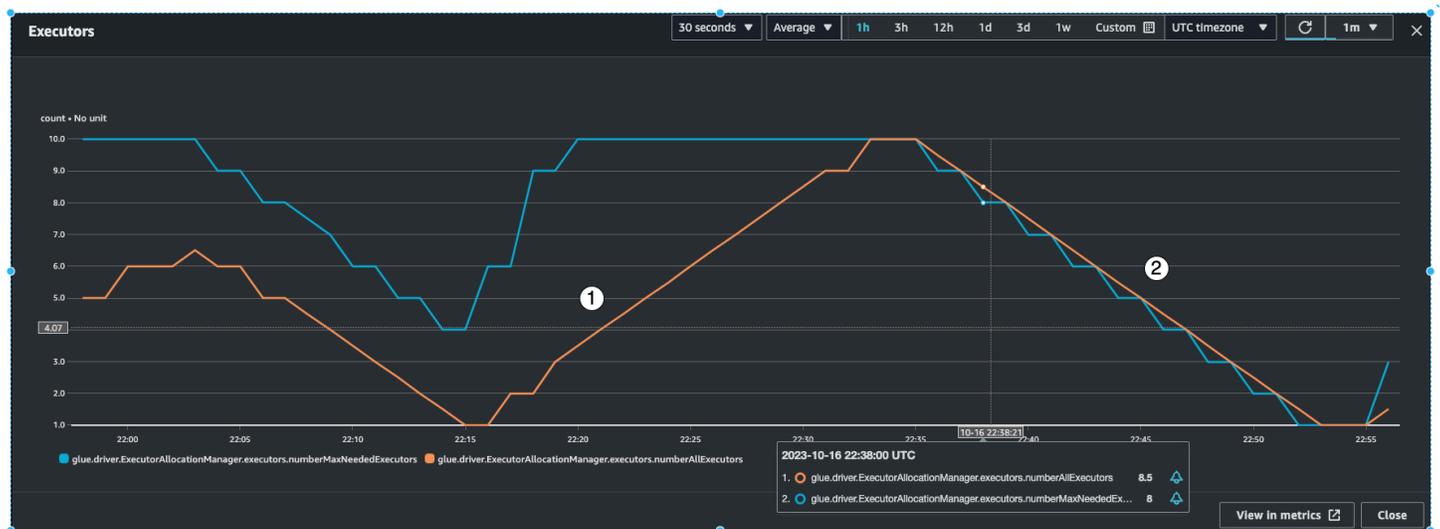
入力速度が処理速度より高いときは、ジョブを処理するための容量を増やすか、並列処理を増やす必要がある場合があります。

自動スケーリングメトリクス

入カトラフィックが急増しているときは、自動スケーリングを有効にすることを検討し、最大ワーカー数を指定する必要があります。そのために、numberAllExecutors と numberMaxNeededExecutors という 2 つのメトリクスが用意されています。

- numberAllExecutors は、アクティブに実行されているジョブエグゼキュターの数です。
- numberMaxNeededExecutors は、現在の負荷を満たすために必要な (アクティブに実行中および保留中の) ジョブエグゼキュターの最大数です。

この2つのメトリクスは、自動スケーリングが正しく機能しているかどうかを判断するのに役立ちます。



AWS Glue は、数マイクロバッチにわたって `batchProcessingTimeInMs` メトリクスをモニタリングし、次の2つのいずれかを実行します。`batchProcessingTimeInMs` が `windowSize` に近い場合はエグゼキュターをスケールアウトし、`batchProcessingTimeInMs` が `windowSize` より比較的低い場合はエグゼキュターをスケールインします。また、エグゼキュターをステップスケールするアルゴリズムも使用します。

- インジケータ #1 は、負荷を処理するために、アクティブなエグゼキュターがどのようにスケールアップし、必要最大数のエグゼキュターに追いついたかを示しています。
- インジケータ #2 は、`batchProcessingTimeInMs` が低くなってから、アクティブなエグゼキュターがどのようにスケールインしたかを示しています。

これらのメトリクスを使用して、現在のエグゼキュターレベルの並列処理をモニタリングし、それに応じて自動スケーリング構成の最大ワーカー数を調整できます。

最高のパフォーマンスを得る方法

Spark は、Amazon Kinesis ストリーム内の読み取りのために、シャードごとに1つのタスクを作成しようとしています。各シャードのデータはパーティションになります。次に、各ワーカーのコア数に応じて、これらのタスクをエグゼキュター/ワーカーに分散します (ワーカーあたりのコア数は、選択したワーカータイプ (G.025X や G.1X など) によって異なります)。ただし、タスクがどのように分散されるかは非決定論的です。すべてのタスクは、それぞれのコアで並列実行されます。使用可能なエグゼキュターコアの数よりも多くのシャードがある場合、タスクはキューに入れられます。

上記のメトリクスとシャード数を組み合わせて、バーストのための余裕がある安定した負荷が実現するようにエグゼキュターをプロビジョニングできます。おおよそのワーカー数を決定するために、ジョブを数回繰り返すことをお勧めします。不安定で急激なワークロードの場合も、自動スケーリングと最大ワーカー数を設定することで同じことを実行できます。

`windowSize` は、ビジネスの SLA 要件に従って設定してください。例えば、処理されたデータが 120 秒を超えてはならないことをビジネスで義務付けている場合は、平均コンシューマーラグが 120 秒未満になるように `windowSize` を少なくとも 60 秒に設定します (前述のコンシューマーラグに関するセクションを参照)。そこから、`numRecords` とシャード数に応じて、DPU の容量を計画し、`batchProcessingTimeInMs` がほとんどいつも `windowSize` の 70% 未満になるようにします。

Note

ホットシャードはデータスキューの原因となる可能性があります。つまり、一部のシャード/パーティションが他のシャード/パーティションよりもはるかに大きくなる可能性があります。これにより、並列実行されている一部のタスクに時間がかかり、ストラグラータスクが発生する場合があります。その結果、前のバッチのすべてのタスクが完了するまで次のバッチを開始できず、これが `batchProcessingTimeInMillis` と最大ラグに影響することになります。

AWS Glue データ品質

AWS Glue Data Quality では、データの品質を測定およびモニタリングできるため、ビジネス上の適切な意思決定を行うことができます。オープンソース DeeQu フレームワーク上に構築された AWS Glue Data Quality は、マネージド型のサーバーレスエクスペリエンスを提供します。AWS Glue Data Quality は、データ品質ルールの定義に使用するドメイン固有の言語である Data Quality Definition Language (DQDL) と連携します。DQDL とサポートされているルールタイプの詳細については、「[データ品質定義言語 \(DQDL\) リファレンス](#)」を参照してください。

製品の詳細および価格に関するその他詳細は、[AWS Glue Data Quality](#) のサービスページを参照してください。

利点と主な特徴

AWS Glue Data Quality の利点と主な機能は次のとおりです。

- サーバーレス — インストール、パッチ適用、メンテナンスは不要です。
- すぐに開始 — AWS Glue Data Quality はデータをすばやく分析し、データ品質ルールを作成します。[Create Data Quality Rules] と [Recommend rules] をクリックするだけですぐに始められます。
- データ品質の問題の検出 — 機械学習 (ML) を使用して異常や hard-to-detect データ品質の問題を検出します。
- ルールを改善 — 開始する 25 以上の out-of-the-box DQ ルールを使用して、特定のニーズに合ったルールを作成できます。
- 品質を評価し、自信を持って判断を下せる — ルールを評価すると、データの状態に関する概要を示した Data Quality スコアを確認できます。この Data Quality スコアを使用すれば、自信を持ってビジネス上の判断を下せます。
- 不良データをゼロにする — AWS Glue Data Quality は、品質スコアの低下の原因となった正確なレコードを特定するのに役立ちます。それらをすぐに特定し、隔離して、修正できます。
- 従量制料金 — AWS Glue Data Quality を使用するために必要な年間ライセンスはありません。
- ロックインなし — AWS Glue Data Quality はオープンソース上に構築されているため DeeQu、オーサリングするルールをオープン言語で保持できます。
- データ品質チェック — AWS Glue データ品質 Data Catalog および ETL AWS Glue パイプラインにデータ品質チェックを適用することで、保管中および転送中のデータ品質を管理できます。

- ML ベースのデータ品質検出 — 機械学習 (ML) を使用して、異常や hard-to-detect データ品質の問題を検出します。

仕組み

AWS Glue Data Quality には、AWS Glue Data Catalog と AWS Glue ETL ジョブの 2 つのエントリポイントがあります。このセクションでは、各エントリポイントがサポートするユースケースと AWS Glue 機能の概要を説明します。

のデータ品質 AWS Glue Data Catalog

AWS Glue Data Quality AWS Glue Data Catalog は、に保存されているオブジェクトを評価します。非コーダーを使用すると、データ品質ルールを簡単にセットアップできます。例えば、データスチュワードやビジネスアナリストなどです。

この方法は、以下のようなユースケースに適しています。

- 既に AWS Glue Data Catalog でカタログ化しているデータセットに対してデータ品質タスクを実行したい場合
- データガバナンスに取り組んでおり、データレイク内のデータ品質に関する問題を、継続的に特定または評価する必要がある場合

データカタログのデータ品質は、次のインターフェイスで管理できます。

- AWS Glue マネジメントコンソール
- AWS Glue APIs

の AWS Glue Data Quality の使用を開始するには、AWS Glue Data Catalog 「」を参照してください [Data Catalog で AWS Glue Data Quality を使用する](#)。

AWS Glue ETL ジョブのデータ品質

AWS Glue AWS Glue ETL ジョブの Data Quality を使用すると、プロアクティブなデータ品質タスクを実行できます。プロアクティブなタスクは、データセットをデータレイクにロードする前に不良データを特定し、除外するのに役立ちます。

[動画: ETL パイプライン AWS Glue のデータ品質の紹介](#)

ETL ジョブ向けの Data Quality は以下のようなユースケースに使用できます。

- データ品質タスクを ETL ジョブに組み込む場合
- データ品質タスクを定義するコードを ETL スクリプトで記述する場合
- ビジュアルデータパイプラインに送信されるデータの品質を管理する場合

ETL ジョブのデータ品質は、次のインターフェイスで管理できます。

- AWS Glue Studio、AWS Glue Studio ノートブック、AWS Glue インタラクティブセッション
- AWS Glue ETL スクリプト用の ライブラリ
- AWS Glue APIs

ETL ジョブ向け Data Quality の使用方法については、「AWS Glue Studio ユーザーガイド」の「[チュートリアル: Data Quality の使用開始](#)」を参照してください。

データカタログ向け Data Quality と ETL ジョブ向け Data Quality の比較

この表は、AWS Glue Data Quality の各エントリポイントがサポートする機能の概要を示しています。

機能	データカタログ向け Data Quality	ETL ジョブ向け Data Quality
データソース	Amazon S3、Amazon Redshift、データカタログと互換性のある JDBC ソース、および Apache Iceberg、Apache Hudi、Delta Lake などトランザクションデータレイク形式。テーブルが AWS Lake Formation 管理されている場合、Iceberg、Delta、HUDI テーブルはサポートされません。でカタログ化されている Amazon Athena ビュー AWS	カスタムコネクタとサードパーティーコネクタを含む AWS Glue、でサポートされているすべてのデータソース。

機能	データカタログ向け Data Quality	ETL ジョブ向け Data Quality
	Glue Data Catalog はサポートされません。	
Data Quality ルールの推奨事項	サポート	サポートされていません
DQDL ルールの作成と実行	サポート	サポート
Auto scaling	サポート外	サポート
AWS Glue Flex サポート	サポート外	サポート
スケジューリング	Data Quality ルールを評価するときおよび Step Functions を使用する場合にサポート	Step Functions とワークフローを使用する場合にサポート
データ品質チェックに失敗したレコードの特定	サポート外	サポート
Amazon EventBridge との統合	サポート	サポート
AWS Cloudwatch との統合	サポート	サポート
データ品質評価の結果を Amazon S3 に書き込む	サポート	サポート
増分的なデータ品質	プッシュダウン述語によりサポート	AWS Glue ブックマークでサポート
AWS CloudFormation サポート	サポート	サポート
ML ベースの異常検出	サポートされていません	プレビュー
動的ルール	サポート外	サポート

考慮事項

AWS Glue Data Quality を使用する前に、次の項目を考慮してください。

- データ品質ルールは、ネストされたデータソースやリストタイプのデータソースを評価することはできません。[ネストされた Struct のフラット化](#) を参照してください。

用語

次のリストでは、AWS Glue Data Quality に関連する用語を定義します。

データ品質定義言語 (DQDL)

AWS Glue Data Quality ルールの記述に使用できるドメイン固有の言語。

DQDL の詳細については、[データ品質定義言語 \(DQDL\) リファレンス](#) のガイドを参照してください。

データ品質

データセットが特定の目的にどの程度役立つかについて説明します。AWS Glue Data Quality は、データセットに対してルールを評価し、データ品質を測定します。各ルールは、データの鮮度や整合性などの特性を確認します。データ品質を数値化するには、データ品質スコアを使用できます。

データ品質スコア

Data Quality でルールセットを評価するときに合格する (true になる) AWS Glue データ品質ルールの割合。

ルール

データにおける特性を確認し、ブール値を返す DQDL 式。詳細については、「[ルールの構造](#)」を参照してください。

analyzer

データ統計を収集する DQDL 式。アナライザーは、ML アルゴリズムが異常やデータ品質の問題を経時的に検出するために使用できる hard-to-detect データ統計を収集します。

ルールセット

一連のデータ品質ルールで構成される AWS Glue リソース。ルールセットは、AWS Glue Data Catalog内のテーブルと関連付ける必要があります。ルールセットを保存すると、AWS Glue は Amazon リソースネーム (ARN) をルールセットに割り当てます。

データ品質スコア

AWS Glue Data Quality でルールセットを評価する際に、合格する (結果が true になる) データ品質ルールの割合。

監視

ルールやアナライザーから収集されたデータ統計を経時的に分析することで、AWS Glue によって得られる未確認のインサイト。

制限

AWS Glue Data Quality サービスの制限：

- ルールセットには 2000 個のルールを含めることができます。ルールセットが大きい場合は、複数のルールセットに分割することをお勧めします。
- ルールセットのサイズは 65KB です。ルールセットが大きい場合は、複数のルールセットに分割することをお勧めします。

AWS Glue Data Quality のリリースノート

このトピックでは、AWS Glue Data Quality で導入された機能について説明します。

一般提供を開始: 新機能

AWS Glue Data Quality の一般提供では、以下の新機能が利用可能です。

- で、失敗したデータ品質チェックのレコードを識別する機能がサポートされるようになりました。
AWS Glue Studio
- 新しいデータ品質のルールタイプ。2 つのデータセット間におけるデータの参照整合性の検証、2 つのデータセット間におけるデータの比較、データタイプのチェックなど
- でのユーザーエクスペリエンスの向上 AWS Glue Data Catalog
- Apache Iceberg、Apache Hudi、Delta Lake のサポート

- Amazon Redshift のサポート
- Amazon による通知の簡素化 EventBridge
- AWS CloudFormation ルールセットの作成のサポート
- パフォーマンスの向上: ETL および のキャッシュオプション AWS Glue Studio により、データ品質を評価する際のパフォーマンスが向上します。

2023 年 11 月 27 日 (プレビュー)

- ML を利用した異常検出機能が AWS Glue ETL と AWS Glue Studio で利用できるようになりました。これにより、異常や hard-to-detect データ品質の問題を検出できるようになりました。
- 動的ルールでは、動的なしきい値 (例: `RowCount > avg(last(10))`) を設定できます。

2024 年 3 月 12 日

- DQDL の改良点
 - NULL、BLANKS、WHITESPACES_ONLY などのキーワードのサポート
 - AWS Glue Data Quality が複合ルールを処理する方法を指定するオプション
 - ColumnValues ルールタイプでは、比較中に NULL 値を渡すことはできません
 - DQDL での NOT 演算子のサポート

2024 年 6 月 26 日

- DQDL の改良点
 - DQDL は where 句をサポートするようになり、DQ ルールを適用する前にデータをフィルタリングできるようになりました。

AWS Glue Data Quality での異常検出

Note

AWS Glue Data Quality は、次のリージョンでプレビューで使用できます。

- 米国東部 (オハイオ、バージニア北部)

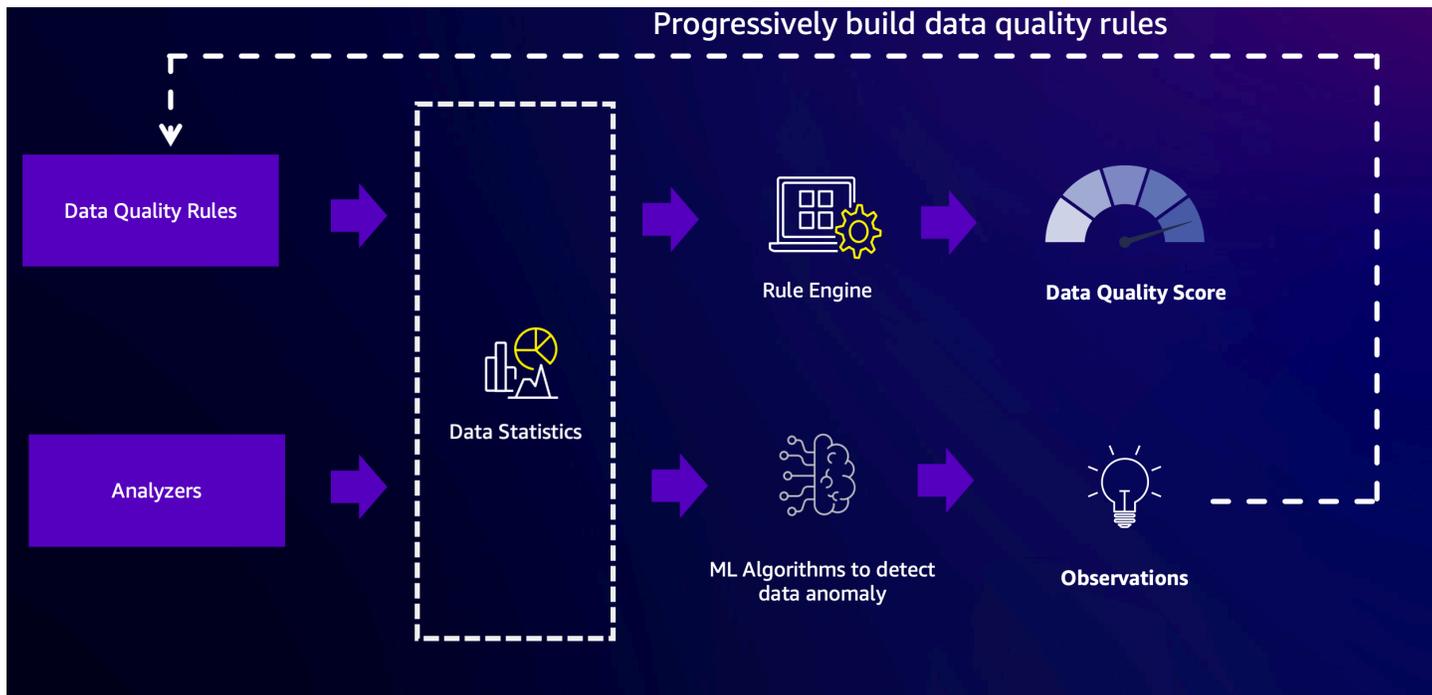
- 米国西部 (オレゴン)
- アジアパシフィック (東京)
- 欧州 (アイルランド)

AWS Glue Data Quality の異常検出では、経時的なデータ統計に機械学習 (ML) アルゴリズムを適用して、ルールでは検出が難しい異常パターンや隠れたデータ品質問題を検出します。現在のところ、異常検出は AWS Glue 4.0 でのみ利用できます。現在、この機能は AWS Glue Studio ビジュアル ETL と AWS Glue ETL でのみ使用できます。この機能は AWS Glue Studio ノートブック、AWS Glue Data Catalog、AWS Glue インタラクティブセッション、AWS Glue データプレビューでは機能しません。

仕組み

Data Quality のルールを評価する際、AWS Glue はデータがルールに準拠しているかどうかを判断するために必要なデータ統計情報を取得します。例えば、Data Quality はデータセット内の異なる値の数を計算し、その値を期待値と比較します。

Data Quality のルールエンジンは、統計値を定義済みのしきい値と比較し、品質要件を評価します。これらの統計は時間をかけて収集されるため、ETL パイプラインの異常検出を有効にすることで、AWS Glue が過去の統計から学習し、隠れたパターンを観察結果として報告できるようにすることができます。観察結果は、AWS Glue の ML アルゴリズムが特定する未確認のインサイトです。これらには Data Quality の推奨ルールが含まれており、それをルールセットに適用することで発見したパターンを監視できます。ジョブは (1 時間ごとや 1 日ごとなど) 定期的に実行することをおすすめします。不規則な実行では、質の悪いインサイトが提供される可能性があります。



アナライザーを使用してデータを検査する

データ品質ルールを作成する時間がない場合もあります。ここでアナライザーが役に立ちます。アナライザーはルールセットの一部であり、設定は非常に簡単です。例えば、ルールセットには次のように記述できます。

```
Analizers = [  
    RowCount,  
    Completeness "AllColumns"  
]
```

これにより以下の統計が収集されます。

- データセット全体の行数
- データセット内のすべての列の完全性

しきい値を気にする必要がないため、アナライザーの使用をお勧めします。データパイプラインを実行すると、3回実行した後から、AWS Glue Data Quality は異常を検出した時点で観察結果とルール推奨事項の生成を開始します。観察結果や関連する統計を確認でき、ルール推奨事項をルールセットに簡単に組み込むことができます。使用開始するには、「[異常検出の設定とインサイトの生成](#)」を

参照してください。アナライザーはデータ品質スコアには影響しないことに注意してください。アナライザーが生成する統計は経時的に分析され、観察結果を生成することができます。

DetectAnomaly ルールの使用

異常を検出したときに、ジョブを失敗させたい場合があります。制約を適用するには、ルールを設定する必要があります。アナライザーはジョブを停止しません。代わりに、統計を収集してデータを分析します。ルールセットのルールセクションで DetectAnomaly ルールを設定すると、ジョブがスキャン内のすべてのルールに合格しなかったことが DQ スキャンによって報告されることを確認できます。

異常検出の利点とユースケース

エンジニアは、いつでも数百のデータパイプラインを管理する可能性があります。各パイプラインは異なるソースからデータを抽出し、データレイクに読み込むことができます。各パイプラインは異なるソースからデータを抽出してデータレイクに読み込むことがあるため、データの形状が大幅に変化したのか、既存の傾向から逸脱したのかなど、データに関するフィードバックをすぐに得ることは困難です。

これまで、アップストリームデータソースはデータエンジニアリングチームに警告することなく変更され、このプロセスに hard-to-track 「データバグ」が導入されていました。Data Quality ノードをジョブに追加することで、問題が発見されるとジョブが失敗するため、作業が大幅に楽になります。ただし、これによってデータチームが懸念している障害モードがすべてなくなるわけではなく、他のデータバグが発生する可能性もあります。

障害モードの 1 つはデータ量に関するものです。企業のデータストアが時間の経過とともに増加するにつれて、データパイプラインによって生成されるレコードの数は指数関数的に増加する可能性があります。データチームは毎週 ETL ジョブを手動で更新して、取り込む行数に制限を設定するデータ品質ルールを増やす必要があるかもしれません。

もう 1 つの障害モードは、トランザクション量が曜日によって異なるという事実に対応するために、データ品質ルールの制限が極めて広範囲に及ぶことです。週末にはほとんどトランザクションがなく、月曜日には他の平日の約 3 倍のトランザクションがあります。データチームには 2 つの選択肢があります。1 つは、その日に合わせてルールセットをその場で変更するロジックを実装するか、もう 1 つは非常に広範な期待値を設定するかです。

最後に、データチームはあまり明確に定義されていないデータバグにも懸念を抱いています。モデルは特定の特性を持つデータに基づいてトレーニングされており、これらのデータが予期せぬ形で歪み

始めたら、チームはそれを知りたいと考えます。例えば、ある企業が 2 月にモンタナ州に進出する可能性があり、そうすると「MT」コードが含まれるトランザクションを開始する頻度が高くなります。これにより ML の推論が崩れる可能性があり、その結果、モデルではモンタナ州のすべての取引が不正であると誤って予測するかもしれません。

このような場合、Data Quality の異常検出がこうした問題の解決に役立ちます。Data Quality の異常検出の利点には次のようなものがあります。

- データスキャンが定期的、イベント主導型、手動ベースで可能。
- 意図しないイベント、季節性または統計上の異常を示す可能性のある異常の検出。
- Data Quality の異常検出が検出した観察結果に対する策を講じるためのルール推奨事項を提示。

これは次のような場合に役立ちます。

- データ品質ルールを記述することなく、データの異常を自動的に検出したい。
- データ品質ルールだけでは検出できない、データ内の潜在的な問題を特定したい。
- データ品質監視用に取り込まれる行数の制限など、時間の経過とともに進化するタスクを自動化したい。

AWS Glue Data Quality の IAM アクセス許可の設定

このトピックでは、IAM 管理者が AWS Glue Data Quality 向けの AWS Identity and Access Management (IAM) ポリシーで使用できるアクションとリソースの理解に役立つ情報を説明します。また、AWS Glue データカタログで AWS Glue Data Quality を使用するために必要な最小限の権限を含むサンプル IAM ポリシーも含まれています。

AWS Glue のセキュリティに関する追加情報については、「[AWS Glue のセキュリティ](#)」を参照してください。

AWS Glue Data Quality の IAM アクセス許可

次の表は、AWS Glue Data Quality で特定のオペレーションを実行する際にユーザーに必要なアクセス許可を一覧化したものです。AWS Glue Data Quality の権限を詳細に設定するときは、IAM ポリシーステートメントの Action 要素でこれらのアクションを指定します。

AWS Glue Data Quality アクション

[アクション]	説明	リソースタイプ
<code>glue:CreateDataQualityRuleset</code>	品質評価のルールセットを、1つ作成するためのアクセス許可を付与します。	<code>::dataQualityRuleset/<name></code>
<code>glue>DeleteDataQualityRuleset</code>	品質評価のルールセットを、1つ削除するためのアクセス許可を付与します。	<code>::dataQualityRuleset/<name></code>
<code>glue:GetDataQualityRuleset</code>	品質評価のルールセットを、1つ取得するためのアクセス許可を付与します。	<code>::dataQualityRuleset/<name></code>
<code>glue:ListDataQualityRulesets</code>	品質評価のルールセットを、すべて削除するためのアクセス許可を付与します。	<code>::dataQualityRuleset/*</code>
<code>glue:UpdateDataQualityRuleset</code>	品質評価のルールセットを、1つ更新するためのアクセス許可を付与します。	<code>::dataQualityRuleset/<name></code>
<code>glue:GetDataQualityResult</code>	品質評価タスクの実行結果を、1つ取得するためのアクセス許可を付与します。	<code>::dataQualityRuleset/<name></code>
<code>glue:ListDataQualityResults</code>	品質評価タスクの実行結果を、すべて取得するためのアクセス許可を付与します。	<code>::dataQualityRuleset/*</code>
<code>glue:CancelDataQualityRuleRecommendationRun</code>	実行中の品質評価推奨タスクを、1つ停止するためのアクセス許可を付与します。	<code>::dataQualityRuleset/*</code>
<code>glue:GetDataQualityRuleRecommendationRun</code>	実行済みの品質評価推奨タスクを、1つ取得するためのアクセス許可を付与します。	<code>::dataQualityRuleset/*</code>

[アクション]	説明	リソースタイプ
<code>glue:ListDataQualityRuleRecommendationRuns</code>	実行中の品質評価推奨タスクを、すべて取得するためのアクセス許可を付与します。	<code>::dataQualityRuleset/*</code>
<code>glue:StartDataQualityRuleRecommendationRun</code>	品質評価推奨タスクを、1つ実行開始するためのアクセス許可を付与します。	<code>::dataQualityRuleset/*</code>
<code>glue:CancelDataQualityRulesetEvaluationRun</code>	実行中の品質評価タスクを、1つ停止するためのアクセス許可を付与します。	<code>::dataQualityRuleset/*</code>
<code>glue:GetDataQualityRulesetEvaluationRun</code>	実行済みの品質評価タスクを、1つ取得するためのアクセス許可を付与します。	<code>::dataQualityRuleset/*</code>
<code>glue:ListDataQualityRulesetEvaluationRuns</code>	実行済みの品質評価タスクを、すべて取得するためのアクセス許可を付与します。	<code>::dataQualityRuleset/*</code>
<code>glue:StartDataQualityRulesetEvaluationRun</code>	品質評価タスクを、1つ実行開始するためのアクセス許可を付与します。	<code>::dataQualityRuleset/<name></code>
<code>glue:PublishDataQuality</code>	品質評価の結果を (複数) 公開するためのアクセス許可を付与します。	<code>::dataQualityRuleset/<name></code>

評価実行のスケジューリングに必要な IAM の設定

IAM アクセス許可

スケジュール済みの Data Quality 評価実行を実行するには、権限ポリシーに `IAM:PassRole` アクションを追加します。

AWS EventBridge スケジューラに必要なアクセス許可

[アクション]	説明	リソースタイプ
iam:PassRole	承認済みのロールをユーザーが渡すことを許可するアクセス許可を、IAM に付与します。	StartDataQualityRulesetEvaluationRun の呼び出しに使用されたロールの ARN

これらのアクセス許可がないと、次のエラーが発生します。

```
"errorCode": "AccessDenied"
"errorMessage": "User: arn:aws:sts::account_id:assumed-role/AWSGlueServiceRole is not
authorized to perform: iam:PassRole on resource: arn:aws:iam::account_id:role/service-
role/AWSGlueServiceRole
because no identity-based policy allows the iam:PassRole action"
```

IAM の信頼されたエンティティ

スケジュール済みの StartDataQualityEvaluationRun を作成して実行するには、AWS Glue と AWS EventBridge スケジューラサービスが、信頼されたエンティティに登録されている必要があります。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "glue.amazonaws.com"
      },
      "Action": "sts:AssumeRole"
    },
    {
      "Effect": "Allow",
      "Principal": {
        "Service": "scheduler.amazonaws.com"
      },
      "Action": "sts:AssumeRole"
    }
  ]
}
```

```
]
}
```

IAM ポリシーの例

AWS Glue Data Quality の IAM ロールには、次のタイプのアクセス許可が必要です。

- AWS Glue Data Quality オペレーションのアクセス許可。これにより、推奨されるデータ品質ルールを取得したり、AWS Glue データカタログ内のテーブルに対してデータ品質タスクを実行したりできます。このセクションの IAM ポリシーの例には、AWS Glue Data Quality オペレーションに必要な最低限のアクセス許可が含まれています。
- データカタログテーブルとその基盤となるデータへのアクセス許可。これらのアクセス許可は、ユースケースによって異なります。例えば、Amazon S3 でカタログ化するデータの場合、アクセス許可には Amazon S3 へのアクセスを含める必要があります。

Note

そのため、このセクションで説明するアクセス許可に加えて、Amazon S3 のアクセス許可を設定する必要があります。

推奨データ品質ルールの取得に最低限必要なアクセス許可

このサンプルポリシーには、推奨データ品質ルールを生成するために必要なアクセス許可が含まれています。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "AllowGlueRuleRecommendationRunActions",
      "Effect": "Allow",
      "Action": [
        "glue:GetDataQualityRuleRecommendationRun",
        "glue:PublishDataQuality",
        "glue:CreateDataQualityRuleset"
      ],
      "Resource": "arn:aws:glue:us-east-1:111122223333:dataQualityRuleset/*"
    },
  ],
}
```

```
"Sid": "AllowCatalogPermissions",
"Effect": "Allow",
"Action": [
  "glue:GetPartitions",
  "glue:GetTable"
],
"Resource": [
  "*"
]
},
{
  "Sid": "AllowS3GetObjectToRunRuleRecommendationTask",
  "Effect": "Allow",
  "Action": [
    "s3:GetObject"
  ],
  "Resource": "arn:aws:s3:::aws-glue-*"
},
{ // Optional for Logs
  "Sid": "AllowPublishingCloudwatchLogs",
  "Effect": "Allow",
  "Action": [
    "logs:CreateLogStream",
    "logs:CreateLogGroup",
    "logs:PutLogEvents"
  ],
  "Resource": "*"
},
]
}
```

品質評価タスクを実行するために最低限必要なアクセス許可

このサンプルポリシーには、データ品質評価タスクを実行するために必要なアクセス許可が含まれています。

以下のポリシーステートメントはオプションで、ユースケースに応じて設定します。

- AllowCloudWatchPutMetricDataToPublishTaskMetrics - Amazon CloudWatchにデータ品質の実行メトリクスを発行する場合は必須です。
- AllowS3PutObjectToWriteTaskResults - Amazon S3 にデータ品質の実行結果を書き込む場合は必須です。

```

{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "AllowGlueGetDataQualityRuleset",
      "Effect": "Allow",
      "Action": [
        "glue:GetDataQualityRuleset"
      ],
      "Resource": "arn:aws:glue:us-east-1:111122223333:dataQualityRuleset/<YOUR-  

RULESET-NAME>"
    },
    {
      "Sid": "AllowGlueRulesetEvaluationRunActions",
      "Effect": "Allow",
      "Action": [
        "glue:GetDataQualityRulesetEvaluationRun",
        "glue:PublishDataQuality"
      ],
      "Resource": "arn:aws:glue:us-east-1:111122223333:dataQualityRuleset/*"
    },
    {
      "Sid": "AllowCatalogPermissions",
      "Effect": "Allow",
      "Action": [
        "glue:GetPartitions",
        "glue:GetTable"
      ],
      "Resource": [
        "*"
      ]
    },
    {
      "Sid": "AllowS3GetObjectForRulesetEvaluationRun",
      "Effect": "Allow",
      "Action": [
        "s3:GetObject"
      ],
      "Resource": "arn:aws:s3:::aws-glue-*"
    },
    {
      "Sid": "AllowCloudWatchPutMetricDataToPublishTaskMetrics",
      "Effect": "Allow",

```

```
"Action": [
  "cloudwatch:PutMetricData"
],
"Resource": "*",
"Condition": {
  "StringEquals": {
    "cloudwatch:namespace": "Glue Data Quality"
  }
}
},
{
  "Sid": "AllowS3PutObjectToWriteTaskResults",
  "Effect": "Allow",
  "Action": [
    "s3:PutObject*"
  ],
  "Resource": "arn:aws:s3:::<YOUR-BUCKET-NAME>/*"
}
]
```

Data Catalog で AWS Glue Data Quality を使用する

この開始方法セクションでは、AWS Glue Data Quality コンソールから AWS Glue を使用する際に役立つ手順をご紹介します。品質評価ルールの推奨事項の生成や、データに照らしたルールセットの評価など、重要なタスクを完了する方法について確認できます。

トピック

- [前提条件](#)
- [手順の例](#)
- [ルールの推奨事項の生成](#)
- [ルールの推奨のモニタリング](#)
- [推奨ルールセットの編集](#)
- [新しいルールセットの作成](#)
- [ルールセットを実行してデータ品質を評価する](#)
- [データ品質スコアと結果を表示する](#)
- [関連トピック](#)

前提条件

AWS Glue Data Quality を使用する前に、Data Catalog と AWS Glue のクローラーの使用に慣れておく必要があります。AWS Glue Data Quality を使用すれば、Data Catalog データベースのテーブルの品質を評価できます。また、以下も必要になります:

- データ品質ルールセットの評価時に対照する Data Catalog のテーブル。
- ルールの推奨事項の生成やデータ品質タスクの実行時に、AWS Glue が使用する IAM ロール。このロールには、さまざまな AWS Glue Data Quality プロセスをユーザーに代わって実行するための、必要なリソースにアクセスする権限が必要です。この対象となるリソースには AWS Glue、Amazon S3、および CloudWatch が含まれます。AWS Glue Data Quality の最低限必要なアクセス許可を含むポリシーの例については、「[IAM ポリシーの例](#)」を参照してください。

AWS Glue 用の IAM ロールの詳細については、「[AWS Glue サービスの IAM ポリシーを作成する](#)」と「[AWS Glue 用の IAM ロールを作成する](#)」を参照してください。「[Authorization for AWS Glue Data Quality actions](#)」のページでは、Data Quality 専用の AWS Glue アクセス許可すべてのリストもご覧いただけます。

- さまざまなデータを含むテーブルが 1 つ以上あるデータベース。このチュートリアルで使用するテーブルには yyz-tickets という名前が付いており、テーブル tickets が付随しています。このデータは、トロント市が公開している駐車違反切符の情報を集めたものです。独自のテーブルを作成するときは、最適な推奨ルールを設定できるよう、さまざまな有効なデータがそこに入力されていることを確認します。

手順の例

サンプルのデータセットを使用した手順の例については、[AWS Glue Data Quality に関するブログ記事](#)を参照してください。

ルールの推奨事項の生成

ルールの推奨事項を利用すると、コードを記述することなく、すぐに Data Quality の使用を始められます。AWS Glue Data Quality を使用すると、データの分析、ルールの特定、品質評価タスクで評価されるルールセットの作成を行えます。推奨事項の実行は、実行から 90 日後に自動的に削除されます。

品質評価ルールの推奨事項を生成するには

1. AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。

2. ナビゲーションペインで、[Tables] (テーブル) を選択します。次に、品質評価ルールの推奨事項を生成する対象となるテーブルを選択します。
3. テーブルの詳細ページで、[Data quality] タブを選択し、そのテーブルの AWS Glue Data Quality ルールにアクセスします。
4. [Data quality] タブで、[Add rules and monitor data quality] を選択します。
5. [Ruleset builder] ページで、ページ上部のアラートに、「ルールの推奨事項が実行されていない場合は推奨事項タスクを開始する」との指示が表示されます。
6. [Recommend rules] を選択してモーダルを開き、推奨事項タスクのパラメータを入力します。
7. AWS Glue へのアクセス許可を持つ IAM ロールを選択します。このロールには、さまざまな AWS Glue Data Quality プロセスを、ユーザーに代わって実行するために必要なリソースに、アクセスするための許可が必要です。
8. 希望する設定に従ってフィールドに入力したら、[Recommend rules] を選択し、推奨事項タスクの実行を開始します。推奨事項の実行が進行中の場合、または完了した場合は、こちらのアラートで実行を管理できます。ステータスの変化を確認する際にページの更新が必要になる場合があります。完了済みと進行中の推奨事項タスクの実行状況は [実行履歴] のページに表示されます。ここには過去 90 日分の推奨事項の実行がすべて表示されます。

推奨ルールの意味

AWS Glue Data Quality は、入力テーブルの各列のデータに基づいてルールを生成します。このルールを使用して、品質要件を維持するために、データのフィルタリングが可能な潜在的境界を特定します。以下は生成されたルールのリストです。ルールの意味と、このルールがデータに適用されたときどのような処理が行われるのか、理解するのに役立つ例が記載されています。

生成されたデータ品質定義言語 (DQDL) のルールタイプの一覧は、「[データ品質定義言語 \(DQDL\)](#)」を参照してください。

- IsComplete "SET_FINE_AMOUNT" – 特定の行に対して列が入力されていることを確認する IsComplete ルールです。データ内で列を非オプションとしてタグ付けするときこのルールを使用します。
- Uniqueness "TICKET_NUMBER" > 0.95 – 列内のデータが一意的なしきい値を満たしていることを確認する Uniqueness ルールです。この例では、"TICKET_NUMBER" の特定の行に入力されたデータは他のすべての行と内容が最大 95% 同じであると判断されましたが、これはこのルールを示唆しています。

- ColumnValues "PROVINCE" in ["ON", "QC", "AB", "NY",...] – ColumnValues ルールは、既存の列の内容に基づいて、この列の有効な値を定義します。この例では、各行のデータは州または地域の 2 文字のナンバープレートコードです。
- ColumnLength "INFRACTION_DESCRIPTION" between 15 and 31 – ColumnLength ルールは、列のデータに文字数の制限を適用します。このルールは、文字列の列に記録された最小長と最大長に基づき、サンプルデータから生成されます。

ルールの推奨のモニタリング

品質評価ルールの推奨事項を実行すると、[Add rules and monitor data quality] ページに、情報と、上部のバーで実行できる追加のアクションが表示されます。

ルールの推奨が進行中の場合、推奨タスクが完了する前に [実行を停止] を選択できます。タスクが進行中の場合、ステータスは [in progress] と表示され、実行が開始された日付と時刻が表示されます。

ルールの推奨が完了すると、ルール推奨のバーに、推奨されたルールの数、最後の推奨実行のステータス、完了した日付とタイムスタンプが表示されます。

推奨ルールを追加するときは、[Insert Rule Recommendation] を選択します。過去に推奨されたルールを閲覧するには、具体的な日付を選択します。新しい推奨を実行するには、[More actions] を選択して、次に [Recommended rules] を選択します。

デフォルトの設定は、[Manage user settings] を選択して設定します。Amazon S3 のデフォルトパスを設定すると、ルールセットを保存したり、データカタログを実行するためのデフォルトロールを設定したりできます。

推奨ルールセットの編集

AWS Glue Data Quality は利用可能な既存のデータに基づいてルールを生成するため、自動化された提案に予期せぬルールや望ましくないルールが表示される場合があります。推奨ルールセットを最大限活用するには、ルールセットを評価し、修正を加えることが必要になります。本チュートリアルこちらのステップでは、前回のステップで生成したルールを修正し、一部のデータに、より制限の厳しい条件を適用します。また、他のルールの制限を緩和し、正確な一意のデータを後で追加できるようにすることも可能です。

推奨されたルールセットを編集する

1. AWS Glue コンソールで [データカタログ] を選択し、次にナビゲーションペインで [Databases tables] を選択します。テーブル tickets を選択します。
2. テーブルの詳細ページで、[Data quality] タブを選択し、そのテーブルの AWS Glue Data Quality のルールと設定にアクセスします。
3. [Rulesets] セクションで、[ルールの推奨事項の生成](#) で生成したルールセットを選択します。
4. [アクション] を選択し、次にコンソールウィンドウで [編集] を選択します。ルールセットエディタがコンソールにロードされます。ここでは、ルールの編集用ペインと DQDL のクイックリファレンスが含まれています。
5. スクリプトの 2 行目を削除します。これにより、データベースのサイズを特定の行数に制限する条件が緩和されます。編集後、ファイルの 1~3 行目に以下の内容が含まれているはずですが、

```
Rules = [  
  IsComplete "TAG_NUMBER_MASKED",  
  ColumnLength "TAG_NUMBER_MASKED" between 6 and 9,
```

6. スクリプトの 25 行目を削除します。これにより、記録された州の 96% が ON であるという条件は緩和されます。編集後、ファイルの 24 行目からルールセットの最後までに以下の内容が含まれているはずですが、

```
ColumnValues "PROVINCE" in ["ON", "QC", "AB", "NY", "AZ", "NS", "BC", "MI", "PQ",  
  "MB", "PA", "FL", "SK", "NJ", "OH", "NB", "IL", "MA", "CA",  
  "VA", "TX", "NF", "MD", "PE", "CT", "NC", "GA", "IN", "OR", "MN", "TN", "WI",  
  "KY", "MO", "WA", "NH", "SC", "CO", "OK", "VT", "RI", "ME", "AL",  
  "YT", "IA", "DE", "AR", "LA", "XX", "WV", "MT", "KS", "NT", "DC", "NV", "NE",  
  "UT", "MS", "NM", "ID", "SD", "ND", "AK", "NU", "GO", "WY", "HI"],  
ColumnLength "PROVINCE" = 2  
]
```

7. 14 行目を以下の内容に変更します。

```
IsComplete "TIME_OF_INFRACTION",
```

これにより、データベースを違反時間が記録されたチケットのみに制限して、列の条件を強化します。違反時間が記録されていないチケットは、このデータセットでは常に無効なデータとみなされます。これは、さらなるデータ使用や品質ルールを決定する際の調査など、パーティショニングや変換の方が適している場合とは異なります。

8. コンソールページの下部で [Update Ruleset] を選択します。

新しいルールセットの作成

ルールセットは、データに対応させて評価するための品質評価ルールをグループ化したものです。AWS Glue コンソールでは、データ品質定義言語 (DQDL) を使用してカスタムルールセットを作成できます。

品質評価ルールセットを作成するには

1. AWS Glue のコンソールで [データカタログ] を選択し、次に [データベース] を選択し、ナビゲーションペインで [テーブル] を選択します。テーブル tickets を選択します。
2. [Data quality] (データ品質) タブを開きます。
3. [Rulesets] セクションで、[Create ruleset] を選択します。DQDL エディタがコンソールで起動します。ここには、直接編集できるテキスト領域と、DQDL ルールおよびテーブルスキーマのクイックリファレンスがあります。
4. DQDL エディタのテキスト領域にルールを追加していきます。チュートリアルから直接ルールを記述することもできますし、データ品質ルールエディタの DQDL ルールビルダー機能を使用することもできます。

Note

DQDL ルールビルダーの使用方法

1. リストからルールタイプを選択し、プラス記号を選択してサンプルの構文をエディタペインに挿入します。
2. プレースホルダーの列名を、実際の列名で置き換えます。テーブルの列名は、[Create ruleset] タブで入手できます。
3. 必要に応じて式のパラメータを更新します。DQDL がサポートする式の完全なリストは、「[表現](#)」でご確認ください。

一例を挙げると、次のルールは tickets テーブルの ticket_number 列のデータ検証に関する制約です。次のルールを追加するには、DQDL ルールビルダーを使用するか、ルールセットを直接編集します。

```
IsComplete "ticket_number",  
IsUnique "ticket_number",  
ColumnValues "ticket_number" > 9000000000
```

5. [Run name] フィールドに新しいルールセットの名前を入力します。
6. [Save ruleset] を選択します。

複数のデータセットのデータ品質を評価する

ReferentialIntegrity ルールセットと DatasetMatch ルールセットを使用することで、複数のデータセットにデータ品質のルールを設定できます。ReferentialIntegrity は、プライマリデータセット内のデータが他のデータセット内に存在するかをチェックします。

参照のデータセットを追加するときは、[スキーマ] タブを選択し、次に [Update reference tables] を選択します。データベースとテーブルを選択するよう要求されます。テーブルを追加したら、データ品質ルールを設定できます。AggregateMatch、RowCountMatch、ReferentialIntegrity、SchemaMatch、DatasetMatch などのルールタイプは、複数のデータセットを横断してデータ品質チェックを実行する機能をサポートしています。

ルールセットを実行してデータ品質を評価する

品質評価タスクを実行すると、AWS Glue Data Quality はデータと対応させてルールセットを評価し、データ品質のスコアを計算します。このスコアは、入力に対して合格したデータ品質ルールの割合を表します。

品質評価タスクを実行するには

1. AWS Glue のコンソールで [データカタログ] を選択し、次に [データベース] を選択し、ナビゲーションペインで [テーブル] を選択します。テーブル tickets を選択します。
2. [Data quality] タブを選択します。
3. [Rulesets] リストから、テーブルに照らして評価を行うルールセットを選択します。このステップでは、生成されたルールではなく、自分で既に作成または変更したルールセットを使用することが推奨されます。[Run] (実行) を選択します。
4. モーダルで、自分の IAM ロールを選択します。このロールには、さまざまな AWS Glue Data Quality プロセスを、ユーザーに代わって実行するために必要なリソースに、アクセスするため

の許可が必要です。この IAM ロールをデフォルトとして保存するか、[Default Setting] のページに進んでこれを変更します。

5. [Data quality actions] (データ品質アクション) で、必要に応じて [Publish metrics to Amazon CloudWatch] (Amazon CloudWatch にメトリクスを公開する) を選択します。これを選択した場合、AWS Glue Data Quality は、評価に合格したルールと不合格だったルール、それぞれの数についてのメトリクスをパブリッシュします。この方法で保存されたメトリクスに対してアクションを実行するときは、CloudWatch アラームを使用します。アラートの設定用に、主要なメトリクスも Amazon EventBridge に公開されます。詳細については、「[Setting up alerts, deployments, and scheduling](#)」を参照してください。
6. [Run Frequency] で、[オンデマンドで実行する] または [Schedule the ruleset] を選択します。ルールセットをスケジュールした場合、タスク名を入力するよう求められます。スケジュールは Amazon EventBridge で作成されます。スケジュールの編集は Amazon EventBridge で行えます。
7. データ品質評価の結果を Amazon S3 に保存するには、[Data quality results location] を選択します。このタスクのために先に選択した IAM ロールには、このロケーションに対する書き込みのアクセス権限が必要です。
8. [Additional Configurations] の [Requested number of workers] に、AWS Glue でデータ品質タスクに割り当てる数を入力します。
9. オプションで、データソースにフィルターを設定できます。これにより、読み取るデータを減らせます。また、フィルターを使用して、パーティション情報を選択し、これを API 呼び出しを介してパラメータとして渡すことにより、増分検証を実行することもできます。パフォーマンスを改善するには、パーティション述語を指定します。
10. [Run] (実行) を選択します。[Data quality task runs] (実行中のデータ品質タスク) リストに、新しいタスクが表示されます。タスクの [実行ステータス] 列に [完了] と表示されたら、品質スコアの結果を確認できます。最新のステータスを確認するため、コンソール画面の更新が必要になる場合があります。
11. 品質評価結果の詳細の列を表示するには、[+] アイコンをクリックしてルールセットを展開します。結果には、評価に合格したルールと不合格だったルール、そして不合格の原因が表示されます。

データ品質スコアと結果を表示する

作成された全ルールセットの最新の実行を確認するには

1. AWS Glue コンソールのナビゲーションペインで、[Tables] (テーブル) を選択します。次に、品質評価タスクを実行する対象のテーブルを選択します。
2. [Data quality] タブを選択します。
3. [Data quality snapshot] には、時間の経過に伴う実行の一般的傾向が示されています。デフォルトでは、すべてのルールセットにおける最後の 10 回の実行が表示されます。ルールセットごとにフィルタリングするには、ドロップダウンリストから希望するルールセットを選択します。実行数が 10 回未満の場合は、完了した実行のうち表示可能なものがすべて表示されます。
4. [Data quality] テーブルには、各ルールセットが、最新の実行 (もしある場合) のスコアと共に表示されます。ルールセットを展開すると、そのルールセットに含まれるルールが、その実行のルール結果と併せて表示されます。

特定のルールセットにおける最新の実行を確認するには

1. AWS Glue コンソールのナビゲーションペインで、[Tables] (テーブル) を選択します。次に、品質評価タスクを実行する対象のテーブルを選択します。
2. [Data quality] タブを選択します。
3. [Data quality] テーブルで、特定のルールセットを選択します。
4. [Ruleset details] ページで、[Run history] タブを選択します。

この特定のルールセットに対するすべての評価実行が、このタブ内のテーブルに一覧表示されます。スコアの履歴と実行のステータスを確認できます。

5. 特定の実行に関する詳細情報を表示するには、[実行 ID] を選択して [Evaluation run details] のページに進みます。このページでは、実行の詳細と、個々のルール結果のステータスに関する詳細を確認できます。

関連トピック

- [DQDL ルールタイプリファレンス](#)
- [データ品質定義言語 \(DQDL\) リファレンス](#)

AWS Glue Studio を使用したデータ品質の評価

AWS Glue Data Quality は、定義したルールに基づき、データ品質の評価とモニタリングを行います。これにより、アクションが必要なデータを簡単に特定できます。AWS Glue Studio では、ビジュアルジョブにデータ品質ノードを追加して、データカタログ内のテーブルにデータ品質ルールを作成できます。これにより、経時的に進化するデータセットの変化をモニタリングし、評価できます。AWS Glue Studio で AWS Glue Data Quality を操作する方法の概要については、次の動画を参照してください。

AWS Glue Data Quality を使用する手順の概要を次に示します。

1. データ品質ルールの作成 – 設定した組み込みルールセットを選択して、DQDL ビルダーを使用してデータ品質ルールのセットを作成します。
2. データ品質ジョブの設定 – データ品質結果と出力オプションに基づいてアクションを定義します。
3. [Save and run a data quality job] – ジョブを作成して実行します。ジョブを保存すると、そのジョブ用に作成したルールセットが保存されます。
4. データ品質結果のモニタリングとレビュー – ジョブの実行が完了した後にデータ品質結果をレビューします。必要に応じて、ジョブを将来の日付にスケジュールすることもできます。

利点

データアナリスト、データエンジニア、データサイエンティストは、AWS Glue Studio でデータ品質評価ノードを使用し、ビジュアルジョブエディタでデータ品質を分析、設定、モニタリング、改善できます。データ品質ノードの使用には、次のような利点があります。

- データ品質の問題を検出可能 - データセットの特性をチェックするルールを作成することで、問題を確認できます。
- 簡単に開始可能 - 事前構築済みのルールとアクションで開始できます。
- 緊密な統合 - AWS Glue Data Quality は AWS Glue データカタログ上で実行されるため、AWS Glue Studio でデータ品質ノードを使用できます。

AWS Glue Studio で ETL ジョブのデータ品質を評価する

このチュートリアルでは、AWS Glue Studio で AWS Glue Data Quality の使用を開始します。次の方法について説明します。

- データ品質定義言語 (DQDL) ルールビルダーを使用して、ルールを作成します。
- データ品質アクション、出力するデータ、およびデータ品質結果を出力するロケーションを指定します。
- データ品質結果を確認する。

例を使用して練習するには、ブログ記事「[Getting started with AWS Data Quality for ETL pipelines](#)」を参照してください。

ステップ 1: データ品質評価変換ノードをビジュアルジョブに追加する

このステップでは、ビジュアルジョブエディタに、データ品質評価ノードを追加します。

データ品質ノードを追加するには

1. AWS Glue Studio コンソールで、[ジョブを作成する] セクションから [Visual with a source and target] を選択し、続いて [作成] を選択します。
2. データ品質変換を適用するノードを選択します。通常これは、変換ノードまたはデータソースです。
3. [+] アイコンを選択して、左側のリソースパネルを開きます。検索バーでデータ品質評価を検索し、取得した検索結果からデータ品質評価を選択します。
4. ビジュアルジョブエディタは、選択したノードから分岐したデータ品質評価変換ノードを表示します。コンソールの右側で、自動的に [Transform] (変換) タブが開きます。親ノードを変更する必要がある場合は、[Node properties] タブを選択し、ドロップダウンメニューで親ノードを選択します。

ノードの親を新たに選択すると、この親ノードとデータ品質評価ノード間に、新しい接続が確立されます。不要な親ノードをすべて削除します。1つのデータ品質評価ノードに接続できるのは、親ノード1つだけです。

5. データ品質評価変換は複数の親をサポートしているため、複数のデータセットにわたってデータ品質ルールを検証できます。複数のデータセットをサポートするルールには、ReferentialIntegrity、DatasetMatch、SchemaMatch、RowCountMatch、AggregateMatch などがあります。

データ品質評価変換に複数の入力を追加する場合は、「プライマリ」入力を選択する必要があります。プライマリ入力は、データ品質を検証するデータセットです。他のすべてのノードまたは入力は、参照として扱われます。

データ品質評価変換を使用して、データ品質チェックで不合格になった特定のレコードを識別できます。不良レコードにフラグを立てる新しい列がプライマリデータセットに追加されるため、プライマリデータセットを選択することをお勧めします。

6. 入力データソースのエイリアスを指定できます。ReferentialIntegrity ルールを使用する場合、エイリアスを使用して入力ソースを参照する別の方法があります。プライマリソースとして指定できるデータソースは 1 つだけのため、追加するデータソースにはそれぞれエイリアスが必要です。

次の例では、ReferentialIntegrity ルールによってエイリアス名で入力データソースを指定し、プライマリデータソースと 1 対 1 の比較を実行します。

```
Rules = [  
  ReferentialIntegrity "Aliasname.name" = 1  
]
```

ステップ 2: DQDL を使用してルールを作成する

このステップでは、DQDL を使用してルールを作成します。このチュートリアルでは、完全性ルールタイプを使用して 1 つのルールを作成します。このルールタイプは、特定の表現に対して列内の完全な (null 以外の) 値の割合を確認します。DQDL の使用の詳細については、「[DQDL](#)」を参照してください。

1. [変換] タブで、[Insert] ボタンを選択し、[Rule type] を追加します。これにより、ルールタイプがルールエディタに追加されるため、そこでルールのパラメータを入力できます。

Note

ルールを編集する際は、ルールがかっこ内にあり、互いにカンマで区切られていることを確認してください。以下は、完全なルール表現の例です。

```
Rules= [  
  Completeness "year">0.8, Completeness "month">0.8  
]
```

この例では、「年」と「月」の名前が付いた列に対して、完全性に関するパラメータを指定しています。このルールがパスするには、これらの列の完全性が 80% を超えているか、各列で 80% を超えるインスタンスにデータが含まれている必要があります。

この例では、[Completeness] (完全性) のルールタイプを検索して挿入します。これにより、ルールタイプがルールエディタに追加されます。このルールタイプの構文は、Completeness <COL_NAME> <EXPRESSION> の形式を取ります。

ほとんどのルールタイプにおいて、ブール型の応答を作成する場合は、パラメータとして表現を指定する必要があります。サポートされている DQDL 表現の詳細については、「[DQDL expressions](#)」を参照してください。次に、列名を追加します。

2. DQDL ルールビルダーで、[スキーマ] タブを選択します。検索バーを使用して、入力スキーマの中の列名を特定します。入力スキーマには、列名とデータ型が表示されます。
3. ルールエディターで、ルールタイプの右側をクリックして、列を挿入する場所にカーソルを移動します。または、ルールの中に列名を入力することもできます。

例えば、入力スキーマリストの列リストで、列の横にある [Insert] ボタンを選択します (この例では [year])。これにより、列がルールに追加されます。

4. 次に、ルールエディターで、ルールを評価するための表現を追加します。完全性ルールタイプは特定の表現に対して列内の完全な (null 以外の) 値の割合を確認するため、> 0.8 のような値を入力します。このルールは、列に完全な値 (null 以外) が 80% を超える割合で含まれているかどうかを確認します。

ステップ 3: データ品質出力の設定

データ品質ルールを作成したら、データ品質ノード出力を指定する追加のオプションを選択できます。

1. [Data quality transform output] (データ品質変換の出力) で、次のオプションから選択します。
 - [Original data] – 元の入力データを出力することを選択します。このオプションを選択すると、新しい子ノード「rowLevelOutcomes」がジョブに追加されます。スキーマは、変換の入力として渡されたプライマリデータセットのスキーマと一致します。このオプションは、品質の問題が発生したときに、データを渡してジョブを失敗させる場合に便利です。

別のユースケースとしては、データ品質チェックに失敗した不良レコードを検出する場合があります。不良レコードを検出するには、オプション [Add new columns to indicate data quality errors] を選択します。このアクションにより、「rowLevelOutcomes」変換のスキーマに 4 つの新しい列が追加されます。

- DataQualityRulesPass (文字列配列) — データ品質チェックに合格したルールの配列を提供します。
 - DataQualityRulesFail (文字列配列) — データ品質チェックに失敗したルールの配列を提供します。
 - DataQualityRulesSkip (文字列配列) — スキップされたルールの配列を提供します。次のルールはデータセットレベルで適用されるため、エラーレコードを識別できません。
 - AggregateMatch
 - ColumnCount
 - ColumnExists
 - ColumnNamesMatchPattern
 - CustomSql
 - RowCount
 - RowCountMatch
 - StandardDeviation
 - 平均値
 - ColumnCorrelation
 - DataQualityEvaluationResult — 行レベルで「合格」または「不合格」のステータスを表示します。全体的な結果は不合格でも、特定の記録は合格である可能性があることに注意してください。例えば、RowCount ルールは失敗でも、その他のルールはすべて成功である可能性があります。このような場合、このフィールドのステータスは「合格」です。
2. [Data quality results] – 設定したルールと、それに関するステータス (合格/不合格) を出力することを選択します。このオプションは、結果を Amazon S3 またはその他のデータベースに書き込む場合に便利です。
 3. [Data quality output settings] (オプション) – [Data quality output settings] を選択し、[Data quality result location] フィールドを表示します。続いて [参照] を選択し、データ品質の出力先として設定する Amazon S3 の場所を検索します。

Step 4. データ品質アクションの設定

このアクションを使用すると、特定の基準に基づいてメトリクスを CloudWatch に公開したり、ジョブを停止したりできます。アクションは、ルールを作成した後にのみ使用できます。このオプションを選択すると、同じメトリクスが Amazon EventBridge にも公開されます。このオプションを使用すると、[通知用アラートを作成できます](#)。

- [On ruleset failure] — ジョブ実行中にルールセットが失敗した場合の対処方法を選択できます。データ品質が低下した場合にジョブを失敗させる場合は、次のいずれかのオプションを選択して、ジョブを失敗させるタイミングを選択します。デフォルトではこのアクションは選択されておらず、データ品質ルールが失敗した場合でもそのジョブの実行は完了します。
- [なし] — [なし] (デフォルト) を選択した場合、ジョブは失敗せず、ルールセットが失敗しても実行が継続されます。
- [ターゲットにデータをロードした後にジョブを失敗させる] — ジョブは失敗し、データは保存されません。結果を保存するには、データ品質の結果を保存する Amazon S3 の場所を選択します。
- [Fail job without loading to target data] — このオプションは、データ品質エラーが発生するとすぐにジョブを失敗させます。データ品質変換の結果を含め、データターゲットはロードされません。

ステップ 5: データ品質結果を表示する

ジョブの実行が完了したら、[Data quality] タブを選択して、データ品質結果を表示します。

1. データ品質結果はジョブを実行するたびに表示されます。各ノードには、データ品質ステータスと、そのステータスに関する詳細が表示されます。ノードを選択すると、すべてのルールと各ルールのステータスが表示されます。
2. [Download results] を選択して、ジョブの実行とデータ品質結果に関する情報が記載された CSV ファイルをダウンロードします。
3. qデータ品質結果と実行されたジョブが複数存在する場合は、日付と時間の範囲を指定することで、結果をフィルタリングできます。[Filter by a date and time range] を選択すると、フィルターウィンドウが展開します。
4. 相対範囲 (relative range) と絶対範囲 (absolute range) のどちらかを選択します。絶対範囲を選択した場合は、カレンダーで日付を選択した後に、開始時刻と終了時刻を入力します。完了したら、[適用] を選択します。

Data Quality ルールビルダー

データ品質定義言語 (DQDL) ルールビルダーを使用すると、データを評価するデータ品質ルールを作成できます。最初にルールタイプを選択し、次にルールエディタでパラメータを指定します。ルールエディタには、ルールを作成する際のエラーや警告も表示されます。

[DQDL のガイド](#)には、DQDL 構文、組み込みのルールタイプ、および例を使用してルールを作成する方法に関する、包括的なドキュメントが用意されています。

データ品質評価ノード

データ品質評価変換ノードと DQDL ルールビルダーを使用する際に、作業スペースを拡張できます。

- [変換] タブを画面全体に展開するには、ノードの詳細パネルの右上隅にある展開アイコンを選択します。
- DQDL ルールエディタを展開するには、[<<] アイコンを選択してルールエディタを展開し、[Rule types] タブと [スキーマ] タブを折りたたみます。

コンポーネント

AWS Glue Studio には 26 のルールタイプが組み込まれています。各ルールタイプには、その使用方法の説明と例があります。

データ品質ルールタイプ

AWS Glue Studio には、ルールを簡単に作成するためのルールタイプが組み込まれています。ルールタイプの詳細については、「[DQDL ルールタイプリファレンス](#)」を参照してください。

Schema

[Schema] (スキーマ) タブには、親ノードの列名とデータ型が表示されます。複数のノードのスキーマが表示されます。入カスキーマを表示したり、列名で検索したり、ルールエディタに列を挿入したりできます。

Node properties | **Transform** | **Output schema** | **Data preview**

Evaluate data quality [Info](#)
Evaluate data quality by defining your data quality rules and actions

Data quality rules [Info](#)
Add rules using DQDL (Data Quality Definition Language)

DQDL rule builder <<

Rule types (18) **Schema**

Search

▼ **Input schema**

- year
int
- month
int
- day
int
- fl_date
string

```
1 Rules= [  
2   Completeness "year" > 0.8  
3 ]
```

Ln 1, Col 1 Errors: 0 Warnings: 0

ルールエディタ

ルールエディタは、ルールを作成および編集できるテキストエディタです。DQDL ルールビルダーからルールタイプを選択すると、そのルールタイプがルールエディタに追加されます。その後、テキストを変更することで、必要に応じてパラメータを指定したり、ルールを追加したり、ルールを編集したりできます。AWS Glue Studio はルールエディタ内のルールを検証し、エラーと警告がある場合はそれらを表示します。

エラーおよび警告

ルールが DQDL ルール構文に従っていない場合、エラーがあることを視覚的に示すいくつかのインジケータがルールエディタに表示されます。

- ルールエディタにエラーアイコンが表示され、エラーのある行が赤くなります。
- ルールエディタでは、赤いエラーアイコンの横にエラーの数が表示されます。
- エラーのある行を選択すると、エラーの説明と場所 (行と列) がルールエディタの下部に表示されます。

Node properties | **Transform** | **Output schema** | **Data preview**

Evaluate data quality [Info](#)
Evaluate data quality by defining your data quality rules and actions

Data quality rules [Info](#)
Add rules using DQDL (Data Quality Definition Language)

DQDL rule builder

Rule types (18) | **Schema**

ColumnCorrelation
column rule
▶ **Description, examples**

ColumnExists
column rule
▶ **Description, examples**

ColumnLength
column rule
▶ **Description, examples**

ColumnCorrelation
Ln 1, Col 18

Ln 1, Col 1 h is null

データ品質アクション

デフォルトでは、このアクションは選択されておらず、データ品質ルールが失敗した場合でもそのジョブの実行は完了します。

次のアクションの中から選択します。アクションを使用して、特定の基準に基づいて結果を CloudWatch に公開したり、ジョブを停止したりできます。アクションは、ルールを作成した後のみ使用できます。

- [Publish results to CloudWatch] – ジョブを実行し、その結果を CloudWatch に追加します。
- [Fail job when data quality fails] – データ品質ルールが失敗すると、結果的にジョブも失敗します。

データ品質変換出力

- [Original data] – 元の入力データを出力することを選択します。このオプションは、品質の問題が検出されたときにジョブを停止する場合に理想的です。
- [Data quality metrics] – 設定したルールとその合格または不合格のステータスを出力するように選択します。このオプションは、カスタムアクションを実行する場合に便利です。

データ品質出力設定

Amazon S3 の場所をデータ品質出力先として指定し、データ品質結果の場所を設定します。

異常検出の設定とインサイトの生成

AWS Glue Data Quality (DQ) は、ユーザーが作成したデータ品質ルールに基づいてデータを評価し、長期にわたるデータに関するインサイトや観察結果を提供するため、ユーザーはすぐに行動を起こすことができます。Data Quality はデータをスキャンするため、行数、最大値、最小値などの統計メトリクスを計算し、それらをしきい値式と比較します。

Data Quality による異常検出には次のようなメリットがあります。

- データの継続的な自動スキャン
- 意図しないイベントや統計上の異常を示唆する異常の検出
- Data Quality の異常検知で明らかになった観察結果に対して対策を講じるためのルールの推奨事項を提示

これは次のような場合に役立ちます。

- データ品質を記述することなく、データの異常を自動的に検出したい
- データをプロファイリングして、データがどのように見えるかを視覚的に表現したい
- データの経時的変化を追跡したい

自分のデータについて、観察結果はどのように表示できますか？

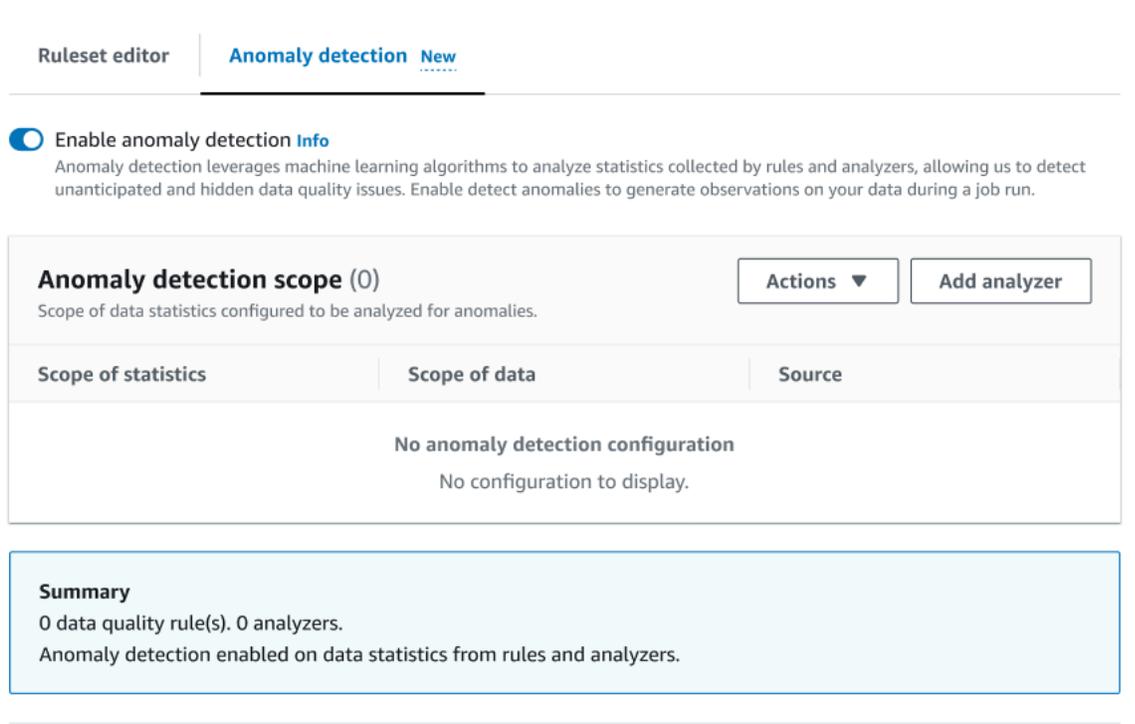
Data Quality は、収集されたデータ統計における外れ値、データ形式の変更、データドリフト、スキーマの変更を特定します。Data Quality は観察結果に基づいて、ユーザーが簡単に運用できるデータ品質ルールを推奨します。統計には、完全性、一意性、平均、合計 StandardDeviation、エントロピー DistinctValuesCount、およびが含まれます UniqueValueRatio。

AWS Glue Studio での異常検出の有効化

異常検出を有効にするには、AWS Glue Studio ジョブを開いて [異常検出を有効にする] をオンにします。これをオンにするとデータの異常検出が可能になり、データを経時的に分析し、データや観察結果に関するデータ統計が提供され、それに基づいて行動できるようになります。

AWS Glue Studio で異常検出を有効にするには:

1. ジョブの [Data Quality] ノードを選択し、[異常検出] タブをクリックします。[異常検出を有効にする] をオンに切り替えます。



Ruleset editor | **Anomaly detection** New

Enable anomaly detection [Info](#)
Anomaly detection leverages machine learning algorithms to analyze statistics collected by rules and analyzers, allowing us to detect unanticipated and hidden data quality issues. Enable detect anomalies to generate observations on your data during a job run.

Anomaly detection scope (0) Actions ▾ Add analyzer
Scope of data statistics configured to be analyzed for anomalies.

Scope of statistics	Scope of data	Source
No anomaly detection configuration No configuration to display.		

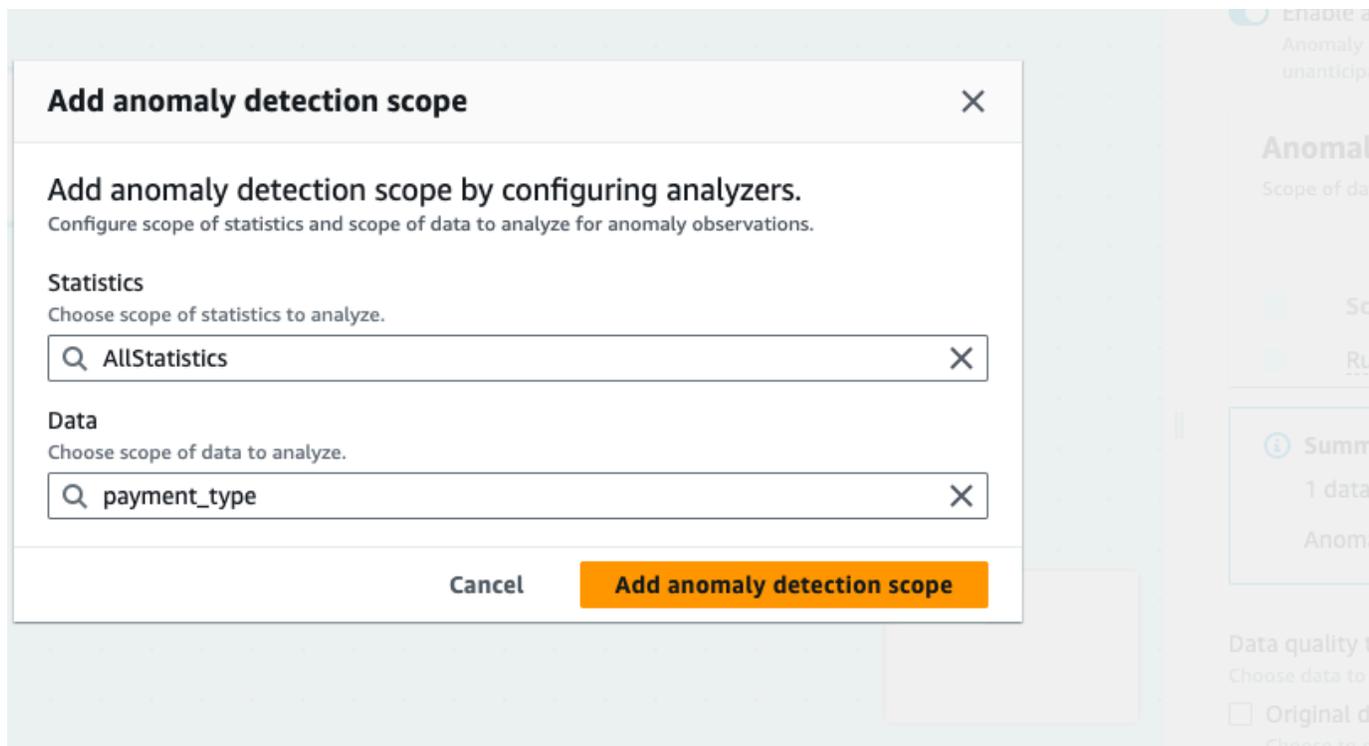
Summary
0 data quality rule(s). 0 analyzers.
Anomaly detection enabled on data statistics from rules and analyzers.

2. [アナライザーを追加] を選択して、異常を監視するデータを定義します。入力できるフィールドには、[統計] と [データ] の 2 つがあります。

統計は、データの形状やその他の特性に関する情報です。一度に 1 つまたは複数の統計情報を選択することも、すべての統計情報を選択することもできます。統計には、完全性、一

意性、平均、合計、エンтроピー StandardDeviation、 DistinctValuesCount、が含まれます UniqueValueRatio。

データはデータセット内の列です。すべての列を選択することも、個々の列を選択することもできます。



3. [異常検出範囲を追加] を選択して変更を保存します。アナライザーを作成したら、[異常検出範囲] セクションで確認できます。

[アクション] メニューを使用してアナライザーを編集したり、[ルールセットエディター] タブを選択してルールセットエディターのメモ帳でアナライザーを直接編集したりすることもできます。保存したアナライザーは、作成したルールのすぐ下に表示されます。

```
Rules = [  
]  
  
Analyzers = [  
  Completeness "id"  
]
```

更新されたルールセットとアナライザーにより、Data Quality は受信データを継続的に監視し、設定に基づいてアラートやジョブ停止を通じて異常を通知します。

Note

観察結果は、データセット内のデータ統計ごとに3つ以上の値が観測された場合に生成されます。観察結果が表示されない場合、[データ品質]には観察結果を生成するのに十分なデータがありません。ジョブを数回実行することにより、[データ品質]でデータに関するインサイトを提供することができ、そのインサイトは[観察結果]セクションに表示されます。

アナライザーはデータ内の異常を検出して観察結果を生成し、ルールを段階的に構築するための推奨事項を提示します。観察結果は[データ品質]タブを選択すると表示できます。観察結果はジョブの実行ごとに異なります。[観察結果]セクションの上部には、特定の Data Quality ノードとジョブ実行が表示されます。新しいノードまたはジョブ実行を選択すると、そのノードとジョブに固有の観察結果が表示されます。

The screenshot displays the AWS Glue Data Quality interface. At the top, there are navigation tabs: Visual, Script, Job details, Runs, Data quality - updated, Schedules, and Version Control. Below this, a section titled 'Getting started with data quality (DQ)' provides introductory text and a 'Give feedback' button. The main area shows 'Data quality at' with a dropdown menu and a 'Go to node' button. Below that, there are tabs for 'Rules (0)' and 'Observations (3) - new'. The 'Observations (3) - new' section is active, showing a list of observations. The first observation is 'RowCount of 19999.0 is lower than the detected lower bound of 61509.0'. It includes related metrics for 'Dataset: RowCount' such as ActualValue (19999.00), ExpectedValue (72964.00), LowerLimit (61509.00), and UpperLimit (84347.00). The rule recommendation is 'RowCount between 61508.00 and 84348.00'. The observed data is 'RowCount'. The second observation is 'Completeness for column fare_amount of 0.96 is lower than the detected lower bound of 1.0'. It includes related metrics for 'Column: fare_amount.Completeness' such as ActualValue (0.96), ExpectedValue (1.00), and LowerLimit (1.00). The rule recommendation is 'Completeness "fare_amount" = 1.0'. The observed data is 'ColumnName: fare_amount'. At the bottom, there is a line chart titled 'Dataset: RowCount' showing the row count over time from Nov 21 10:10 to Nov 21 10:18. The y-axis ranges from 0.00 to 80K. The line shows a peak around 70K at 10:15 and a sharp decline to around 20K by 10:18.

観察結果 — 各インサイトは、指定したルールセットとアナライザーによって設定された特定のジョブ実行に基づいています。

関連メトリクス — 観察結果が生成されると、[関連メトリクス]列には、ルール、実際の値、期待値、下限と上限が表示されます。

ルールの推奨事項 — AWS Glue はその後、これに対処するためのルールを推奨します。推奨されている各ルールは、コピーアイコンをクリックしてコピーできます。推奨ルールをすべてコピーする

には、各ルールの横にあるコピーアイコンをクリックし、[コピーしたルールを適用] をクリックします。

監視対象データ — [監視対象データ] 列には、監視対象となり観察結果をもたらす要因となった列または行が表示されます。

データ品質ノードに推奨ルールを適用する

観察結果が生成され、推奨ルールが提供されたら、そのルールをデータ品質ノードに適用できます。これを実行するには:

1. 各推奨ルールの横にあるコピーアイコンをクリックします。これでルール推奨事項がメモ帳に追加され、後で読み込むことができます。
2. [推奨ルールを適用] をクリックします。メモ帳が開き、以前にコピーしたルールを確認できます。
3. [ルールをコピー] を選択します。
4. [ルールセットエディターに適用] を選択します。ルールセットエディターが開き、コピーしたルールを貼り付けることができます。
5. コピーしたルールをルールセットエディターに貼り付けます。

AWS Glue Studio ノートブックの ETL ジョブのデータ品質

このチュートリアルでは、AWS Glue Studio ノートブックで、抽出、変換、ロード (ETL) ジョブに AWS Glue Data Quality を使用方法について説明します。

AWS Glue Studio でノートブックを使用してジョブスクリプトを編集し、完全なジョブを実行せずに出力を表示できます。マークダウンを追加したり、ノートブックを .ipynb ファイルやジョブスクリプトとして保存したりすることもできます。ソフトウェアをローカルにインストールしたり、サーバーを管理したりすることなくノートブックを起動できます。コードに問題がなければ、AWS Glue Studio を使用して、ノートブックを AWS Glue ジョブに簡単に変換できます。

この例で使用されているデータセットは、2 つの Data.CMS.gov データセット (Inpatient Prospective Payment System Provider Summary for the Top 100 Diagnosis-Related Groups - FY2011 および Inpatient Charge Data FY 2011) からダウンロードされた、メディケアプロバイダーの支払いデータで構成されています。

ダウンロードした後、データセットを修正してファイルの最後にいくつかの誤ったレコードを追加しました。この変更されたファイルは、パブリックな Amazon S3 バケット (s3://awsglue-datasets/examples/medicare/Medicare_Hospital_Provider.csv) 内に置かれています。

前提条件

- 宛先 Amazon S3 バケットに書き込むための Amazon S3 アクセス許可を持つ AWS Glue ロール
- 新しいノートブック (「[AWS Glue Studio 中でのノートブックの使用開始](#)」を参照)

AWS Glue Studio での ETL ジョブの作成

ETL ジョブを作成するには

1. セッションバージョンを AWS Glue 3.0 に変更します。

変更するには、以下の方法で定型コードのセルをすべて削除し、セルを実行します。この定型コードは、新しいノートブックを作成すると、最初のセルに自動的に入力されることに注意してください。

```
%glue_version 3.0
```

2. 次のコードをコピーして、セルで実行します。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
```

3. 次のセルで、AWS Glue Data Quality を評価する EvaluateDataQuality クラスをインポートします。

```
from awsgluedq.transforms import EvaluateDataQuality
```

- 次のセルで、公開 Amazon S3 バケットに保存されている .csv ファイルを使用して、ソースデータを読み取ります。

```
medicare = spark.read.format(
    "csv").option(
    "header", "true").option(
    "inferSchema", "true").load(
    's3://aws glue-datasets/examples/medicare/Medicare_Hospital_Provider.csv')
medicare.printSchema()
```

- データを変換しますAWS Glue DynamicFrame。

```
from awsglue.dynamicframe import DynamicFrame
medicare_dyf = DynamicFrame.fromDF(medicare,glueContext,"medicare_dyf")
```

- データ品質定義言語 (DQDL) を使用して、ルールセットを作成します。

```
EvaluateDataQuality_ruleset = """
    Rules = [
        ColumnExists "Provider Id",
        IsComplete "Provider Id",
        ColumnValues " Total Discharges " > 15
    ]
    """
```

- データセットをルールセットに照らして検証します。

```
EvaluateDataQualityMultiframe = EvaluateDataQuality().process_rows(
    frame=medicare_dyf,
    ruleset=EvaluateDataQuality_ruleset,
    publishing_options={
        "dataQualityEvaluationContext": "EvaluateDataQualityMultiframe",
        "enableDataQualityCloudWatchMetrics": False,
        "enableDataQualityResultsPublishing": False,
    },
    additional_options={"performanceTuning.caching": "CACHE_NOTHING"},
)
```

8. 結果を確認します。

```
ruleOutcomes = SelectFromCollection.apply(
  dfc=EvaluateDataQualityMultiframe,
  key="ruleOutcomes",
  transformation_ctx="ruleOutcomes",
)

ruleOutcomes.toDF().show(truncate=False)
```

出力:

```
-----+-----
+-----+-----
+-----+-----
|Rule                                     |Outcome|FailureReason
          |EvaluatedMetrics                               |
+-----+-----
+-----+-----
|ColumnExists "Provider Id"             |Passed |null
          |{}                                             |
|IsComplete "Provider Id"               |Passed |null
          |{Column.Provider Id.Completeness -> 1.0}    |
|ColumnValues " Total Discharges " > 15|Failed |Value: 11.0 does not meet the
  constraint requirement!|{Column. Total Discharges .Minimum -> 11.0}|
+-----+-----
+-----+-----
+-----+-----
```

9. 合格した行をフィルタリングし、データ品質の行レベルの結果から不合格の行を確認します。

```
rowLevelOutcomes = SelectFromCollection.apply(
  dfc=EvaluateDataQualityMultiframe,
  key="rowLevelOutcomes",
  transformation_ctx="rowLevelOutcomes",
)
```

```

rowLevelOutcomes_df = rowLevelOutcomes.toDF() # Convert Glue DynamicFrame to
SparkSQL DataFrame
rowLevelOutcomes_df_passed =
rowLevelOutcomes_df.filter(rowLevelOutcomes_df.DataQualityEvaluationResult ==
"Passed") # Filter only the Passed records.
rowLevelOutcomes_df.filter(rowLevelOutcomes_df.DataQualityEvaluationResult ==
"Failed").show(5, truncate=False) # Review the Failed records

```

出力:

```

+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
|DRG Definition          |Provider Id|Provider Name
      |Provider Street Address |Provider City|Provider State|Provider Zip
Code|Hospital Referral Region Description| Total Discharges | Average Covered
Charges | Average Total Payments |Average Medicare Payments|DataQualityRulesPass
      |DataQualityRulesFail      |DataQualityRulesSkip      |
DataQualityEvaluationResult|
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+
|039 - EXTRACRANIAL PROCEDURES W/O CC/MCC|10005      |MARSHALL MEDICAL CENTER SOUTH
      |2505 U S HIGHWAY 431 NORTH|BOAZ          |AL          |35957
|AL - Birmingham          |14          |$15131.85
|$5787.57                  |$4976.71    |[[IsComplete "Provider Id"]|
[ColumnValues " Total Discharges " > 15]]|ColumnExists "Provider Id"]|Failed
      |
|039 - EXTRACRANIAL PROCEDURES W/O CC/MCC|10046      |RIVERVIEW REGIONAL MEDICAL
CENTER |600 SOUTH THIRD STREET |GADSDEN      |AL          |35901
|AL - Birmingham          |14          |$67327.92
|$5461.57                  |$4493.57    |[[IsComplete "Provider Id"]|
[ColumnValues " Total Discharges " > 15]]|ColumnExists "Provider Id"]|Failed
      |

```

```
|039 - EXTRACRANIAL PROCEDURES W/O CC/MCC|10083      |SOUTH BALDWIN REGIONAL
MEDICAL CENTER|1613 NORTH MCKENZIE STREET|FOLEY      |AL          |36535
      |AL - Mobile          |15          |$25411.33
      |$5282.93            |$4383.73    |[IsComplete "Provider
Id"]|[ColumnValues " Total Discharges " > 15]|[ColumnExists "Provider Id"]|Failed
      |
|039 - EXTRACRANIAL PROCEDURES W/O CC/MCC|30002      |BANNER GOOD SAMARITAN MEDICAL
CENTER |1111 EAST MCDOWELL ROAD |PHOENIX      |AZ          |85006
      |AZ - Phoenix          |11          |$34803.81
      |$7768.90            |$6951.45    |[IsComplete "Provider Id"]|
[ColumnValues " Total Discharges " > 15]|[ColumnExists "Provider Id"]|Failed
      |
|039 - EXTRACRANIAL PROCEDURES W/O CC/MCC|30010      |CARONDELET ST MARYS HOSPITAL
      |1601 WEST ST MARY'S ROAD |TUCSON       |AZ          |85745
      |AZ - Tucson          |12          |$35968.50
      |$6506.50            |$5379.83    |[IsComplete "Provider Id"]|
[ColumnValues " Total Discharges " > 15]|[ColumnExists "Provider Id"]|Failed
      |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 5 rows
```

AWS Glue Data Quality は 4 つの新しい列 (DataQualityRulesPass、 DataQualityRulesFail、 DataQualityRulesSkip、 および) を追加したことに注意してください
DataQualityEvaluationResult。これには、合格したレコード、不合格のレコード、行レベルの評価でスキップされたルール、および行レベルの全体的な結果が表示されます。

10. Amazon S3 バケットに出力を書き込んでデータを分析し、結果を視覚化します。

```
#Write the Passed records to the destination.

glueContext.write_dynamic_frame.from_options(
    frame = rowLevelOutcomes_df_passed,
    connection_type = "s3",
    connection_options = {"path": "s3://glue-sample-target/output-dir/
medicare_parquet"},
    format = "parquet")
```

データ品質定義言語 (DQDL) リファレンス

Data Quality Definition Language (DQDL) は、AWS Glue Data Quality のルールを定義するために使用するドメイン固有の言語です。

このガイドでは、DQDL の理解を助けるために、この言語の主要な概念について解説します。また、DQDL ルールタイプのリファレンスの中で、構文や例も提供しています。このガイドを使用する前に、AWS Glue Data Quality に精通しておくことをお勧めします。詳細については、「[AWS Glue データ品質](#)」を参照してください。

Note

DynamicRules は ETL AWS Glue でのみサポートされています。

目次

- [DQDL 構文](#)
 - [ルールの構造](#)
 - [複合ルール](#)
 - [複合ルールの仕組み](#)
 - [表現](#)
 - [NULL、EMPTY、WHITESPACES_ONLY のキーワード](#)
 - [Where 句によるフィルタリング](#)
 - [動的ルール](#)
 - [アナライザー](#)
 - [コメント](#)
- [DQDL ルールタイプリファレンス](#)
 - [AggregateMatch](#)
 - [ColumnCorrelation](#)
 - [ColumnCount](#)
 - [ColumnDataType](#)
 - [ColumnExists](#)

- [ColumnLength](#)
- [ColumnNamesMatchPattern](#)
- [ColumnValues](#)
- [Completeness](#)
- [CustomSQL](#)
- [DataFreshness](#)
- [DatasetMatch](#)
- [DistinctValuesCount](#)
- [Entropy](#)
- [IsComplete](#)
- [IsPrimaryKey](#)
- [IsUnique](#)
- [平均値](#)
- [ReferentialIntegrity](#)
- [RowCount](#)
- [RowCountMatch](#)
- [StandardDeviation](#)
- [Sum](#)
- [SchemaMatch](#)
- [Uniqueness](#)
- [UniqueValueRatio](#)
- [DetectAnomalies](#)

DQDL 構文

DQDL のドキュメントでは大文字と小文字が区別されます。また、個々の品質評価ルールをグループ化したルールセットが含まれます。ルールセットを作成するには、Rules という名前 (大文字) の、角括弧で区切られたリストを作成する必要があります。リストの中には、次の例のようにカンマで区切られた 1 つ以上の DQDL ルールを含めます。

```
Rules = [  
    IsComplete "order-id",
```

```
    IsUnique "order-id"  
  ]
```

ルールの構造

DQDL ルールの構造は、そのルールタイプによって異なります。ただし、DQDL の各ルールは、以下のような基本的な形式を取ります。

```
<RuleType> <Parameter> <Parameter> <Expression>
```

RuleType は設定するルールタイプの名前で、大文字と小文字が区別されます。例えば、IsComplete、IsUnique、または CustomSql などです。ルールタイプごとに、ルールのパラメータが異なります。DQDL のルールタイプと、そのパラメータの詳細なリファレンスについては、「[DQDL ルールタイプリファレンス](#)」を参照してください。

複合ルール

DQDL では、ルールを組み合わせる際に、以下の論理演算子を使用できます。これらのルールは複合ルールと呼ばれます。

and

and 論理演算子は、接続するルールがすべて true である場合に限り、結果に true を返します。それ以外の場合、ルールの組み合わせ結果は false になります。and 演算子前後に配置する各ルールは、丸括弧で囲む必要があります。

次の例では、and 演算子を使用して 2 つの DQDL ルールを組み合わせています。

```
(IsComplete "id") and (IsUnique "id")
```

or

or 論理演算子は、接続するルールが 1 つ以上 true であれば、結果に true を返します。or 演算子前後に配置する各ルールは、丸括弧で囲む必要があります。

次の例では、or 演算子を使用して 2 つの DQDL ルールを組み合わせています。

```
(RowCount "id" > 100) or (IsPrimaryKey "id")
```

同じ演算子を使用して複数のルールを接続でき、次のようなルールの組み合わせが可能です。

```
(Mean "Star_Rating" > 3) and (Mean "Order_Total" > 500) and (IsComplete "Order_Id")
```

ただし、論理演算子を 1 つの表現内でまとめることはできません。例えば、以下のように組み合わせることはできません。

```
(Mean "Star_Rating" > 3) and (Mean "Order_Total" > 500) or (IsComplete "Order_Id")
```

複合ルールの仕組み

デフォルトでは、複合ルールはデータセットまたはテーブル全体の個別のルールとして評価されたら、結果が結合されます。つまり、最初に列全体を評価してから演算子を適用します。このデフォルト動作は以下の例で説明します。

```
# Dataset

+-----+-----+
|myCol1|myCol2|
+-----+-----+
|    2|    1|
|    0|    3|
+-----+-----+

# Overall outcome

+-----+-----+-----+
|Rule                                         |Outcome|
+-----+-----+-----+
|(ColumnValues "myCol1" > 1) OR (ColumnValues "myCol2" > 2)|Failed |
+-----+-----+-----+
```

上の例では、AWS Glue Data Quality は最初に失敗に終わる (`ColumnValues "myCol1" > 1`) を評価します。次に (`ColumnValues "myCol2" > 2`) を評価しますが、これも失敗します。両方の結果の組み合わせは `FAILED` として表示されます。

ただし、行全体を評価する必要がある SQL のような動作を好む場合、以下のコードスニペットの `additionalOptions` で示すように `ruleEvaluation.scope` パラメータを明示的に設定する必要があります。

```
object GlueApp {
  val datasource = glueContext.getCatalogSource(
    database="<db>",
    tableName="<table>",
    transformationContext="datasource"
  ).getDynamicFrame()

  val ruleset = """
    Rules = [
      (ColumnValues "age" >= 26) OR (ColumnLength "name" >= 4)
    ]
  """

  val dq_results = EvaluateDataQuality.processRows(
    frame=datasource,
    ruleset=ruleset,
    additionalOptions=JsonOptions("""
      {
        "compositeRuleEvaluation.method":"ROW"
      }
    """)
  )
}
```

AWS Glue Studio と AWS Glue Data Catalog では、次に示すように、ユーザーインターフェイスでこのオプションを簡単に設定できます。

▼ Composite rule settings - *new*

Rule evaluation configuration [Info](#)

Configure how composite rules should work. [Learn more](#) 

Row

The composite rules will behave as single rule evaluating entire row.

Column

The composite rules will evaluate individual rules across the entire dataset and combine the results.

いったん設定すると、複合ルールは行全体を評価する単一のルールとして動作します。次の例では、この動作を示します。

```
# Row Level outcome

+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+
|myCol1|myCol2|DataQualityRulesPass                                     |
DataQualityEvaluationResult|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+
|2      |1      |[(ColumnValues "myCol1" > 1) OR (ColumnValues "myCol2" > 2)]|Passed
|      |      |
|0      |3      |[(ColumnValues "myCol1" > 1) OR (ColumnValues "myCol2" > 2)]|Passed
|      |      |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+
```

一部のルールは全体的な結果がしきい値または比率に依存するため、この機能ではサポートできません。以下に一覧表示します。

比率に依存するルール:

- Completeness
- DatasetMatch
- ReferentialIntegrity
- Uniqueness

しきい値に依存するルール:

次のルールにしきい値が含まれる場合、ルールはサポートされません。ただし、with thresholdに関連しないルールは引き続きサポートされます。

- ColumnDataType
- ColumnValues
- CustomSQL

表現

ルールタイプがブール型で応答を生成しない場合、この応答を取得するには、パラメータで表現を指定する必要があります。例えば、次のルールでは、表現に対応させながら列内のすべての値の mean 値 (平均) をチェックし、真または偽の結果を返します。

```
Mean "colA" between 80 and 100
```

一部のルールタイプ (IsUnique や IsComplete など) は、既にブール型の応答を返しています。

次の表に、DQDL ルールで使用できる式を一覧で示します。

サポートされている DQDL 表現

式	説明	例
<code>=x</code>	ルールタイプの応答が <code>x</code> と等しい場合、true を返します。	<pre>Completeness "colA" = "1.0", ColumnValues "colA" = "2022-06-30"</pre>
<code>!=x</code>	ルールタイプの応答が <code>x</code> と等しくない場合、x は true を返します。	<pre>ColumnValues "colA" != "a", ColumnValues "colA" != "2022-06-30"</pre>
<code>> x</code>	ルールタイプの応答が <code>x</code> より大きい場合、true を返します。	<pre>ColumnValues "colA" > 10</pre>
<code>< x</code>	ルールタイプの応答が <code>x</code> より小さい場合、true を返します。	<pre>ColumnValues "colA" < 1000, ColumnValues "colA" < "2022-06-30"</pre>
<code>>= x</code>	ルールタイプの応答が <code>x</code> 以上の場合、true を返します。	<pre>ColumnValues "colA" >= 10</pre>
<code><= x</code>	ルールタイプの応答が <code>x</code> 以下の場合、true を返します。	<pre>ColumnValues "colA" <= 1000</pre>

式	説明	例
<code>between x and y</code>	ルールタイプの応答が指定された範囲内 (除外) にある場合、true を返します。この式のタイプは、数値および日付タイプにのみ使用できます。	<pre>Mean "colA" between 8 and 100, ColumnValues "colA" between "2022-05-31" and "2022-06-30"</pre>
<code>x と y の間ではない</code>	ルールタイプの応答が指定された範囲外 (含有) にある場合、true を返します。この式のタイプは、数値および日付のタイプにのみ使用できます。	<pre>ColumnValues "colA" not between "2022-05-31" and "2022-06-30"</pre>
<code>in [a, b, c, ...]</code>	ルールタイプの応答が指定されたセットに含まれている場合に、true を返します。	<pre>ColumnValues "colA" in [1, 2, 3], ColumnValues "colA" in ["a", "b", "c"]</pre>
<code>not in [a, b, c, ...]</code>	ルールタイプの応答が指定されたセットに含まれていない場合、true を返します。	<pre>ColumnValues "colA" not in [1, 2, 3], ColumnValues "colA" not in ["a", "b", "c"]</pre>
<code>matches /ab+c/i</code>	ルールタイプの応答が正規表現と一致する場合に、true を返します。	<pre>ColumnValues "colA" matches "[a-zA-Z]*"</pre>
<code>not matches /ab+c/i</code>	ルールタイプのレスポンスが正規表現と一致しtrueない場合、に解決されます。	<pre>ColumnValues "colA" not matches "[a-zA-Z]*"</pre>
<code>now()</code>	日付表現を作成する ColumnValues ルールタイプでのみ機能します。	<pre>ColumnValues "load_date" > (now() - 3 days)</pre>

式	説明	例
<code>matches/in [...]/not matches/not in [...] with threshold</code>	<p>ルール条件に一致する値をパーセンテージで指定します。ColumnValues、ColumnDataType、CustomSQL のルールタイプのみで使用できます。</p>	<pre>ColumnValues "colA" in ["A", "B"] with threshold > 0.8, ColumnValues "colA" matches "[a-zA-Z]*" with threshold between 0.2 and 0.9 ColumnDataType "colA" = "Timestamp" with threshold > 0.9</pre>

NULL、EMPTY、WHITESPACES_ONLY のキーワード

文字列の列に null、空、空白のみの文字列があるかどうかを検証する場合、次のキーワードを使用できます。

- NULL / null – 文字列の列の null 値には、このキーワードは true を返します。

50% 以上のデータに null 値が含まれていない場合、`ColumnValues "colA" != NULL with threshold > 0.5` は true を返します。

null 値または長さが 5 を超えるすべての行に対し、`(ColumnValues "colA" = NULL) or (ColumnLength "colA" > 5)` は true を返します。これには compositeRuleEvaluation 「.method」 = 「ROW」 オプションの使用が必要であることに注意してください。

- EMPTY / empty – 文字列の列に空の文字列 ("") の値には、このキーワードは true を返します。一部のデータ形式は、文字列の列に null を空の文字列に変換します。このキーワードは、データの空の文字列を除外するのに便利です。

行が空、「a」、「b」のいずれかの場合、`(ColumnValues "colA" = EMPTY) or (ColumnValues "colA" in ["a", "b"])` は true を返します。これには compositeRuleEvaluation 「.method」 = 「ROW」 オプションの使用が必要であることに注意してください。

- WHITESPACES_ONLY / whitespaces_only – 文字列の列で空白 (" ") のみの値が含まれる文字列には、このキーワードは true を返します。

行が「a」、「b」、空白のいずれかでない場合、ColumnValues "colA" not in ["a", "b", WHITESPACES_ONLY] は true を返します。

サポートされているルール:

- [ColumnValues](#)

数値または日付ベースの式には、列に null があるかどうかを検証する場合、次のキーワードを使用できます。

- NULL / null – 文字列の列の null 値には、このキーワードは true を返します。

列の日付が 2023-01-01 または null の場合、ColumnValues "colA" in [NULL, "2023-01-01"] は true を返します。

null 値または 1 から 9 までの値が含まれているすべての行に対し、(ColumnValues "colA" = NULL) or (ColumnValues "colA" between 1 and 9) は true を返します。これには、compositeRuleEvaluation 「.method」 = 「ROW」 オプションを使用する必要があります。

サポートされているルール:

- [ColumnValues](#)

Where 句によるフィルタリング

ルールの作成時にデータをフィルタリングできます。これは、条件付きルールを適用する場合に便利です。

```
<DQDL Rule> where "<valid SparkSQL where clause> "
```

フィルターは where キーワードで指定し、その後に引用符 で囲まれた有効な SparkSQL ステートメントを指定する必要があります(“”)。

しきい値を持つルールに where 句を追加するルールの場合、しきい値条件の前に where 句を指定する必要があります。

```
<DQDL Rule> where "valid SparkSQL statement"> with threshold <threshold condition>
```

この構文では、次のようなルールを記述できます。

```
Completeness "colA" > 0.5 where "colB = 10"
ColumnValues "colB" in ["A", "B"] where "colC is not null" with threshold > 0.9
ColumnLength "colC" > 10 where "colD != Concat(colE, colF)"
```

提供された SparkSQL ステートメントが有効であることを確認します。無効な場合、ルールの評価は失敗し、次の形式の `IllegalArgumentException` がスローされます。

```
Rule <DQDL Rule> where "<invalid SparkSQL>" has provided an invalid where clause :
<SparkSQL Error>
```

行レベルのエラーレコード識別が有効になっている場合の句の動作

AWS Glue Data Quality を使用すると、失敗した特定のレコードを特定できます。行レベルの結果をサポートするルールに `where` 句を適用すると、`where` 句で除外された行に というラベルが付けられます `Passed`。

フィルタリングされた行に個別に というラベルを付ける場合は `SKIPPED`、ETL ジョブ `additionalOptions` に以下を設定できます。

```
object GlueApp {
  val datasource = glueContext.getCatalogSource(
    database="<db>",
    tableName="<table>",
    transformationContext="datasource"
  ).getDynamicFrame()

  val ruleset = """
    Rules = [
      IsComplete "att2" where "att1 = 'a'"
    ]
  """

  val dq_results = EvaluateDataQuality.processRows(
    frame=datasource,
    ruleset=ruleset,
    additionalOptions=JsonOptions("""
      {
        "rowLevelConfiguration.filteredRowLabel":"SKIPPED"
      }
    """)
  )
}
```

```
}

```

例として、次のルールとデータフレームを参照してください。

```
IsComplete att2 where "att1 = 'a'"

```

id	att1	att2	行レベルの結果 (デフォルト)	行レベルの結果 (スキップオプション)	コメント
1	a	f	合格	合格	
2	b	d	合格	SKIPPED	はではないため、行att1は除外されます "a"
3	a	null	FAILED	FAILED	
4	a	f	合格	合格	
5	b	null	合格	SKIPPED	はではないため、行att1は除外されます "a"
6	a	f	合格	合格	

動的ルール

動的ルールを作成して、ルールによって生成された現在のメトリクスと履歴値を比較できるようになりました。これらの履歴比較は、式に `last()` 演算子を使用することで可能になります。例えば、現在の実行で行数が同じデータセットの最新の行数よりも多い場合、`RowCount > last()` ルールは成功します。`last()` は、考慮すべき事前メトリクスの数を記述する自然数引数をオプションで取り、`last(k)` では $k \geq 1$ が直近の k メトリクスを参照します。

- データポイントがない場合、`last(k)` はデフォルト値 0.0 を返します。
- 利用可能な `k` メトリクスよりも少ない場合は、`last(k)` はそれ以前のメトリクスをすべて返します。

有効な式を作成するには `last(k)` を使用します。 `k > 1` では、複数の履歴結果を 1 つの数値にまとめる集計関数が必要です。例えば `RowCount > avg(last(5))` は、現在のデータセットの行数が、同じデータセットの直近の 5 つの行数の平均より真に大きいかどうかを確認します。現在のデータセットの行数をリストと有意に比較できないため、`RowCount > last(5)` はエラーを生成します。

サポートされている集計関数:

- `avg`
- `median`
- `max`
- `min`
- `sum`
- `std` (標準偏差)
- `abs` (絶対値)
- `index(last(k), i)` では直近の `k` から `i` 番目に新しい値を選択できます。 `i` はゼロインデックスであるため、`index(last(3), 0)` は最新のデータポイントを返しますが、データポイントが 3 つしかなく、4 番目に新しいデータポイントにインデックスを付けようとするため、`index(last(3), 3)` はエラーになります。

式の例

ColumnCorrelation

- `ColumnCorrelation "colA" "colB" < avg(last(10))`

DistinctValuesCount

- `DistinctValuesCount "colA" between min(last(10))-1 and max(last(10))+1`

数値条件またはしきい値を使用するほとんどのルールタイプは動的ルールをサポートします。使用しているルールタイプで動的ルールがサポートされているかどうかを確認するには、提供されている「[アナライザーとルール](#)」の表を参照してください。

アナライザー

Note

アナライザーは AWS Glue データカタログではサポートされていません。

DQDL ルールは、アナライザーと呼ばれる関数を使用してデータに関する情報を収集します。この情報をルールのブール式に使用して、ルールが成功するか失敗するかを決定します。例えば、RowCount ルール `RowCount > 5` は行数アナライザーを使用してデータセット内の行数を検出し、その数を式と比較して `> 5`、現在のデータセットに 5 行以上が存在するかどうかを確認します。

ルールを作成する代わりに、アナライザーを作成して、異常の検出に使用できる統計を生成させることを推奨する場合があります。このような場合は、アナライザーを作成できます。アナライザーは、次の点でルールとは異なります。

特徴	アナライザー	ルール
ルールセットの一部	はい	あり
統計を生成	はい	あり
観察結果を生成	はい	あり
条件を評価してアサートできる	なし	あり
障害発生時にジョブを停止したり、ジョブの処理を続行したりするなどのアクションを設定できます。	なし	あり

アナライザーはルールなしで独立して存在できるため、アナライザーをすばやく構成し、データ品質ルールを段階的に構築できます。

ルールセットの Analyzers ブロックに入力できるルールタイプもあります。これにより、アナライザーに必要なルールを実行し、条件のチェックを行わずに情報を収集できます。アナライザーの中には、ルールに関連付けられておらず、Analyzers ブロックにしか入力できないものもあります。次の表は、各項目がルールとしてサポートされているか、スタンドアロンアナライザーとしてサポートされているか、および各ルールタイプの詳細を示しています。

アナライザーを使用したルールセットの例

以下のルールセットでは以下のものが使用されます。

- 直近の 3 回のジョブ実行で、データセットが平均値を上回っているかどうかを確認する動的ルール
- データセットの Name 列内の異なる値の数を記録する DistinctValuesCount アナライザー
- Name の最小サイズと最大サイズを経時的に追跡する ColumnLength アナライザー

アナライザーメトリクスの結果は、ジョブ実行の [データ品質] タブで確認できます。

```
Rules = [  
  RowCount > avg(last(3))  
]  
Analyzers = [  
  DistinctValuesCount "Name",  
  ColumnLength "Name"  
]
```

コメント

「#」文字を使用して DQDL ドキュメントにコメントを追加できます。「#」文字の後から行末までに含まれるすべての文字は DQDL に無視されます。

```
Rules = [  
  # More items should generally mean a higher price, so correlation should be  
  positive  
  ColumnCorrelation "price" "num_items" > 0  
]
```

DQDL ルールタイプリファレンス

このセクションでは、AWS Glue Data Quality がサポートする各ルールタイプのリファレンスを提供します。

Note

- 現在、DQDL はリストタイプまたはネストされた列データをサポートしていません。
- 以下の表の括弧内の値は、ルール引数で指定された情報に置き換えられます。
- ルールでは通常、式に追加の引数が必要です。

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
Aggregate Match	売上総額などのサマリーメトリクスを比較して、2つのデータセットが一致しているかをチェックしま	1つ以上の集計	1番目と2番目の集計列名が一致する場合: Column. [Column]. aggregate	あり	いいえ	いいえ	いいえ	いいえ	なし

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
	す。金融機関が、すべてのデータがソースシステムから取り込まれているかを比較する際などに便利です。		目の集計列の名前が異なる場合: Column. [Column1, Column2]. aggregate						

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
AllStatistics	指定した列またはデータセット内のすべての列の複数のメトリクスを収集するスタンドアロンアナライザー。	1つの列名、またはAllColumns「」	すべてのタイプの列の場合: Dataset. .RowColumns Column. [Column]. completers Column. [Column]. liquenes 文字列 値列の その他 のメト リク ス:	なし	はい	いいえ	いいえ	いいえ	なし

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
			Columnlength metrics 数値列のその他のメトリクス: Columnvalues metrics						
ColumnCorrelation	2つの列にどの程度の相関性があるかを確認します。	列名はちょうど2つです	MultiColumnCorrelation. [Column1], [Column2].ColumnCorrelation	はい	はい	いいえ	はい	いいえ	あり
ColumnCount	抜け落ちた列がないかを確認します。	なし	Dataset.ColumnCount	あり	いいえ	いいえ	はい	はい	なし

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
ColumnC aType	列がデータ型に準拠しているかをチェックします。	列名は1つだけです	Column. [C olumn]. lumnDat ype.Con iance	あり	いいえ	なし	はい (行レベルのしきい値式の場合)	なし	あり

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
ColumnEs	データセットに列が存在するかをチェックします。これにより、セルフサービスのデータプラットフォームを構築しているユーザーは、特定の列が利用可能であることを確	列名は1つだけです	該当なし	はい	いいえ	いいえ	いいえ	いいえ	なし

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
	認できません。								

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
ColumnLength	データの長さが一貫しているかをチェックします。	列名は1つだけです	Column.[Column].maximumLength Column.[Column].minimumLength 行レベルのしきい値が指定されている場合のその他のメトリクス: Column.[Column].ColumnValues.Compance	はい	あり	はい (行レベルのしきい値が指定されている場合)	なし	はい。最小長と最大長を分析して観察結果のみを生成します。	あり

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
ColumnNamesMatchPattern	列名が定義済みのパターンと一致しているかをチェックします。ガバナンスチームが列名の一貫性を保つ際に便利です。	列名の正規表現	Dataset.ColumnNamesMatchFunction	あり	いいえ	いいえ	いいえ	いいえ	なし

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
ColumnValues	データが定義済みの値と一致しているかをチェックします。このルールは正規表現に対応しています。	列名は1つだけです	ColumnMaximum ColumnMinimum 行レベルのしきい値が指定されている場合のその他のメトリクス: ColumnMaximumValue.Complexity	はい	あり	はい (行レベルのしきい値が指定されている場合)	なし	はい。最小値と最大値を分析して観察結果のみを生成します。	あり

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
Completeness	データに空白またはNULLがないかをチェックします。	列名は1つだけです	ColumnCompleteness	はい	はい	はい	はい	はい	あり
CustomSQL	ユーザーは、ほぼすべてのタイプのデータ品質チェックをSQLに実装できます。	SQLステートメント (オプション) 行レベルのしきい値	DatasetCustomSQL 行レベルのしきい値が指定されている場合のその他のメトリクス: DatasetCustomSQLCompleteness	あり	なし	はい (行レベルのしきい値が指定されている場合)	あり	いいえ	なし

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
DataFreshness	データが最新であることをチェックします。	列名は1つだけです	Column.[Column].DataFreshness.Compliance	あり	いいえ	はい	いいえ	いいえ	あり
DatasetMatch	2つのデータセットを比較して、同期しているかを識別します。	参照データセットの名前の列のマッピング (オプション)一致を確認する列	Dataset[RefererDatasetias].DatasetMatch	あり	いいえ	はい	はい	いいえ	なし

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
DistinctValuesCount	重複する値がないかをチェックします。	列名は1つだけです	Column.[Column].distinctValuesCount	はい	はい	はい	はい	はい	あり
DetectAnomalies	別のルールタイプで報告されたメトリクスに異常がないかチェックします。	ルールタイプ	ルールタイプ引数で報告されたメトリクス(1つまたは複数)	あり	いいえ	いいえ	いいえ	いいえ	なし
エントロピー	データのエントロピーをチェックします。	列名は1つだけです	Column.[Column].entropy	はい	はい	いいえ	はい	いいえ	あり

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
IsComplete	すべてのデータが完全であることをチェックします。	列名は1つだけです	Column.[Column].mpleters	あり	いいえ	はい	いいえ	いいえ	あり
IsPrimaryKey	列がプライマリキー (NULL および一意ではない) であるかをチェックします。	列名は1つだけです	1列の場合: Column.[Column].iquenes 複数列の場合: Multicolumn[Column].elimitedcolumns].iquenes	あり	いいえ	はい	いいえ	いいえ	あり

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
IsUnique	データがすべて一意であるかをチェックします。	列名は1つだけです	Column.[Column].iquenes	あり	いいえ	はい	いいえ	いいえ	あり
平均値	平均値が、設定済みのしきい値と一致するかをチェックします。	列名は1つだけです	Column.[Column].an	はい	はい	はい	はい	いいえ	あり

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
ReferentialIntegrity	2つのデータセットに参照整合性があるかをチェックします。	データセットの1つまたは複数の列名 参照データセットの1つまたは複数の列名	Column.[ReferentialIntegrity]	あり	いいえ	はい	はい	いいえ	なし
RowCount	レコード数がしきい値と一致するかをチェックします。	なし	Dataset.RowCount	はい	はい	いいえ	はい	はい	あり

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
RowCountMatch	2つのデータセットのレコード数が一致するかをチェックします。	参照データセットのエイリアス	Dataset [ReferencedDataset].RowCountMatch	あり	いいえ	いいえ	はい	いいえ	なし
StandardDeviation	標準偏差がしきい値と一致するかをチェックします。	列名は1つだけです	Column [Column].StandardDeviation	はい	はい	はい	はい	いいえ	あり

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
SchemaMatch	2つのデータセットのスキーマが一致するかをチェックします。	参照データセットのエイリアス	Dataset [RefererDataset].SchemaMatch	あり	いいえ	いいえ	はい	いいえ	なし
合計	合計が、設定済みのしきい値と一致するかをチェックします。	列名は1つだけです	Column [Column].m	はい	はい	いいえ	はい	いいえ	あり

Ruletype	説明	引数	報告されたメトリクス	ルールとしてサポートされていますか？	アナライザーとしてサポートされていますか？	行レベルの結果を返しますか？	動的ルールをサポートしますか？	観察結果を生成	句の構文をサポートしますか？
Uniques	データセットの一貫性がしきい値と一致するかをチェックします。	列名は1つだけです	ColumnUniques	はい	はい	はい	はい	いいえ	あり
UniqueValueRatio	一意の値の比率がしきい値と一致するかをチェックします。	列名は1つだけです	ColumnUniqueValueRatio	はい	はい	はい	はい	いいえ	あり

トピック

- [AggregateMatch](#)
- [ColumnCorrelation](#)
- [ColumnCount](#)

- [ColumnDataType](#)
- [ColumnExists](#)
- [ColumnLength](#)
- [ColumnNamesMatchPattern](#)
- [ColumnValues](#)
- [Completeness](#)
- [CustomSQL](#)
- [DataFreshness](#)
- [DatasetMatch](#)
- [DistinctValuesCount](#)
- [Entropy](#)
- [IsComplete](#)
- [IsPrimaryKey](#)
- [IsUnique](#)
- [平均値](#)
- [ReferentialIntegrity](#)
- [RowCount](#)
- [RowCountMatch](#)
- [StandardDeviation](#)
- [Sum](#)
- [SchemaMatch](#)
- [Uniqueness](#)
- [UniqueValueRatio](#)
- [DetectAnomalies](#)

AggregateMatch

特定の式を参照して、2つの列の集計の比率をチェックします。このルールタイプは、複数のデータセットで機能します。2つの列の集計が評価され、1番目の列の集計結果を2番目の列の集計結果で割ることで比率が算出されます。比率と指定された式が照合され、ブール型応答が生成されます。

構文

列の集計

```
ColumnExists <AGG_OPERATION> (<OPTIONAL_REFERENCE_ALIAS>.<COL_NAME>)
```

- AGG_OPERATION – 集計に使用する操作。現在サポートされている形式は、sum と avg です。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- OPTIONAL_REFERENCE_ALIAS – 列がプライマリデータセットではなく参照データセットのものである場合は、このパラメータを指定する必要があります。AWS Glue データカタログでこのルールを使用している場合、参照エイリアスは「<database_name>.<table_name>.<column_name>」の形式に従う必要があります。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- COL_NAME — 集計する列の名前。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

例: 平均

```
"avg(rating)"
```

例: 合計

```
"sum(amount)"
```

例: 参照データセットの列の平均

```
"avg(reference.rating)"
```

ルール

```
AggregateMatch <AGG_EXP_1> <AGG_EXP_2> <EXPRESSION>
```

- AGG_EXP_1 – 最初の列の集計。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- AGG_EXP_2 – 2 番目の列の集計。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 合計を使用した Aggregate Match

次のルール例では、amount 列の値の合計が、total_amount 列の値の合計と完全に等しいかをチェックします。

```
AggregateMatch "sum(amount)" "sum(total_amount)" = 1.0
```

例: 平均を使用した Aggregate Match

次のルール例では、ratings 列の値の平均が、reference データセットの ratings 列における値の平均の、90% 以上と等しいかをチェックします。参照データセットは、ETL またはデータカタログ体験の、追加のデータソースとして提供されます。

AWS Glue ETL では、以下を使用できます。

```
AggregateMatch "avg(ratings)" "avg(reference.ratings)" >= 0.9
```

AWS Glue データカタログでは、以下を使用できます。

```
AggregateMatch "avg(ratings)" "avg(database_name.tablename.ratings)" >= 0.9
```

Null 動作

AggregateMatch ルールは、集計方法 (合計/平均) の計算で NULL 値の行を無視します。例:

```
+---+-----+
|id |units   |
+---+-----+
|100|0      |
|101|null  |
|102|20     |
|103|null  |
|104|40     |
+---+-----+
```

units 列の平均は $(0 + 20 + 40) / 3 = 20$ になります。行 101 および 103 はこの計算の対象になりません。

ColumnCorrelation

2 つの列間の相関関係を特定の式と照合します。AWS Glue Data Quality は、ピアソン相関係数を使用して 2 つの列間の線形相関を測定します。この結果では、-1 から 1 までの数値により、相関関係の強さと方向性が示されます。

[Syntax (構文)]

```
ColumnCorrelation <COL_1_NAME> <COL_2_NAME> <EXPRESSION>
```

- COL_1_NAME – 品質評価ルールの評価に対応させる、最初の列の名前。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- COL_2_NAME – 品質評価ルールの評価に対応させる、2 番目の列の名前。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 列の相関

次のルール例では、列 height と weight の間の相関係数が、強い正の (係数値が 0.8 より大きい) 相関を示しているかをチェックします。

```
ColumnCorrelation "height" "weight" > 0.8
```

```
ColumnCorrelation "weightinkgs" "Salary" > 0.8 where "weightinkgs" > 40
```

動的ルールの例

- ColumnCorrelation "colA" "colB" between min(last(10)) and max(last(10))
- ColumnCorrelation "colA" "colB" < avg(last(5)) + std(last(5))

Null 動作

ColumnCorrelation ルールは、相関の計算で NULL 値を含む行を無視します。例:

```
+---+-----+
|id |units  |
+---+-----+
|100|0      |
|101|null  |
|102|20     |
|103|null  |
|104|40     |
+---+-----+
```

行 101 および 103 は無視され、ColumnCorrelation は 1.0 になります。

ColumnCount

特定の式を参照しながら、プライマリデータセットの列数を確認します。この式では、> や < などの演算子を使用して、行数または行の範囲を指定できます。

[Syntax (構文)]

```
ColumnCount <EXPRESSION>
```

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 列数の数値による確認

次のルール例では、列数が特定の範囲内にあるかを確認します。

```
ColumnCount between 10 and 20
```

動的ルールの例

- `ColumnCount >= avg(last(10))`
- `ColumnCount between min(last(10))-1 and max(last(10))+1`

ColumnDataType

特定の列における値の固有のデータ型を、指定の、想定された型と照合します。with threshold 式を承諾し、列内の値のサブセットをチェックします。

[Syntax (構文)]

```
ColumnDataType <COL_NAME> = <EXPECTED_TYPE>  
ColumnDataType <COL_NAME> = <EXPECTED_TYPE> with threshold <EXPRESSION>
```

- **COL_NAME** – データ品質ルールを評価する対象となる列の名前。

サポートされている列の型: String 型

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- **EXPECTED_TYPE** — 列内の想定されている値の型。

サポートされている値: Boolean、Date、Timestamp、Integer、Double、Float、Long

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- **EXPRESSION** – 想定されている型の、値の割合を指定するオプションの式。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

例: 文字列としての、列データ型の整数。

次のルール例では、所定の列の値 (string 型) が本当に整数かをチェックします。

```
ColumnDataType "colA" = "INTEGER"
```

例: 文字列としての列データ型の整数が、値のサブセットをチェックします。

次のルール例は、所定の列の値 (string 型) の 90% 以上が本当に整数であるかを確認します。

```
ColumnDataType "colA" = "INTEGER" with threshold > 0.9
```

ColumnExists

列が存在するかどうかを確認します。

[Syntax (構文)]

```
ColumnExists <COL_NAME>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: 任意の型

例: 列の存在

次のルール例では、Middle_Name という名前の列が存在するかどうかを確認します。

```
ColumnExists "Middle_Name"
```

ColumnLength

列内にある各行の長さが、特定の表現に適合しているかどうかを確認します。

[Syntax (構文)]

```
ColumnLength <COL_NAME><EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされる型: String

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 列内の行の長さ

次のルール例では、Postal_Code という名前の列内にある各行の値が、5 文字の長さであるかどうかを確認します。

```
ColumnLength "Postal_Code" = 5
ColumnLength "weightinkgs" = 2 where "weightinkgs" > 10"
```

Null 動作

ColumnLength ルールは、NULL を長さ 0 の文字列として扱います。NULL 行には、次の内容が適用されます。

```
ColumnLength "Postal_Code" > 4 # this will fail
```

```
ColumnLength "Postal_Code" < 6 # this will succeed
```

次の複合ルールの例では、NULL 値を明示的に満たさない方法が示されています。

```
(ColumnLength "Postal_Code" > 4) AND (ColumnValues != NULL)
```

ColumnNamesMatchPattern

プライマリデータセットの全列の名前が、所定の正規表現と一致するかをチェックします。

[Syntax (構文)]

```
ColumnNamesMatchPattern <PATTERN>
```

- PATTERN – データ品質ルールを評価する際の基準となるパターン。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

例: 列名がパターンと一致する

次のルール例では、すべての列がプレフィックス「aws_」から始まっているかを確認します。

```
ColumnNamesMatchPattern "aws_.*"
```

```
ColumnNamesMatchPattern "aws_.*" where "weightinkgs > 10"
```

ColumnValues

列内の値に対して式を実行します。

[Syntax (構文)]

```
ColumnValues <COL_NAME> <EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: 任意の型

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 許可される値

次のルール例では、指定された列の各値が、許可された値 (null、空、空白のみの文字列を含む) のセットの中に含まれているかどうかを確認します。

```
ColumnValues "Country" in [ "US", "CA", "UK", NULL, EMPTY, WHITESPACES_ONLY ]  
ColumnValues "gender" in ["F", "M"] where "weightinkgs < 10"
```

例: 正規表現

次のルール例では、列内の値が正規表現と照合されます。

```
ColumnValues "First_Name" matches "[a-zA-Z]*"
```

例: 日付値

次のルール例では、日付列内の値が日付表現に適合しているかをチェックします。

```
ColumnValues "Load_Date" > (now() - 3 days)
```

例: 数値

次のルール例では、列内の値が特定の数値制約に適合しているかどうかをチェックします。

```
ColumnValues "Customer_ID" between 1 and 2000
```

Null 動作

すべての ColumnValues ルール (!= および NOT IN 以外) では、NULL 行はルールを満たしません。null 値が原因でルールが失敗した場合、失敗理由には次の内容が表示されます。

```
Value: NULL does not meet the constraint requirement!
```

次の複合ルールの例では、NULL 値を明示的に満たす方法が示されています。

```
(ColumnValues "Age" > 21) OR (ColumnValues "Age" = NULL)
```

!= および not in 構文を使用する否定された ColumnValues ルールは、NULL 行に対して渡されません。例:

```
ColumnValues "Age" != 21
```

```
ColumnValues "Age" not in [21, 22, 23]
```

次の例では、NULL 値を明示的に満たさない方法を示されています。

```
(ColumnValues "Age" != 21) AND (ColumnValues "Age" != NULL)
```

```
ColumnValues "Age" not in [21, 22, 23, NULL]
```

Completeness

列内の完全な (NULL 以外の) 値を指定された式と照合し、パーセンテージで示します。

[Syntax (構文)]

```
Completeness <COL_NAME> <EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: 任意の型

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: NULL 値のパーセンテージ

次のルール例では、列内に完全な値が 95% 以上存在するかどうかをチェックします。

```
Completeness "First_Name" > 0.95  
Completeness "First_Name" > 0.95 where "weightinkgs" > 10"
```

動的ルールの例

- Completeness "colA" between min(last(5)) - 1 and max(last(5)) + 1
- Completeness "colA" <= avg(last(10))

Null 動作

CSV データ形式に関するメモ: CSV 列の空白行は、複数の動作を表示することがあります。

- 列が String 型の場合、空白行は空の文字列として認識され、Completeness ルールを満たしません。
- 列が Int のような別のデータ型の場合、空白行は NULL として認識され、Completeness ルールを満たしません。

CustomSQL

このルールタイプは、次の 2 つのユースケースをサポートするように拡張されました。

- データセットに対してカスタム SQL ステートメントを実行し、その戻り値を指定された式と対応させて確認します。
- SELECT ステートメントで列名を指定するカスタム SQL ステートメントを実行し、それを何らかの条件と比較して行レベルの結果を取得します。

[Syntax (構文)]

```
CustomSql <SQL_STATEMENT> <EXPRESSION>
```

- SQL_STATEMENT – 二重引用符で囲まれた単一の数値を返す SQL ステートメント。
- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: ルール全体の結果を取得するカスタム SQL

このルール例では、データセットのレコード数を取得するための、SQL ステートメントを使用します。その後、このルールは、レコード数が 10 から 20 の間であることを確認します。

```
CustomSql "select count(*) from primary" between 10 and 20
```

例: 行レベルの結果を取得するカスタム SQL

このサンプルルールでは、SELECT ステートメントで列名を指定し、それを何らかの条件と比較して行レベルの結果を取得する SQL ステートメントを使用しています。しきい値条件式は、ルール全体が失敗するまでに失敗するレコード数のしきい値を定義します。ルールには条件とキーワードの両方を一緒に含めることはできないことに注意してください。

```
CustomSql "select Name from primary where Age > 18"
```

または

```
CustomSql "select Name from primary where Age > 18" with threshold > 3
```

Important

primary エイリアスは、評価するデータセットの名前の代替です。コンソールでビジュアル ETL ジョブを使用する場合、primary は常に、EvaluateDataQuality.apply() 変換に渡されている DynamicFrame を表します。AWS Glue Data Catalog を使用してテーブルに対してデータ品質タスクを実行する場合、はテーブルprimaryを表します。

AWS Glue データカタログを使用している場合は、実際のテーブル名を使用することもできます。

```
CustomSql "select count(*) from database.table" between 10 and 20
```

複数のテーブルを結合して、異なるデータ要素を比較することもできます。

```
CustomSql "select count(*) from database.table inner join database.table2 on id1 = id2"
          between 10 and 20
```

AWS Glue ETL では、CustomSQL はデータ品質チェックに失敗したレコードを特定できます。これを機能させるには、データ品質を評価する主テーブルの一部であるレコードを返す必要があります。クエリの一部として返されたレコードは成功と見なされ、返されなかったレコードは不合格と見なされます。

次のルールにより、経過日数が 100 未満のレコードは成功と見なされ、それ以上のレコードは不合格とマークされます。

```
CustomSql "select id from primary where age < 100"
```

次の CustomSQL ルールは、レコードの 50% の経過日数が 10 を超えた場合は合格とし、失敗したレコードも特定します。この CustomSQL によって返されたレコードは合格と見なされ、返されなかったレコードは不合格と見なされます。

```
CustomSQL "select ID, CustomerID from primary where age > 10" with threshold > 0.5
```

メモ: データセットにないレコードを返すと、CustomSQL ルールは失敗します。

DataFreshness

現在の時刻と日付列の値との差を評価して、列内のデータがどの程度新しいかをチェックします。このルールタイプで時間ベースの式を指定することで、列の値を最新に保つことができます。

[Syntax (構文)]

```
DataFreshness <COL_NAME> <EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされる型: Date

- EXPRESSION – 時間または日付の数値表現。表現の中では、時間単位を指定する必要があります。

例: データの新しさ

次のルール例では、データの新しさをチェックします。

```
DataFreshness "Order_Date" <= 24 hours
DataFreshness "Order_Date" between 2 days and 5 days
```

Null 動作

DataFreshness ルールは、NULL 値がある行を満たしません。null 値が原因でルールが失敗した場合、失敗理由には次の内容が表示されます。

```
80.00 % of rows passed the threshold
```

満たさなかった 20% の行には、NULL を含む行が対象となる場合。

次の複合ルールの例では、NULL 値を明示的に満たす方法が示されています。

```
(DataFreshness "Order_Date" <= 24 hours) OR (ColumnValues "Order_Date" = NULL)
```

Amazon S3 オブジェクトのデータの鮮度

Amazon S3 ファイルの作成時間に基づいて、データの鮮度を検証する必要がある場合があります。これを行うには、次のコードを使用してタイムスタンプを取得してデータフレームに追加し、データの鮮度のチェックを適用します。

```
df = glueContext.create_data_frame.from_catalog(database = "default", table_name =
  "mytable")
df = df.withColumn("file_ts", df["_metadata.file_modification_time"])

Rules = [
  DataFreshness "file_ts" < 24 hours
]
```

DatasetMatch

プライマリデータセットのデータが参照データセットのデータと一致するかをチェックします。これら 2 つのデータセットは、入力されたキー列マッピングを使用して結合されます。これらの列のみでデータが等しいかを確認する場合は、追加の列マッピングを提供できます。が機能DataSetMatchするためには、結合キーは一意で、NULL (プライマリキー) であってはなりません。これらの条件を満たしていないと、「Provided key map not suitable for given data frames」というエラーメッセージが表示されます。一意の結合キーを使用できない場合は、などの他のルールタイプを使用して概要データにAggregateMatch一致させることを検討してください。

[Syntax (構文)]

```
DatasetMatch <REFERENCE_DATASET_ALIAS> <JOIN_CONDITION_WITH_MAPPING> <OPTIONAL_MATCH_COLUMN_MAPPINGS> <EXPRESSION>
```

- REFERENCE_DATASET_ALIAS – プライマリデータセットのデータと比較する、参照データセットのエイリアス。
- KEY_COLUMN_MAPPINGS – データセットのキーとなる列名のカンマ区切りリスト。列名が両方のデータセットで同一でない場合は、-> で区切る必要があります。
- OPTIONAL_MATCH_COLUMN_MAPPINGS – 特定の列のみでデータの一致を確認したい場合はこのパラメータを指定します。これは、キー列のマッピングと同じ構文を使用します。このパラメータを指定しないと、残りすべての列のデータと照合されます。残りの非キー列は、両方のデータセットで同じ名前を持っている必要があります。
- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: ID 列を使用したセットデータセットの照合

次のルール例では、[ID] 列を使用して 2 つのデータセットを結合し、プライマリデータセットの 90% 以上が参照データセットと一致していることを確認します。この場合、すべての列を比較します。

```
DatasetMatch "reference" "ID" >= 0.9
```

例: 複数のキー列を使用したセットデータセットの照合

次の例では、プライマリデータセットと参照データセットで、キー列の名前が異なります。ID_1 と ID_2 は、結合して、プライマリデータセットの複合キーを構成します。ID_ref1 と ID_ref2 は、結合して、参照データセットの複合キーを構成します。このシナリオでは、特殊な構文を使用して列名を指定できます。

```
DatasetMatch "reference" "ID_1->ID_ref1,ID_ref2->ID_ref2" >= 0.9
```

例: 複数のキー列を使用してセットデータセットを照合し、特定の列が一致していることを確認する

この例は、前述の例に基づいています。金額を含む列のみが一致しているかを、確認しようとしています。この列の名前は、プライマリデータセットでは Amount1、参照データセットでは Amount2 です。これらを完全に一致させたいと思います。

```
DatasetMatch "reference" "ID_1->ID_ref1,ID_ref2->ID_ref2" "Amount1->Amount2" >= 0.9
```

DistinctValuesCount

列内の異なる値の数を、指定された式に対応させて確認します。

[Syntax (構文)]

```
DistinctValuesCount <COL_NAME> <EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: 任意の型

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 列内の異なる値のカウント

次のルール例では、State という名前の列に 3 つ以上の異なる値が含まれていることを確認します。

```
DistinctValuesCount "State" > 3  
DistinctValuesCount "Customer_ID" < 6 where "Customer_ID < 10"
```

動的ルールの例

- DistinctValuesCount "colA" between avg(last(10))-1 and avg(last(10))+1
- DistinctValuesCount "colA" <= index(last(10),2) + std(last(5))

Entropy

列のエントロピー値が、特定の式と一致するかどうかを確認します。エントロピーは、メッセージに含まれる情報のレベルを測定します。エントロピーでは、列内の値に関する特定の確率分布に基づき、値を区別するのに必要なビット数を表します。

[Syntax (構文)]

```
Entropy <COL_NAME> <EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: 任意の型

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 列のエントロピー

次のルール例では、Feedback という名前の列のエントロピー値が、1 より大きいことを確認します。

```
Entropy "Star_Rating" > 1
Entropy "First_Name" > 1 where "Customer_ID < 10"
```

動的ルールの例

- Entropy "colA" < max(last(10))
- Entropy "colA" between min(last(10)) and max(last(10))

IsComplete

列のすべての値が完全 (NULL 以外) かどうかをチェックします。

[Syntax (構文)]

```
IsComplete <COL_NAME>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: 任意の型

例: NULL 値

次の例では、email という名前が付いた列内で、値がすべて NULL 以外であるかどうかをチェックします。

```
IsComplete "email"
IsComplete "Email" where "Customer_ID between 1 and 50"
```

```
IsComplete "Customer_ID" where "Customer_ID < 16 and Customer_ID != 12"  
IsComplete "passenger_count" where "payment_type<>0"
```

Null 動作

CSV データ形式に関するメモ: CSV 列の空白行は、複数の動作を表示することがあります。

- 列が String 型の場合、空白行は空の文字列として認識され、Completeness ルールを満たしません。
- 列が Int のような別のデータ型の場合、空白行は NULL として認識され、Completeness ルールを満たしません。

IsPrimaryKey

列にプライマリキーが含まれているかどうかを確認します。列内の値がすべて一意であり、かつそれらすべてが完全な (NULL 以外) 場合、その列にはプライマリキーが含まれます。

[Syntax (構文)]

```
IsPrimaryKey <COL_NAME>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: 任意の型

例; プライマリキー

次のルール例では、Customer_ID という名前の列にプライマリキーが含まれているかどうかを確認します。

```
IsPrimaryKey "Customer_ID"  
IsPrimaryKey "Customer_ID" where "Customer_ID < 10"
```

例:複数の列を持つ主プライマリキー。以下のいずれの例も有効です。

```
IsPrimaryKey "colA" "colB"  
IsPrimaryKey "colA" "colB" "colC"  
IsPrimaryKey colA "colB" "colC"
```

IsUnique

列のすべての値が一意であるかどうかをチェックし、ブール型の値を返します。

[Syntax (構文)]

```
IsUnique <COL_NAME>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: 任意の型

例: 列での一意の値

次のルール例では、email という名前の列で、すべての値が一意であるかどうかをチェックします。

```
IsUnique "email"  
IsUnique "Customer_ID" where "Customer_ID < 10"]
```

平均値

列内にあるすべての値の mean 値 (平均) が、特定の表現に適合するかどうかを確認します。

[Syntax (構文)]

```
Mean <COL_NAME> <EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 平均値

次のルール例では、列のすべての値の平均が、しきい値を超えているかどうかを確認します。

```
Mean "Star_Rating" > 3
Mean "Salary" < 6200 where "Customer_ID < 10"
```

動的ルールの例

- `Mean "colA" > avg(last(10)) + std(last(2))`
- `Mean "colA" between min(last(5)) - 1 and max(last(5)) + 1`

Null 動作

Mean ルールは、平均の計算で NULL 値を含む行を無視します。例:

```
+---+-----+
|id |units  |
+---+-----+
|100|0      |
|101|null  |
|102|20     |
|103|null  |
|104|40     |
+---+-----+
```

units 列の平均は $(0 + 20 + 40) / 3 = 20$ になります。行 101 および 103 はこの計算の対象になりません。

ReferentialIntegrity

プライマリデータセットに存在する一連の列の値が、どの程度、参照データセットの一連の列における値のサブセットになっているかをチェックします。

[Syntax (構文)]

```
ReferentialIntegrity <PRIMARY_COLS> <REFERENCE_DATASET_COLS> <EXPRESSION>
```

- PRIMARY_COLS — プライマリデータセット内の列名のカンマ区切りリスト。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- REFERENCE_DATASET_COLS – このパラメータには、ピリオドで区切られた 2 つの要素が含まれます。前半の部分は、参照データセットのエイリアスです。後半の部分は、中かっこで囲まれた参照データセット内の列名のカンマ区切りリストです。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 郵便番号列の参照整合性の確認

次のルール例では、zipcode 列の値の 90% 以上が reference データセットの zipcode 列にあることを確認します。

```
ReferentialIntegrity "zipcode" "reference.zipcode" >= 0.9
```

例: 都市列と州列の参照整合性の確認

次の例では、都市と州の情報を含む列が、プライマリデータセットと参照データセットの中にあります。列の名前は、両方のデータセットで異なっています。このルールは、プライマリデータセットの列の値セットが、参照データセットの列の値セットと完全に等しいかをチェックします。

```
ReferentialIntegrity "city,state" "reference.{ref_city,ref_state}" = 1.0
```

動的ルールの例

- ReferentialIntegrity "city,state" "reference.{ref_city,ref_state}" > avg(last(10))
- ReferentialIntegrity "city,state" "reference.{ref_city,ref_state}" between min(last(10)) - 1 and max(last(10)) + 1

RowCount

特定の表現を参照しながら、データセットの行数を確認します。この式では、> や < などの演算子を使用して、行数または行の範囲を指定します。

[Syntax (構文)]

```
RowCount <EXPRESSION>
```

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 行数の数値による確認

次のルール例では、行数が特定の範囲内にあるかどうかを確認します。

```
RowCount between 10 and 100  
RowCount between 1 and 50 where "Customer_ID < 10"
```

動的ルールの例

```
RowCount > avg(lats(10)) *0.8
```

RowCountMatch

プライマリデータセットの行数と参照データセットの行数の比率を、所定の式に照らして確認します。

[Syntax (構文)]

```
RowCountMatch <REFERENCE_DATASET_ALIAS> <EXPRESSION>
```

- REFERENCE_DATASET_ALIAS – 行数の比較に使用する参照データセットのエイリアス。
列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)
- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 参照データセットに対する行数の確認

次のルール例は、プライマリデータセットの行数が、参照データセットの行数の 90% 以上あるかをチェックします。

```
RowCountMatch "reference" >= 0.9
```

StandardDeviation

列のすべての値の標準偏差が、特定の表現に適合するかを確認します。

[Syntax (構文)]

```
StandardDeviation <COL_NAME> <EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 標準偏差

次のルール例では、colA という名前の列内の値の標準偏差が、指定された値より小さいかどうかを確認します。

```
StandardDeviation "Star_Rating" < 1.5
StandardDeviation "Salary" < 3500 where "Customer_ID < 10"
```

動的ルールの例

- StandardDeviation "colA" > avg(last(10)) + 0.1
- StandardDeviation "colA" between min(last(10)) - 1 and max(last(10)) + 1

Null 動作

StandardDeviation ルールは、標準偏差の計算で NULL 値を含む行を無視します。例:

```
+---+-----+-----+
|id |units1      |units2      |
+---+-----+-----+
|100|0           |0           |
|101|null      |0           |
|102|20          |20          |
|103|null      |0           |
```

```
|104|40      |40      |
+---+-----+-----+
```

units1 列の標準偏差は行 101 および 103 を対象にせず、結果は 16.33 になります。units2 列の標準偏差は 16 になります。

Sum

列のすべての値の合計値を、特定の表現を参照しながら確認します。

[Syntax (構文)]

```
Sum <COL_NAME> <EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 合計

次のルール例では、列のすべての値の合計が、特定のしきい値を超えているかどうかを確認します。

```
Sum "transaction_total" > 500000
Sum "Salary" < 55600 where "Customer_ID < 10"
```

動的ルールの例

- Sum "ColA" > avg(last(10))
- Sum "colA" between min(last(10)) - 1 and max(last(10)) + 1

Null 動作

Sum ルールは、合計の計算で NULL 値を含む行を無視します。例:

```
+---+-----+
|id |units      |
```

```
+---+-----+
|100|0      |
|101|null   |
|102|20     |
|103|null   |
|104|40     |
+---+-----+
```

units 列の合計は行 101 および 103 を対象にせず、 $(0 + 20 + 40) = 60$ になります。

SchemaMatch

プライマリデータセットのスキーマが参照データセットのスキーマと一致しているかをチェックします。スキーマのチェックは列ごとに行われます。名前と型が同じであれば、2つの列のスキーマは一致します。行の順序は関係ありません。

[Syntax (構文)]

```
SchemaMatch <REFERENCE_DATASET_ALIAS> <EXPRESSION>
```

- REFERENCE_DATASET_ALIAS – スキーマの比較に使用する参照データセットのエイリアスです。

列でサポートされている型: Byte (バイト)、Decimal (十進数)、Double (倍精度浮動小数点数)、Float (浮動小数点数)、Integer (整数)、Long (整数)、Short (整数)

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: SchemaMatch

次のルール例は、プライマリデータセットのスキーマが参照データセットのスキーマと厳密に一致しているかをチェックします。

```
SchemaMatch "reference" = 1.0
```

Uniqueness

特定の表現と対応させて、列内にある一意の値の個数を、パーセンテージにより確認します。一意の値とは、1つのみ存在する値です。

[Syntax (構文)]

```
Uniqueness <COL_NAME> <EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: 任意の型

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 一意性の割り合い

次のルール例では、列内での一意の値の割合が、特定の数値基準と一致するかどうかを確認します。

```
Uniqueness "email" = 1.0  
Uniqueness "Customer_ID" != 1.0 where "Customer_ID < 10"
```

動的ルールの例

- Uniqueness "colA" between min(last(10)) and max(last(10))
- Uniqueness "colA" >= avg(last(10))

UniqueValueRatio

特定の表現を参照して、列での一意な値の比率をチェックします。個別値の比率は、列内の一意の値の個数を、個別なすべての値の個数で割ったものです。一意の値とは 1 つだけ含まれる値であり、一方、個別の値は少なくとも 1 つ含まれている値です。

例えば、[a, a, b] のセットには 1 つの一意の値 (b) と 2 つの個別の値 (a と b) が含まれています。したがって、このセットの一意な値の比率は $\frac{1}{2} = 0.5$ です。

[Syntax (構文)]

```
UniqueValueRatio <COL_NAME> <EXPRESSION>
```

- COL_NAME – データ品質ルールを評価する対象となる列の名前。

列でサポートされている型: 任意の型

- EXPRESSION – ルールタイプの応答に対して実行し、論地値を生成するための式。詳細については、「[表現](#)」を参照してください。

例: 一意な値の比率

この例では、列の一意な値の比率を、値の範囲と比較します。

```
UniqueValueRatio "test_score" between 0 and 0.5
UniqueValueRatio "Customer_ID" between 0 and 0.9 where "Customer_ID < 10"
```

動的ルールの例

- UniqueValueRatio "colA" > avg(last(10))
- UniqueValueRatio "colA" <= index(last(10),2) + std(last(5))

DetectAnomalies

特定のデータ品質ルールの異常を検出します。DetectAnomalies ルールを実行するたびに、特定のルールの評価値が保存されます。十分なデータが収集されると、異常検出アルゴリズムはその特定のルールのすべての履歴データを取得し、異常検出を実行します。異常が検出されるとDetectAnomalies、ルールは失敗します。どのような異常が検出されたかについての詳細は、「[観察結果](#)」で確認できます。

[Syntax (構文)]

```
DetectAnomalies <RULE_NAME> <RULE_PARAMETERS>
```

RULE_NAME – 異常の評価、検出を希望するルールの名前。サポートされているルール:

- "RowCount"
- "Completeness"
- "Uniqueness"
- "Mean"
- "Sum"

- "StandardDeviation"
- "Entropy"
- "DistinctValuesCount"
- "UniqueValueRatio"
- "ColumnLength"
- "ColumnValues"
- "ColumnCorrelation"

RULE_PARAMETERS —一部のルールでは、実行に追加のパラメータが必要です。必要なパラメータについては、該当するルールのドキュメントを参照してください。

例: の異常 RowCount

例えば、 RowCount 異常を検出する場合は、ルール名 RowCount として を指定します。

```
DetectAnomalies "RowCount"
```

例: の異常 ColumnLength

例えば、 ColumnLength 異常を検出する場合は、ルール名と列名 ColumnLength として を指定します。

```
DetectAnomalies "ColumnLength" "id"
```

API を使用したデータ品質の測定および管理

このトピックでは、API を使用してデータ品質を測定および管理する方法について説明します。

目次

- [前提条件](#)
- [AWS Glue Data Quality 推奨事項の操作](#)
- [AWS Glue Data Quality ルールセットの操作](#)
- [AWS Glue Data Quality 実行の操作](#)
- [AWS Glue Data Quality 評価結果の操作](#)

前提条件

- お使いの boto3 のバージョンが最新で、最新の AWS Glue Data Quality API が含まれていることを確認します。
- お使いの AWS CLI のバージョンが最新で、最新の CLI が含まれていることを確認します。

AWS Glue ジョブを使用してこれらの API を実行している場合は、次の方法で boto3 ライブラリを最新バージョンに更新できます。

```
-additional-python-modules boto3==<version>
```

AWS Glue Data Quality 推奨事項の操作

AWS Glue Data Quality 推奨事項の実行を開始するには

```
class GlueWrapper:
    """Encapsulates AWS Glue actions."""
    def __init__(self, glue_client):
        """
        :param glue_client: A Boto3 Glue client.
        """
        self.glue_client = glue_client

    def start_data_quality_rule_recommendation_run(self, database_name, table_name,
role_arn):
        """
        Starts a recommendation run that is used to generate rules when you don't
know what rules to write. AWS Glue Data Quality analyzes the data and comes up with
recommendations for a potential ruleset. You can then triage the ruleset and modify
the generated ruleset to your liking.

        :param database_name: The name of the AWS Glue database which contains the
dataset.
        :param table_name: The name of the AWS Glue table against which we want a
recommendation
        :param role_arn: The Amazon Resource Name (ARN) of an AWS Identity and Access
Management (IAM) role that grants permission to let AWS Glue access the resources it
needs.

        """
        try:
```

```

        response = self.client.start_data_quality_rule_recommendation_run(
            DataSource={
                'GlueTable': {
                    'DatabaseName': database_name,
                    'TableName': table_name
                }
            },
            Role=role_arn
        )
    except ClientError as err:
        logger.error(
            "Couldn't start data quality recommendation run %s. Here's why: %s:
%s", name,
            err.response['Error']['Code'], err.response['Error']['Message'])
        raise
    else:
        return response['RunId']

```

推奨事項の実行では、`pushDownPredicates` または `catalogPartitionPredicates` を使用することで、パフォーマンスを高め、カタログソースの特定のパーティションのみに推奨事項を実行できます。

```

client.start_data_quality_rule_recommendation_run(
    DataSource={
        'GlueTable': {
            'DatabaseName': database_name,
            'TableName': table_name,
            'AdditionalOptions': {
                'pushDownPredicate': "year=2022"
            }
        }
    },
    Role=role_arn,
    NumberOfWorkers=2,
    CreatedRulesetName='<rule_set_name>'
)

```

AWS Glue Data Quality 推奨事項実行の結果を取得するには

```

class GlueWrapper:
    """Encapsulates AWS Glue actions."""
    def __init__(self, glue_client):

```

```

        """
        :param glue_client: A Boto3 AWS Glue client.
        """
        self.glue_client = glue_client

    def get_data_quality_rule_recommendation_run(self, run_id):
        """
        Gets the specified recommendation run that was used to generate rules.

        :param run_id: The id of the data quality recommendation run

        """
        try:
            response =
self.client.get_data_quality_rule_recommendation_run(RunId=run_id)
        except ClientError as err:
            logger.error(
                "Couldn't get data quality recommendation run %. Here's why: %s: %s",
run_id,
                err.response['Error']['Code'], err.response['Error']['Message'])
            raise
        else:
            return response

```

上記の反応により得られたオブジェクトから、実行時に推奨された RuleSet を抽出し、以降のステップで使用できます。

```

print(response['RecommendedRuleset'])

Rules = [
    RowCount between 2000 and 8000,
    IsComplete "col1",
    IsComplete "col2",
    StandardDeviation "col3" between 58138330.8 and 64258155.09,
    ColumnValues "col4" between 1000042965 and 1214474826,
    IsComplete "col5"
]

```

推奨事項の実行 (フィルタリングと一覧表示が可能) の一覧を表示するには

```

response = client.list_data_quality_rule_recommendation_runs(
    Filter={
        'DataSource': {

```

```
        'GlueTable': {
            'DatabaseName': '<database_name>',
            'TableName': '<table_name>'
        }
    }
)
```

既存の AWS Glue Data Quality の推奨事項タスクをキャンセルするには

```
response = client.cancel_data_quality_rule_recommendation_run(
    RunId='dqrun-d4b6b01957fdd79e59866365bf9cb0e40fxxxxxxx'
)
```

AWS Glue Data Quality ルールセットの操作

AWS Glue Data Quality ルールセットを作成するには

```
response = client.create_data_quality_ruleset(
    Name='<ruleset_name>',
    Ruleset='Rules = [IsComplete "col1", IsPrimaryKey "col2", RowCount between 2000 and 8000]',
    TargetTable={
        'TableName': '<table_name>',
        'DatabaseName': '<database_name>'
    }
)
```

データ品質ルールセットを取得するには

```
response = client.get_data_quality_ruleset(
    Name='<ruleset_name>'
)
print(response)
```

この API を使用すると、ルールセットを抽出できます。

```
print(response['Ruleset'])
```

テーブルのデータ品質ルールセットをすべて一覧表示するには

```
response = client.list_data_quality_rulesets()
```

API 内のフィルター条件を使用すると、特定のデータベースまたはテーブルに添付されたすべてのルールセットをフィルタリングできます。

```
response = client.list_data_quality_rulesets(
    Filter={
        'TargetTable': {
            'TableName': '<table_name>',
            'DatabaseName': '<database_name>'
        }
    },
)
```

データ品質ルールセットを更新するには

```
class GlueWrapper:
    """Encapsulates AWS Glue actions."""
    def __init__(self, glue_client):
        """
        :param glue_client: A Boto3 AWS Glue client.
        """
        self.glue_client = glue_client

    def update_data_quality_ruleset(self, ruleset_name, ruleset_string):
        """
        Update an AWS Glue Data Quality Ruleset

        :param ruleset_name: The name of the AWS Glue Data Quality ruleset to update
        :param ruleset_string: The DQDL ruleset string to update the ruleset with

        """
        try:
            response = self.client.update_data_quality_ruleset(
                Name=ruleset_name,
                Ruleset=ruleset_string
            )
        except ClientError as err:
            logger.error(
                "Couldn't update the AWS Glue Data Quality ruleset. Here's why: %s:
                %s",
                err.response['Error']['Code'], err.response['Error']['Message'])
```

```
        raise
    else:
        return response
```

データ品質ルールセットを削除するには

```
class GlueWrapper:
    """Encapsulates AWS Glue actions."""
    def __init__(self, glue_client):
        """
        :param glue_client: A Boto3 AWS Glue client.
        """
        self.glue_client = glue_client

    def delete_data_quality_ruleset(self, ruleset_name):
        """
        Delete a AWS Glue Data Quality Ruleset

        :param ruleset_name: The name of the AWS Glue Data Quality ruleset to delete

        """
        try:
            response = self.client.delete_data_quality_ruleset(
                Name=ruleset_name
            )
        except ClientError as err:
            logger.error(
                "Couldn't delete the AWS Glue Data Quality ruleset. Here's why: %s:
%s",
                err.response['Error']['Code'], err.response['Error']['Message'])
            raise
        else:
            return response
```

AWS Glue Data Quality 実行の操作

AWS Glue Data Quality 実行を開始するには

```
class GlueWrapper:
    """Encapsulates AWS Glue actions."""
    def __init__(self, glue_client):
        """
```

```

:param glue_client: A Boto3 AWS Glue client.
"""
self.glue_client = glue_client

def start_data_quality_ruleset_evaluation_run(self, database_name, table_name,
role_name, ruleset_list):
    """
    Start an AWS Glue Data Quality evaluation run

    :param database_name: The name of the AWS Glue database which contains the
dataset.
    :param table_name: The name of the AWS Glue table against which we want to
evaluate.
    :param role_arn: The Amazon Resource Name (ARN) of an AWS Identity and Access
Management (IAM) role that grants permission to let AWS Glue access the resources it
needs.
    :param ruleset_list: The list of AWS Glue Data Quality ruleset names to
evaluate.

    """
    try:
        response = client.start_data_quality_ruleset_evaluation_run(
            DataSource={
                'GlueTable': {
                    'DatabaseName': database_name,
                    'TableName': table_name
                }
            },
            Role=role_name,
            RulesetNames=ruleset_list
        )
    except ClientError as err:
        logger.error(
            "Couldn't start the AWS Glue Data Quality Run. Here's why: %s: %s",
            err.response['Error']['Code'], err.response['Error']['Message'])
        raise
    else:
        return response['RunId']

```

`pushDownPredicate` または `catalogPartitionPredicate` パラメータに合格すると、データ品質の実行の対象を、カタログテーブル内の特定のパーティションセットのみにできます。例:

```
response = client.start_data_quality_ruleset_evaluation_run(
```

```
DataSource={
  'GlueTable': {
    'DatabaseName': '<database_name>',
    'TableName': '<table_name>',
    'AdditionalOptions': {
      'pushDownPredicate': 'year=2023'
    }
  }
},
Role='<role_name>',
NumberOfWorkers=5,
Timeout=123,
AdditionalRunOptions={
  'CloudWatchMetricsEnabled': False
},
RulesetNames=[
  '<ruleset_name>',
]
)
```

AWS Glue Data Quality の実行に関する情報を取得するには

```
class GlueWrapper:
    """Encapsulates AWS Glue actions."""
    def __init__(self, glue_client):
        """
        :param glue_client: A Boto3 AWS Glue client.
        """
        self.glue_client = glue_client

    def get_data_quality_ruleset_evaluation_run(self, run_id):
        """
        Get details about an AWS Glue Data Quality Run

        :param run_id: The AWS Glue Data Quality run ID to look up

        """
        try:
            response = self.client.get_data_quality_ruleset_evaluation_run(
                RunId=run_id
            )
        except ClientError as err:
            logger.error(
```

```

        "Couldn't look up the AWS Glue Data Quality run ID. Here's why: %s:
%s",
        err.response['Error']['Code'], err.response['Error']['Message'])
    raise
else:
    return response

```

AWS Glue Data Quality 実行の結果を取得するには

所定の AWS Glue Data Quality 実行では、次の方法を使用して実行の評価の結果を抽出できます。

```

response = client.get_data_quality_ruleset_evaluation_run(
    RunId='d4b6b01957fdd79e59866365bf9cb0e40fxxxxxxx'
)

resultID = response['ResultIds'][0]

response = client.get_data_quality_result(
    ResultId=resultID
)

print(response['RuleResults'])

```

AWS Glue Data Quality 実行のすべてを一覧表示するには

```

class GlueWrapper:
    """Encapsulates AWS Glue actions."""
    def __init__(self, glue_client):
        """
        :param glue_client: A Boto3 AWS Glue client.
        """
        self.glue_client = glue_client

    def list_data_quality_ruleset_evaluation_runs(self, database_name, table_name):
        """
        Lists all the AWS Glue Data Quality runs against a given table

        :param database_name: The name of the database where the data quality runs
        :param table_name: The name of the table against which the data quality runs
        were created

        """
        try:

```

```
response = self.client.list_data_quality_ruleset_evaluation_runs(
    Filter={
        'DataSource': {
            'GlueTable': {
                'DatabaseName': database_name,
                'TableName': table_name
            }
        }
    }
)
except ClientError as err:
    logger.error(
        "Couldn't list the AWS Glue Quality runs. Here's why: %s: %s",
        err.response['Error']['Code'], err.response['Error']['Message'])
    raise
else:
    return response
```

フィルター句を変更すると、特定の時間内の結果のみ、あるいは特定のテーブルに対する実行のみを表示できます。

進行中の AWS Glue Data Quality 実行を停止するには

```
class GlueWrapper:
    """Encapsulates AWS Glue actions."""
    def __init__(self, glue_client):
        """
        :param glue_client: A Boto3 AWS Glue client.
        """
        self.glue_client = glue_client

    def cancel_data_quality_ruleset_evaluation_run(self, result_id):
        """
        Cancels a given AWS Glue Data Quality run

        :param result_id: The result id of a AWS Glue Data Quality run to cancel

        """
        try:
            response = self.client.cancel_data_quality_ruleset_evaluation_run(
                ResultId=result_id
            )
        except ClientError as err:
```

```
        logger.error(
            "Couldn't cancel the AWS Glue Data Quality run. Here's why: %s: %s",
            err.response['Error']['Code'], err.response['Error']['Message'])
        raise
    else:
        return response
```

AWS Glue Data Quality 評価結果の操作

AWS Glue Data Quality 実行の結果を取得するには

```
class GlueWrapper:
    """Encapsulates AWS Glue actions."""
    def __init__(self, glue_client):
        """
        :param glue_client: A Boto3 AWS Glue client.
        """
        self.glue_client = glue_client

    def get_data_quality_result(self, result_id):
        """
        Outputs the result of an AWS Glue Data Quality Result

        :param result_id: The result id of an AWS Glue Data Quality run

        """
        try:
            response = self.client.get_data_quality_result(
                ResultId=result_id
            )
        except ClientError as err:
            logger.error(
                "Couldn't get the AWS Glue Data Quality result. Here's why: %s: %s",
                err.response['Error']['Code'], err.response['Error']['Message'])
            raise
        else:
            return response
```

既存の AWS Glue Data Quality の推奨事項タスクをキャンセルするには

AWS Glue Data Quality 実行の ID を指定すると、次のように、結果 ID を抽出して実際の結果を取得できます。

```
response = client.get_data_quality_ruleset_evaluation_run(
    RunId='dqr-un-abca77ee126abe1378c1da1ae0750xxxxxxxx'
)

resultID = response['ResultIds'][0]

response = client.get_data_quality_result(
    ResultId=resultID
)

print(resp['RuleResults'])
```

アラート、デプロイ、スケジュールの設定

このトピックでは、AWS Glue Data Quality のアラート、デプロイ、スケジューリングを設定する方法について説明します。

目次

- [Amazon EventBridge 統合でのアラートと通知の設定](#)
 - [イベントパターンの追加設定のオプション](#)
 - [通知をメールとしてフォーマットする](#)
- [CloudWatch 統合でアラートと通知を設定する](#)
- [データ品質評価の結果をクエリしてダッシュボードを作成する](#)
- [を使用したデータ品質ルールのデプロイ AWS CloudFormation](#)
- [データ品質ルールのスケジューリング](#)

Amazon EventBridge 統合でのアラートと通知の設定

AWS Glue Data Quality は、Data Quality ルールセットの評価実行の完了時に発行される EventBridge イベントの発行をサポートします。これにより、Data Quality のルールが失敗した場合のアラートを、簡単に設定できます。

以下は、データカタログの Data Quality ルールセットを評価する際の、イベントのサンプルです。この情報を使用して、Amazon で利用できるデータを確認できます EventBridge。追加の API 呼び出しを発行すると、さらに詳細を入手できます。例えば、特定の実行の詳細を取得するには、結果 ID を使用して `get_data_quality_result` API を呼び出します。

```
{
  "version": "0",
  "id": "abcdef00-1234-5678-9abc-def012345678",
  "detail-type": "Data Quality Evaluation Results Available",
  "source": "aws.glue-dataquality",
  "account": "123456789012",
  "time": "2017-09-07T18:57:21Z",
  "region": "us-west-2",
  "resources": [],
  "detail": {
    "context": {
      "contextType": "GLUE_DATA_CATALOG",
      "runId": "dqrun-12334567890",
      "databaseName": "db-123",
      "tableName": "table-123",
      "catalogId": "123456789012"
    },
    "resultID": "dqresult-12334567890",
    "rulesetNames": ["rulset1"],
    "state": "SUCCEEDED",
    "score": 1.00,
    "rulesSucceeded": 100,
    "rulesFailed": 0,
    "rulesSkipped": 0
  }
}
```

Glue ETL または AWS Glue Studio ノートブックでデータ品質ルールセットを評価するときに発行されるイベントの例を次に示します。

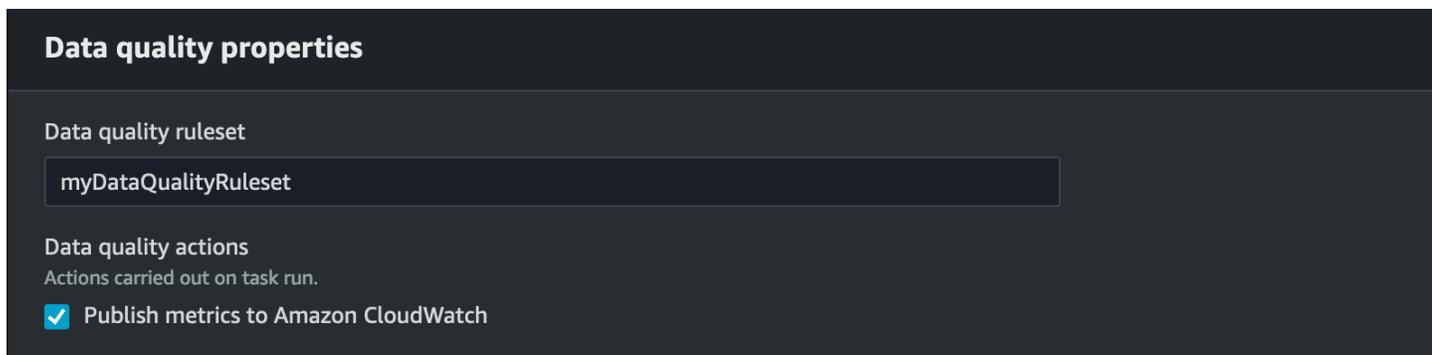
```
{
  "version": "0",
  "id": "abcdef00-1234-5678-9abc-def012345678",
  "detail-type": "Data Quality Evaluation Results Available",
  "source": "aws.glue-dataquality",
  "account": "123456789012",
  "time": "2017-09-07T18:57:21Z",
  "region": "us-west-2",
  "resources": [],
  "detail": {
    "context": {
      "contextType": "GLUE_JOB",
      "jobId": "jr-12334567890",

```

```
        "jobName": "dq-eval-job-1234",
        "evaluationContext": "",
      }
    "resultID": "dqresult-12334567890",
    "rulesetNames": ["rulset1"],
    "state": "SUCCEEDED",
    "score": 1.00
    "rulesSucceeded": 100,
    "rulesFailed": 0,
    "rulesSkipped": 0
  }
}
```

Data Quality 評価は、Data Catalog ジョブと ETL ジョブの両方で実行されるため、デフォルトで選択されている へのメトリクスの発行 Amazon CloudWatch オプションは、EventBridge 発行が機能するためには選択したままである必要があります。

EventBridge 通知の設定



The screenshot shows the 'Data quality properties' configuration page in the AWS Glue console. It includes a section for 'Data quality ruleset' with a text input field containing 'myDataQualityRuleset'. Below that is a section for 'Data quality actions' with the subtext 'Actions carried out on task run.' and a checked checkbox for 'Publish metrics to Amazon CloudWatch'.

出力されたイベントを受信してターゲットを定義するには、Amazon EventBridge ルールを設定する必要があります。ルールを作成するには

1. Amazon EventBridge コンソールを開きます。
2. ナビゲーションバーの [バス] セクションで [ルール] を選択します。
3. [Create Rule] (ルールの作成) を選択します。
4. [ルールの詳細を定義] では以下のとおり入力および選択します。
 - a. [名前] に myDQRule と入力します。
 - b. 説明を記入します (任意)。
 - c. [イベントバス] では、お使いのイベントバスを選択します。お持ちでない場合は、デフォルトのままにしておきます。

- d. [ルールタイプ] では [イベントパターンを持つルール] を選択し、続いて [次へ] を選択します。
5. [イベントパターンを構築] では以下のとおり入力および選択します。
 - a. イベントソースの場合は、AWS イベントまたは EventBridge パートナーイベントを選択します。
 - b. [サンプルイベント] セクションは飛ばします。
 - c. [作成のメソッド] では [パターンフォームを使用する] を選択します。
 - d. イベントパターンの場合
 - i. イベントソースに [AWS のサービス] を選択します。
 - ii. AWS サービスの Glue Data Quality を選択します。
 - iii. イベントタイプに [Data Quality Evaluation Results Available] を選択します。
 - iv. 具体的な状態に [FAILED] を選択します。そうすると、次のようなイベントパターンが表示されます。

```
{
  "source": ["aws.glue-dataquality"],
  "detail-type": ["Data Quality Evaluation Results Available"],
  "detail": {
    "state": ["FAILED"]
  }
}
```

- v. 設定オプションの詳細については、「[イベントパターンの追加設定のオプション](#)」を参照してください。
6. [ターゲットを選択] では以下のとおり入力および選択します。
 - a. [ターゲットタイプ] で [AWS のサービス] を選択します。
 - b. ターゲットの選択ドロップダウンを使用して、接続する AWS サービス (SNS、Lambda、SQS など) を選択し、次へ を選択します。
 7. [Configure tag(s)] で、[Add new tags] をクリックし、オプションのタグを追加して、[次へ] を選択します。
 8. すべての選択を示す概要ページが表示されます。一番下で [ルールの作成] を選択します。

イベントパターンの追加設定のオプション

イベントを、成功または失敗に応じてフィルタリングするだけでなく、さまざまなパラメータに応じてさらにフィルタリングするとよい場合があります。

これを行うには、[イベントパターン]のセクションに進み、[パターンを編集]を選択して追加のパラメータを指定します。イベントパターンのフィールドは大文字と小文字を区別する必要があります。以下は、イベントパターンの設定の例です。

特定のルールセットを評価する特定のテーブルからイベントを取り込むときは、次のタイプのパターンを使用します。

```
{
  "source": ["aws.glue-dataquality"],
  "detail-type": ["Data Quality Evaluation Results Available"],
  "detail": {
    "context": {
      "contextType": ["GLUE_DATA_CATALOG"],
      "databaseName": "db-123",
      "tableName": "table-123",
    },
    "rulesetNames": ["ruleset1", "ruleset2"]
  }
  "state": ["FAILED"]
}
```

ETL 体験の特定のジョブからイベントを取り込むときは、次のタイプのパターンを使用します。

```
{
  "source": ["aws.glue-dataquality"],
  "detail-type": ["Data Quality Evaluation Results Available"],
  "detail": {
    "context": {
      "contextType": ["GLUE_JOB"],
      "jobName": ["dq_evaluation_job1", "dq_evaluation_job2"]
    },
    "state": ["FAILED"]
  }
}
```

スコアが特定のしきい値 (70% など) を下回るイベントを取り込むときは、次のタイプのパターンを使用します。

```
{
  "source": ["aws.glue-dataquality"],
  "detail-type": ["Data Quality Evaluation Results Available"],
  "detail": {
```

```

"score": [{
  "numeric": ["<=", 0.7]
}]
}
}

```

通知をメールとしてフォーマットする

自社のビジネスチームに、適切にフォーマットしたメール通知を送信することが必要になる場合があります。Amazon EventBridge と AWS Lambda を使用してこれを実現できます。

Glue Data Quality rulesets **Glue_DQ_RULESET_CUSTOM_20de29c13537** run details



AWS Notifications <no-reply@sns.amazonaws.com>

Thursday, 11. May 2023 at 15:01

To: [REDACTED]

Glue Data Quality run details:

```

ruleset_name:  Glue_DQ_RULESET_CUSTOM_20de29c13537
glue_table_name:  devprod_tbl_nxc_taxi_data
glue_database_name:  devprod_db_nyc_taxi_data
run_id:  dqrun-066b41002a56921f9163a4e9156a4f6e20ce47a8
result_id:  dqresult-cd03a2e91c9114b611f6f79363b2288133fc96c0
state:  FAILED
score:  0.5
numRulesSucceeded:  1
numRulesFailed:  1
numRulesSkipped:  0

```

The subject of the email contains the name of the ruleset

Body of email with statistics from the Glue Data Quality Ruleset.

Here are the results of the ruleset evaluation steps

ruleset details evaluation steps results:

Name: Rule_1	Result: PASS	Description: IsComplete "vendorid"	
Name: Rule_2	Result: FAIL	EvaluationMessage: Value: 0.0 does not meet the constraint requirement!	Description: IsPrimaryKey "vendorid"

If you wish to stop receiving notifications from this topic, please click or visit the link below to unsubscribe:

[REDACTED] > [https://sns.us-east-1.amazonaws.com/unsubscribe.html?SubscriptionArn=arn:aws:sns:us-east-1:728060703200:SNSSStandardGlueDataQualityBlogAlertNotification:9d82097d-06f6-4c11-951a-3c0e1d9748f2&Endpoint=\[REDACTED\]](https://sns.us-east-1.amazonaws.com/unsubscribe.html?SubscriptionArn=arn:aws:sns:us-east-1:728060703200:SNSSStandardGlueDataQualityBlogAlertNotification:9d82097d-06f6-4c11-951a-3c0e1d9748f2&Endpoint=[REDACTED])

Please do not reply directly to this email. If you have any questions or comments regarding this email, please contact us at <https://aws.amazon.com/support>

次のサンプルコードを使用すると、データ品質通知をフォーマットしてメールを生成できます。

```

import boto3
import json
from datetime import datetime

```

```
sns_client = boto3.client('sns')
glue_client = boto3.client('glue')

sns_topic_arn = 'arn:aws:sns:<region-code>:<account-id>:<sns-topic-name>'

def lambda_handler(event, context):
    log_metadata = {}
    message_text = ""
    subject_text = ""

    if event['detail']['context']['contextType'] == 'GLUE_DATA_CATALOG':
        log_metadata['ruleset_name'] = str(event['detail']['rulesetNames'][0])
        log_metadata['tableName'] = str(event['detail']['context']['tableName'])
        log_metadata['databaseName'] = str(event['detail']['context']['databaseName'])
        log_metadata['runId'] = str(event['detail']['context']['runId'])
        log_metadata['resultId'] = str(event['detail']['resultId'])
        log_metadata['state'] = str(event['detail']['state'])
        log_metadata['score'] = str(event['detail']['score'])
        log_metadata['numRulesSucceeded'] = str(event['detail']['numRulesSucceeded'])
        log_metadata['numRulesFailed'] = str(event['detail']['numRulesFailed'])
        log_metadata['numRulesSkipped'] = str(event['detail']['numRulesSkipped'])

        message_text += "Glue Data Quality run details:\n"
        message_text += "ruleset_name: {}\n".format(log_metadata['ruleset_name'])
        message_text += "glue_table_name: {}\n".format(log_metadata['tableName'])
        message_text += "glue_database_name: {}\n".format(log_metadata['databaseName'])
        message_text += "run_id: {}\n".format(log_metadata['runId'])
        message_text += "result_id: {}\n".format(log_metadata['resultId'])
        message_text += "state: {}\n".format(log_metadata['state'])
        message_text += "score: {}\n".format(log_metadata['score'])
        message_text += "numRulesSucceeded:
{}\n".format(log_metadata['numRulesSucceeded'])
        message_text += "numRulesFailed: {}\n".format(log_metadata['numRulesFailed'])
        message_text += "numRulesSkipped: {}\n".format(log_metadata['numRulesSkipped'])

        subject_text = "Glue Data Quality ruleset {} run
details".format(log_metadata['ruleset_name'])

    else:
        log_metadata['ruleset_name'] = str(event['detail']['rulesetNames'][0])
        log_metadata['jobName'] = str(event['detail']['context']['jobName'])
```

```

log_metadata['jobId'] = str(event['detail']['context']['jobId'])
log_metadata['resultId'] = str(event['detail']['resultId'])
log_metadata['state'] = str(event['detail']['state'])
log_metadata['score'] = str(event['detail']['score'])

log_metadata['numRulesSucceeded'] = str(event['detail']['numRulesSucceeded'])
log_metadata['numRulesFailed'] = str(event['detail']['numRulesFailed'])
log_metadata['numRulesSkipped'] = str(event['detail']['numRulesSkipped'])

message_text += "Glue Data Quality run details:\n"
message_text += "ruleset_name: {}\n".format(log_metadata['ruleset_name'])
message_text += "glue_job_name: {}\n".format(log_metadata['jobName'])
message_text += "job_id: {}\n".format(log_metadata['jobId'])
message_text += "result_id: {}\n".format(log_metadata['resultId'])
message_text += "state: {}\n".format(log_metadata['state'])
message_text += "score: {}\n".format(log_metadata['score'])
message_text += "numRulesSucceeded:
{}\n".format(log_metadata['numRulesSucceeded'])
message_text += "numRulesFailed: {}\n".format(log_metadata['numRulesFailed'])
message_text += "numRulesSkipped: {}\n".format(log_metadata['numRulesSkipped'])

subject_text = "Glue Data Quality ruleset {} run
details".format(log_metadata['ruleset_name'])

resultID = str(event['detail']['resultId'])
response = glue_client.get_data_quality_result(ResultId=resultID)
RuleResults = response['RuleResults']
message_text += "\n\nruleset details evaluation steps results:\n\n"
subresult_info = []

for dic in RuleResults:
    subresult = "Name: {}\t\tResult: {}\t\tDescription: \t{}".format(dic['Name'],
dic['Result'], dic['Description'])
    if 'EvaluationMessage' in dic:
        subresult += "\t\tEvaluationMessage: {}".format(dic['EvaluationMessage'])
    subresult_info.append({
        'Name': dic['Name'],
        'Result': dic['Result'],
        'Description': dic['Description'],
        'EvaluationMessage': dic.get('EvaluationMessage', '')
    })
    message_text += "\n" + subresult

log_metadata['resultrun'] = subresult_info

```

```
sns_client.publish(  
    TopicArn=sns_topic_arn,  
    Message=message_text,  
    Subject=subject_text  
)  
  
return {  
    'statusCode': 200,  
    'body': json.dumps('Message published to SNS topic')  
}
```

CloudWatch 統合でアラートと通知を設定する

Amazon では顧客にアラートを送信するために 1 回限りの設定 EventBridge が必要なため EventBridge、Amazon を使用してデータ品質アラートを設定することをお勧めします。ただし、使い慣れ CloudWatch しているため、一部のお客様は Amazon を好みます。このようなお客様には、Amazon との統合を提供しています CloudWatch。

各 AWS Glue Data Quality 評価は、データ品質の実行ごとに `glue.data.quality.rules.passed` (合格したルールの数を示す) と `glue.data.quality.rules.failed` (失敗したルールの数を示す) という名前のメトリクスのペアを出力します。所定のデータ品質の実行がしきい値を下回った場合、発行されたこのメトリクスを使用してアラームを作成し、ユーザーに警告できます。Amazon SNS 通知からメールを送信するアラームをセットアップするには、以下の手順に従います。

Amazon SNS 通知からメールを送信するアラームをセットアップするには、以下の手順に従います。

1. Amazon CloudWatch コンソールを開きます。
2. [メトリクス] で、[すべてのメトリクス] を選択します。[カスタム名前空間] に、Glue Data Quality というタイトルの追加の名前空間が表示されています。

 Note

AWS Glue Data Quality の実行を開始するときは、Amazon へのメトリクスの発行 CloudWatchチェックボックスが有効になっていることを確認します。そうしないと、その特定の実行のメトリクスは Amazon に公開されません CloudWatch。

Glue Data Quality 名前空間で、テーブルごと、ルールセットごとに発行されるメトリクスを確認できます。このトピックでは `glue.data.quality.rules.failed` ルールを使用し、この値が 1 を超えた場合にアラームを発行する (つまり、不合格と評価されたルール数が 1 を超えた場合に通知する) ようにします。

3. アラームを作成するには、[アラーム] で、[すべてのアラーム] を選択します。
4. [アラームを作成] を選択します。
5. メトリクスの選択 を選択します。
6. 作成したテーブルに対応する `glue.data.quality.rules.failed` メトリクスを選択し、次に [メトリクスの選択] を選択します。
7. [メトリクスと条件の指定] タブの [メトリクス] セクションで、以下のとおり選択します。
 - a. [統計] で、[合計] を選択します。
 - b. [期間] は、[1 minute] を選択します。
8. [条件] セクションで、以下のとおり選択および入力します。
 - a. [Threshold type] で [静的] を選択します。
 - b. [Whenever glue.data.quality.rules.failed is...] では、[以上] を選択します。
 - c. [than...] に、しきい値として 1 と入力します。

上記を選択したことにより、`glue.data.quality.rules.failed` メトリクスで 1 以上の値が発行されると、アラームがトリガーされることとなります。ただし、データがない場合は許容範囲として処理されます。

9. [次へ] をクリックします。
10. [アクションの設定] で、以下のとおり選択および入力します。
 - a. [アラーム状態トリガー] セクションで、[アラーム状態] を選択します。
 - b. [Send a notification to the following SNS topic] セクションで、[Create a new topic to send a notification via a new SNS topic] を選択します。

- c. [Email endpoints that will receive the notification] に、自身のメールアドレスを入力します。[トピックの作成] をクリックします。
- d. [次へ] をクリックします。

11[アラーム名] に myFirstDQAlarm と入力し、[次へ] を選択します。

12.すべての選択を示す概要ページが表示されます。一番下で [アラームを作成] を選択します。

これで、Amazon アラーム CloudWatch ダッシュボードからアラームが作成されていることを確認できます。

データ品質評価の結果をクエリしてダッシュボードを作成する

データ品質評価の結果を表示する、ダッシュボードの作成が必要になる場合があります。これには、以下の2つの方法があります。

Amazon S3 にデータを書き込むには、次のコード EventBridge で Amazon を設定します。Amazon S3

```
import boto3
import json
from datetime import datetime

s3_client = boto3.client('s3')
glue_client = boto3.client('glue')

s3_bucket = 's3-bucket-name'

def write_logs(log_metadata):
    try:
        filename = datetime.now().strftime("%m%d%Y%H%M%S") + ".json"
        key_opts = {
            'year': datetime.now().year,
            'month': "{:02d}".format(datetime.now().month),
            'day': "{:02d}".format(datetime.now().day),
            'filename': filename
        }
        s3key = "gluedataqualitylogs/year={year}/month={month}/day={day}/"
        {filename}").format(**key_opts)
```

```
s3_client.put_object(Bucket=s3_bucket, Key=s3key,
Body=json.dumps(log_metadata))
except Exception as e:
    print(f'Error writing logs to S3: {e}')

def lambda_handler(event, context):
    log_metadata = {}
    message_text = ""
    subject_text = ""

    if event['detail']['context']['contextType'] == 'GLUE_DATA_CATALOG':
        log_metadata['ruleset_name'] = str(event['detail']['rulesetNames'][0])
        log_metadata['tableName'] = str(event['detail']['context']['tableName'])
        log_metadata['databaseName'] = str(event['detail']['context']['databaseName'])
        log_metadata['runId'] = str(event['detail']['context']['runId'])
        log_metadata['resultId'] = str(event['detail']['resultId'])
        log_metadata['state'] = str(event['detail']['state'])
        log_metadata['score'] = str(event['detail']['score'])
        log_metadata['numRulesSucceeded'] = str(event['detail']['numRulesSucceeded'])
        log_metadata['numRulesFailed'] = str(event['detail']['numRulesFailed'])
        log_metadata['numRulesSkipped'] = str(event['detail']['numRulesSkipped'])

        message_text += "Glue Data Quality run details:\n"
        message_text += "ruleset_name: {}\n".format(log_metadata['ruleset_name'])
        message_text += "glue_table_name: {}\n".format(log_metadata['tableName'])
        message_text += "glue_database_name: {}\n".format(log_metadata['databaseName'])
        message_text += "run_id: {}\n".format(log_metadata['runId'])
        message_text += "result_id: {}\n".format(log_metadata['resultId'])
        message_text += "state: {}\n".format(log_metadata['state'])
        message_text += "score: {}\n".format(log_metadata['score'])
        message_text += "numRulesSucceeded:
{}\n".format(log_metadata['numRulesSucceeded'])
        message_text += "numRulesFailed: {}\n".format(log_metadata['numRulesFailed'])
        message_text += "numRulesSkipped: {}\n".format(log_metadata['numRulesSkipped'])

        subject_text = "Glue Data Quality ruleset {} run
details".format(log_metadata['ruleset_name'])

    else:
        log_metadata['ruleset_name'] = str(event['detail']['rulesetNames'][0])
        log_metadata['jobName'] = str(event['detail']['context']['jobName'])
        log_metadata['jobId'] = str(event['detail']['context']['jobId'])
        log_metadata['resultId'] = str(event['detail']['resultId'])
```

```
log_metadata['state'] = str(event['detail']['state'])
log_metadata['score'] = str(event['detail']['score'])

log_metadata['numRulesSucceeded'] = str(event['detail']['numRulesSucceeded'])
log_metadata['numRulesFailed'] = str(event['detail']['numRulesFailed'])
log_metadata['numRulesSkipped'] = str(event['detail']['numRulesSkipped'])

message_text += "Glue Data Quality run details:\n"
message_text += "ruleset_name: {}\n".format(log_metadata['ruleset_name'])
message_text += "glue_job_name: {}\n".format(log_metadata['jobName'])
message_text += "job_id: {}\n".format(log_metadata['jobId'])
message_text += "result_id: {}\n".format(log_metadata['resultId'])
message_text += "state: {}\n".format(log_metadata['state'])
message_text += "score: {}\n".format(log_metadata['score'])
message_text += "numRulesSucceeded:
{}\n".format(log_metadata['numRulesSucceeded'])
message_text += "numRulesFailed: {}\n".format(log_metadata['numRulesFailed'])
message_text += "numRulesSkipped: {}\n".format(log_metadata['numRulesSkipped'])

subject_text = "Glue Data Quality ruleset {} run
details".format(log_metadata['ruleset_name'])

resultID = str(event['detail']['resultId'])
response = glue_client.get_data_quality_result(ResultId=resultID)
RuleResults = response['RuleResults']
message_text += "\n\nruleset details evaluation steps results:\n\n"
subresult_info = []

for dic in RuleResults:
    subresult = "Name: {}\t\tResult: {}\t\tDescription: \t{}".format(dic['Name'],
dic['Result'], dic['Description'])
    if 'EvaluationMessage' in dic:
        subresult += "\t\tEvaluationMessage: {}".format(dic['EvaluationMessage'])
    subresult_info.append({
        'Name': dic['Name'],
        'Result': dic['Result'],
        'Description': dic['Description'],
        'EvaluationMessage': dic.get('EvaluationMessage', '')
    })
    message_text += "\n" + subresult

log_metadata['resultrun'] = subresult_info

write_logs(log_metadata)
```

```
return {
    'statusCode': 200,
    'body': json.dumps('Message published to SNS topic')
}
```

Amazon S3 に書き込んだ後、AWS Glue クローラを使用して Athena に登録し、テーブルをクエリできます。

データ品質評価中に Amazon S3 のロケーションを設定する

AWS Glue Data Catalog または AWS Glue ETL でデータ品質タスクを実行する場合、Amazon S3 の場所を指定して、データ品質結果を Amazon S3 に書き込むことができます。以下の構文を使用すると、ターゲットを参照してテーブルを作成し、データ品質評価の結果を読み取れます。

CREATE EXTERNAL TABLE と MSCK REPAIR TABLE のクエリは個別に実行する必要がありますので注意してください。

```
CREATE EXTERNAL TABLE <my_table_name>(
    catalogid string,
    databasename string,
    tablename string,
    dqrunid string,
    evaluationstartedon timestamp,
    evaluationcompletedon timestamp,
    rule string,
    outcome string,
    failurereason string,
    evaluatedmetrics string)
PARTITIONED BY (
    `year` string,
    `month` string,
    `day` string)
ROW FORMAT SERDE 'org.openx.data.jsonserde.JsonSerDe'
WITH SERDEPROPERTIES (

    'paths'='catalogId,databaseName,dqRunId,evaluatedMetrics,evaluationCompletedOn,evaluationStart
STORED AS INPUTFORMAT 'org.apache.hadoop.mapred.TextInputFormat'
OUTPUTFORMAT 'org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat'
LOCATION 's3://glue-s3-dq-bucket-us-east-2-results/'
TBLPROPERTIES (
```

```
'classification'='json',  
'compressionType'='none',  
'typeOfData'='file');
```

```
MSCK REPAIR TABLE <my_table_name>;
```

上記のテーブルを作成すれば、Amazon Athena を使用して分析クエリを実行できるようになります。

を使用したデータ品質ルールのデプロイ AWS CloudFormation

を使用して AWS CloudFormation、データ品質ルールを作成できます。詳細については、「[for AWS CloudFormation AWS Glue](#)」を参照してください。

データ品質ルールのスケジューリング

データ品質ルールのスケジューリングには、次の方法を使用します。

- データカタログからデータ品質ルールをスケジュールする: このオプションを使用してデータ品質スキャンを簡単にスケジュールできるコードユーザーはいません。AWS Glue Data Quality は Amazon でスケジュールを作成します EventBridge。データ品質ルールをスケジュールするには、以下の手順に従います。
 - ルールセットに進み、[実行] をクリックします。
 - [Run frequency] で希望するスケジュールを選択し、タスク名を入力します。このタスク名は、スケジュールの名前です EventBridge。
- Amazon EventBridge と AWS Step Functions を使用して、データ品質ルールの評価とレコメンデーションを調整します。

AWS Glue Data Quality エラーのトラブルシューティング

AWS Glue Data Quality でエラーが発生した場合は、次の解決策を使用して問題の原因を見つけて修正してください。

目次

- [エラー: AWS Glue Data Quality モジュールがありません](#)
- [エラー: AWS Lake Formation のアクセス許可が不十分です](#)

- [エラー: ルールセットに一意の名前が付いていない](#)
- [エラー: テーブルに特殊文字が含まれている](#)
- [エラー: 大規模なルールセットによるオーバーフローエラー](#)
- [エラー: ルール全体のステータスが失敗である](#)
- [AnalysisException: デフォルトデータベースの存在を確認できません](#)
- [エラーメッセージ: Provided key map not suitable for given data frames](#)
- [ユーザークラス の例外: java.lang.RuntimeException : データの取得に失敗しました。ログインを確認して詳細 CloudWatch を確認する](#)
- [LAUNCH ERROR: Error downloading from S3 for bucket](#)
- [InvalidInputException \(ステータス: 400\): DataQuality ルールを解析できません](#)
- [エラー: EventBridge が、設定したスケジュールに基づいて Glue DQ ジョブをトリガーしない](#)
- [CustomSQL エラー](#)
- [動的ルール](#)
- [ユーザークラス: org.apache.spark.sql.AnalysisException: org.apache.hadoop.hive.ql.metadata の例外。HiveException](#)
- [UNCLASSIFIED_ERROR; IllegalArgumentException: 解析エラー: ルールまたはアナライザーが指定されていません。、入力時に実行可能な代替手段がありません](#)

エラー: AWS Glue Data Quality モジュールがありません

エラーメッセージ: No module named 'awsgluedq'.

解決策: このエラーは、サポートされていないバージョンで AWS Glue Data Quality を実行したときに発生します。AWS Glue Data Quality は Glue バージョン 3.0 以降でのみサポートされています。

エラー: AWS Lake Formation のアクセス許可が不十分です

エラーメッセージ: ユーザークラスの例外:

com.amazonaws.services.glue.model.AccessDeniedException: impact_sdg_involvement (サービス: 、ステータスコード: 400AWS Glue、エラーコード: 、AccessDeniedExceptionリクエスト ID: 465ae693-b7ba-4df0-a4e4-6b17xxxx、プロキシ: null) に対する Lake Formation 許可が不十分です。

解決策: AWS Lake Formation で十分なアクセス許可を付与する必要があります。

エラー: ルールセットに一意の名前が付いていない

エラーメッセージ: ユーザークラスの例外: ...services.glue.model.AlreadyExistsException: 同じ名前の別のルールセットが既に存在します。

解決方法: ルールセットはグローバルであり、一意である必要があります。

エラー: テーブルに特殊文字が含まれている

エラーメッセージ: ユーザークラス: org.apache.spark.sql.AnalysisException: は、指定された入力列 [primary.data_end_time, primary.data_start_time, primary.end_time, primary.last_updated, primary.message, primary.process_date, primary.rowhash, primary.run_by, primary.run_id, primary.start_time, primary.status] の「C」を解決できません。

解決策: 現在、「」などの特殊文字を含むテーブルでは AWS Glue Data Quality を実行できないという制限があります。

エラー: 大規模なルールセットによるオーバーフローエラー

エラーメッセージ: ユーザークラスの例外: java.lang.StackOverflowError。

解決方法: 2,000 以上のルールを含む大規模なルールセットを使用すると、このエラーが発生する可能性があります。ルールを複数のルールセットに分割します。

エラー: ルール全体のステータスが失敗である

エラーの状態: My Ruleset is successful, but my overall rule status is failed.

解決策: このエラーは、発行 CloudWatch 中にメトリクスを Amazon に発行するオプションを選択したために発生した可能性があります。データセットが VPC にある場合、VPC は AWS Glue が Amazon にメトリクスを発行することを許可しない場合があります CloudWatch。この場合、>VPC が Amazon にアクセスするためのエンドポイントを設定する必要があります CloudWatch。

AnalysisException: デフォルトデータベースの存在を確認できません

エラー条件: AnalysisException: デフォルトデータベースの存在を確認できません:
com.amazonaws.services.glue.model.AccessDeniedException: デフォルトに対する Lake Formation のアクセス許可が不十分です (サービス: AWS Glue、ステータスコード: 400、エラーコード:

AccessDeniedException、リクエスト ID: XXXXXXXX-XXXX-XXXX-XXXXXXXXXXXX、プロキシ:
null)

解決方法: AWS Glue ジョブのカタログ統合では、AWS Glue は AWS Glue GetDatabase API を使用してデフォルトのデータベースが存在するかを常にチェックします。DESCRIBE Lake Formation アクセス許可が付与されないとき、または GetDatabase IAM アクセス許可が付与されたとき、デフォルトのデータベースの存在を確認するときのジョブは失敗します。

解決するには:

1. Lake Formation にデフォルトのデータベース用の DESCRIBE アクセス許可を追加します。
2. AWS Glue ジョブにアタッチされている IAM ロール を Lake Formation でデータベース作成者として設定します。これにより、デフォルトのデータベースが自動的に作成され、必要な Lake Formation へのアクセス許可がロールに付与されます。
3. `--enable-data-catalog` オプションを無効にします (AWS Glue Studio では [Use Data Catalog as the Hive metastore] と表示されます)。

ジョブに Spark SQL Data Catalog 統合を必要としない場合は、無効にできます。

エラーメッセージ: Provided key map not suitable for given data frames

エラー状態: Provided key map not suitable for given data frames.

解決策: DataSetMatchルールタイプを使用していて、結合キーが重複しています。結合キーは一意でなければならず、また NULL であってはなりません。一意の結合キーが使用できない場合は、などの他のルールタイプを使用して概要データにAggregateMatch一致させることを検討してください。

ユーザークラス の例外: java.lang.RuntimeException : データの取得に失敗しました。ログインを確認して詳細 CloudWatch を確認する

エラー条件: ユーザークラス の例外: java.lang.RuntimeException : データの取得に失敗しました。詳細については、ログイン CloudWatch を確認してください。

解決策: これは、Amazon RDS または と比較する Amazon S3 ベースのテーブルに DQ ルールを作成するときに発生します Amazon Redshift。この場合、AWS Glue は接続をロードできません。代わりに、Amazon Redshift または Amazon RDS データセットに DQ ルールを設定してみてください。これは既知のバグです。

LAUNCH ERROR: Error downloading from S3 for bucket

エラー状態: LAUNCH ERROR: Error downloading from S3 for bucket: aws-glue-ml-data-quality-assets-us-east-1, key: jars/aws-glue-ml-data-quality-etl.jar.Access Denied (Service: Amazon S3; Status Code: 403; Please refer logs for details) .

解決策 : AWS Glue Data Quality に渡されるロールのアクセス許可は、前の Amazon S3 の場所からの読み取りを許可する必要があります。この IAM ポリシーを、次のロールに添付します。

```
{
  "Sid": "allowS3",
  "Effect": "Allow",
  "Action": "s3:GetObject",
  "Resource": "arn:aws:s3:::aws-glue-ml-data-quality-assets-<region>/*"
}
```

アクセス許可の詳細については、「[Data Quality authorization](#)」を参照してください。これらのライブラリは、データセットのデータ品質を評価する際に必要になります。

InvalidInputException (ステータス: 400): DataQuality ルールを解析できません

エラー条件 : InvalidInputException (ステータス: 400): DataQuality ルールを解析できません。

解決方法: このエラーには、さまざまな原因が考えられます。その 1 つが、ルールに一重引用符が含まれていることです。二重引用符で囲まれていることを確認してください。例:

```
Rules = [
  ColumnValues "tipo_vinculo" in ["COD0", "DOC0", "COC0", "DOD0"] AND "categoria" = 'ES'
  AND "cod_bandera" = 'CEP'
```

上記を下記のように変更します。

```
Rules = [
  (ColumnValues "tipovinculo" in [ "COD0", "DOC0", "COC0", "DOD0"]) AND (ColumnValues
  "categoria" = "ES")
```

```
AND (ColumnValues "codbandera" = "CEP")
]
```

エラー: EventBridge が、設定したスケジュールに基づいて Glue DQ ジョブをトリガーしない

エラー状態: Eventbridge is not triggering AWS Glue Data Quality jobs based on the schedule I setup.

解決方法: ジョブをトリガーするロールに適切なアクセス許可が付与されていない可能性があります。ジョブのトリガーに使用するロールに、「[評価実行のスケジュールリングに必要な IAM の設定](#)」に記載されているアクセス許可が付与されていることを確認します。

CustomSQL エラー

エラー状態: The output from CustomSQL must contain at least one column that matches the input dataset for AWS Glue Data Quality to provide row level results. The SQL query is a valid query but no columns from the SQL result are present in the Input Dataset. Ensure that matching columns are returned from the SQL。

解決策: SQL クエリは有効ですが、主テーブルの列のみを選択していることを確認してください。主テーブルから列の sum、count などの集計関数を選択すると、このエラーが発生する可能性があります。

エラー状態: There was a problem when executing your SQL statement: cannot resolve "Col"。

解決策: この列は主テーブルには存在しません。

エラー状態: The columns that are returned from the SQL statement should only belong to the primary table. "In this case, some columns (Col) belong to reference table"。

解決策: SQL クエリで主テーブルを他の参照テーブルと結合する場合は、主テーブルの行レベルの結果を生成するために、SELECT ステートメントに主テーブルの列名のみが含まれていることを確認してください。

動的ルール

エラー状態: Dynamic rules require job context, and cannot be evaluated in interactive session or data preview..

原因: ルールセットに動的 DQ ルールが含まれている場合、このエラーメッセージがデータプレビュー結果や他のインタラクティブセッションに表示されることがあります。動的ルールは特定のジョブ名と評価コンテキストに関連する履歴メトリクスを参照するため、インタラクティブセッションでは評価できません。

解決策: AWS Glue ジョブを実行すると履歴メトリクスが生成され、同じジョブについて後でジョブを実行する際に参照できます。

エラー状態:

- [RuleType] rule only supports simple atomic operands in thresholds..
- Function last not yet implemented for [RuleType] rule.

解決策: 動的ルールは通常、数値式に含まれるすべての DQDL ルールタイプでサポートされています ([「DQDL リファレンス」](#)を参照)。ただし、複数のメトリクス ColumnValues および を生成する一部のルール ColumnLength はまだサポートされていません。

エラー状態: Binary expression operands must resolve to a single number..

原因: 動的ルールは $\text{RowCount} > \text{avg}(\text{last}(5)) * 0.9$ のようなバイナリ式をサポートしています。ここでのバイナリ式とは、 $\text{avg}(\text{last}(5)) * 0.9$ のことです。このルールは、オペランド $\text{avg}(\text{last}(5))$ と 0.9 の両方が単一の数値に解決されるため有効です。間違っただけの例としては、 $\text{RowCount} > \text{last}(5) * 0.9$ などがあります。last(5) は現在の行数と有意義な比較ができないリストが生成されるからです。

解決策: 集計関数を使用して、リスト値のオペランドを 1 つの数値に減らします。

エラー状態:

- Rule threshold results in list, and a single value is expected. Use aggregation functions to produce a single value. Valid example: `sum(last(10)), avg(last(10))`.
- Rule threshold results in empty list, and a single value is expected.

原因: 動的ルールを使用すると、データセットの一部の機能とその履歴値を比較できます。最後の関数では、正の整数の引数が指定されていれば、複数の履歴値を取得できます。例えば `last(5)` は、ルールのジョブ実行で確認された直近の 5 つの値を取得します。

解決策: 集計関数を使用してこれらの値を 1 つの数値に減らし、現在のジョブ実行で測定された値と有意義な比較を行う必要があります。

有効な例:

- `RowCount >= avg(last(5))`
- `RowCount > last(1)`
- `RowCount < last()`

無効な例: `RowCount > last(5)`。

エラー状態:

- `Function index used in threshold requires positive integer argument.`
- `Index argument must be an integer. Valid syntax example: RowCount > index(last(10), 2)), which means RowCount must be greater than third most recent execution from last 10 job runs.`

解決策: 動的ルールを作成する場合、`index` 集計機能を使用してリストから 1 つの履歴値を選択できます。例えば、`RowCount > index(last(5), 1)` は、現在のジョブで観測された行数が、そのジョブで観測された行数の中で 2 番目に新しい行数より真に大きいかどうかを確認します。`index` のインデックスは 0 です。

エラー状態: `IllegalArgumentException: Parsing Error: Rule Type: DetectAnomalies is not valid.`

解決策: 異常検出は AWS Glue 4.0 でのみ使用できます。

エラー状態: `IllegalArgumentException: Parsing Error: Unexpected condition for rule of type ... no viable alternative at input ...`。

注: ... は動的です。例えば、`IllegalArgumentException: Parsing Error: Unexpected condition for rule of type RowCount with number return type, line 4:19 no viable alternative at input '>last'` などです。

解決策: 異常検出は AWS Glue 4.0 でのみ使用できます。

ユーザークラス: org.apache.spark.sql.AnalysisException:
org.apache.hadoop.hive.ql.metadata の例外。HiveException

エラー状態: Exception in User Class: org.apache.spark.sql.AnalysisException:
org.apache.hadoop.hive.ql.metadata.HiveException: Unable to fetch table
mailpiece_submitted. StorageDescriptor#InputFormat cannot be null for
table: mailpiece_submitted (Service: null; Status Code: 0; Error Code:
null; Request ID: null; Proxy: null)

原因: AWS Glue Data Catalog で Apache Iceberg を使用していて、AWS Glue Data Catalog の入
力形式属性が空です。

解決策: この問題は、DQ ルールで CustomSQL ルールタイプを使用しているときに発生します。
この問題を解決する 1 つの方法は、「プライマリ」を使用するか、<database>.<table> in
Custom ruletype に glue_catalog. のカタログ名を追加することです。

UNCLASSIFIED_ERROR; IllegalArgumentException: 解析エラー: ルールま
たはアナライザーが指定されていません。、入力時に実行可能な代替手段
がありません

エラー状態: UNCLASSIFIED_ERROR; IllegalArgumentException: Parsing Error: No
rules or analyzers provided., no viable alternative at input

解決策: DQDL は解析できません。これが発生する可能性があるインスタンスがいくつかあります。
複合ルールを使用している場合は、それらに適切な括弧があることを確認してください。

```
(RowCount >= avg(last(10)) * 0.6) and (RowCount <= avg(last(10)) * 1.4) instead of  
RowCount >= avg(last(10)) * 0.6 and RowCount <= avg(last(10)) * 1.4
```

Amazon Q でのデータインテグレーション AWS Glue

Amazon Q AWS Glue データ統合は、データエンジニアと ETL 開発者が自然言語を使用してデータ統合ジョブを構築できるようにする、新しいジェネレーティブ AI 機能です。AWS Glue エンジニアと開発者は Amazon Q にジョブの作成、問題のトラブルシューティング、AWS Glue データ統合に関する質問への回答を依頼できます。

Amazon Q とは

Note

Amazon Bedrock を搭載: AWS [不正行為の自動検出機能を実装しています](#)。Amazon Q データ統合は Amazon Bedrock に構築されているため、ユーザーは Amazon Bedrock に実装されているコントロールを最大限に活用し、安全性、セキュリティ、人工知能 (AI) の責任ある使用を実施できます。

Amazon Q は生成型人工知能 (AI) を活用した会話型アシスタントで、アプリケーションの理解、構築、拡張、運用を支援します。AWS Amazon Q AWS を支えるこのモデルには高品質のコンテンツが追加され、より完全で実用的で参考になる回答が得られるため、構築を加速できます。AWS 詳細については、「[Amazon Q とは](#)」を参照してください。

AWS Glue の Amazon Q データ統合とは

の Amazon Q AWS Glue データ統合には以下の機能が含まれます。

- **チャット** — Amazon Q Data Integration in AWS Glue は、AWS Glue AWS Glue ソースコネクタとデステイネーションコネクタ、AWS Glue ETL ジョブ、データカタログ、クローラー、およびその他の機能ドキュメント、ベストプラクティスなどのデータ統合ドメインに関する自然言語の質問に英語で回答できます。AWS Lake Formation Amazon Q AWS Glue データインテグレーションでは、step-by-step 情報ソースへの参照を含む指示が送られてきます。
- **データ統合コード生成** — Amazon Q データ統合では、AWS Glue ETL AWS Glue スクリプトに関する質問に回答したり、英語で自然言語で質問された場合に新しいコードを生成したりできます。
- **トラブルシューティング** — Amazon Q AWS Glue のデータ統合は、AWS Glue ジョブのエラーを理解するのに役立つように設計されており、step-by-step 問題の根本原因と解決のための指示を提供します。

Note

の Amazon Q データ統合では、会話中の会話のコンテキストを使用して future AWS Glue 応答を通知することはありません。での Amazon Q AWS Glue データ統合に関する各会話は、以前または future 会話とは無関係です。

AWS Glue の Amazon Q データ統合の操作について

Amazon Q パネルでは、Amazon Q に AWS Glue ETL スクリプト用のコード生成をリクエストしたり、AWS Glue 機能に関する質問に答えたり、エラーのトラブルシューティングを行うことができます。レスポンスは ETL スクリプトで、スクリプトのカスタマイズ、確認、PySpark step-by-step 実行の手順が記載されています。質問に対して、データ統合ナレッジベースに基づいて、概要と参照用のソース URL を含む回答が生成されます。

たとえば、Amazon Q に「Snowflakeから読み取り、フィールドの名前を変更し、Redshift に書き込む Glue eスクリプトを提供してください」とリクエストすると、それに応じて、Amazon Q のデータ統合により、AWS Glue AWS Glue 要求されたアクションを実行できるジョブスクリプトが返されます。生成されたコードをレビューして、要求した意図を満たしていることを確認できます。問題がなければ、AWS Glue それをジョブとして本番環境にデプロイできます。統合に、エラーや障害の説明や解決策の提案を依頼し、解決策を提案することで、ジョブのトラブルシューティングを行うことができます。Amazon Q は、AWS Glue データ統合のベストプラクティスに関する質問にお答えします。

The screenshot shows the AWS Glue Studio interface. On the left is a navigation sidebar with categories like 'Getting started', 'ETL jobs', 'Data Catalog', and 'Legacy pages'. The main content area is titled 'AWS Glue Studio' and includes a 'Create job' section with three options: 'Visual ETL', 'Notebook', and 'Script editor'. Below this is an 'Example jobs' section with a 'Create example job' button. The 'Your jobs' section shows a search for 'demo' with 2 matches, listing jobs like 'q-demo-taxi' and 'q-demo' with their respective types and last modified dates.

以下は、Amazon Q AWS Glue でのデータ統合がどのように構築に役立つかを示す質問例です AWS Glue。

AWS Glue ETL コード生成:

- S3 から JSON を読み取り、適用マッピングを使用してフィールドを変換し、Amazon Redshift AWS Glue に書き込むスクリプトを記述します。
- DynamoDBから読み込み、DropNullFields 変換を適用し、AWS Glue ParquetとしてS3に書き込むためのスクリプトを作成する方法を教えてください。
- MySQLから読み取り、ビジネスロジックに基づいていくつかのフィールドを削除し、AWS Glue Snowflakeに書き込むスクリプトを教えてください
- DynamoDB から読み取り、S3 に JSON AWS Glue として書き込むジョブを記述する
- S3 AWS Glue AWS Glue へのデータカタログのスクリプトの開発を手伝ってください
- S3 から JSON を読み取り、ヌルをドロップして Redshift AWS Glue に書き込むジョブを作成する

AWS Glue 機能の説明:

- AWS Glue データクオリティを使用するにはどうすればいいですか？
- AWS Glue ジョブブックマークの使い方は？
- AWS Glue オートスケーリングを有効にするにはどうすればいいですか？

- AWS Glue ダイナミックフレームと Spark データフレームの違いは？
- AWS Glueではどのような種類の接続がサポートされていますか？

AWS Glue トラブルシューティング:

- AWS Glue ジョブのメモリ不足 (OOM) エラーのトラブルシューティング方法は？
- AWS Glue Data Quality の設定時に表示されるエラーメッセージにはどのようなものがありますか。また、その修正方法を教えてください。
- Amazon S3 AWS Glue アクセスが拒否されたというエラーが発生したジョブを修正する方法を教えてください。
- AWS Glue ジョブのデータシャッフルに関する問題を解決する方法を教えてください。

Amazon Q データ統合とのやりとりのベストプラクティス

Amazon Q データ統合を操作するためのベストプラクティスを次に示します。

- Amazon Q データインテグレーションを利用するときは、具体的な質問をし、複雑なリクエストがある場合は繰り返し、回答が正確かどうかを検証します。
- データ統合のプロンプトを自然言語で提供する場合は、必要なものをアシスタントが正確に理解できるように、できるだけ具体的にしてください。「S3 からデータを抽出」と尋ねる代わりに、「S3 から JSON AWS Glue ファイルを抽出するスクリプトを書く」などの詳細な情報を提供してください。
- 生成されたスクリプトを実行する前に確認して、正確であることを確認してください。生成されたスクリプトにエラーがあったり、意図したものと一致しない場合は、アシスタントに修正方法を伝えてください。
- 生成 AI は新しいテクノロジーであり、応答にはハルシネーションと呼ばれる誤りがある場合があります。現在の環境やワークロードで使用する前に、すべてのコードをテストしてエラーや脆弱性がないかを確認する必要があります。

AWS Glue サービス改善における Amazon Q データ統合

Amazon Q Data Integration in AWS Glue AWS がサービスに関する最も関連性の高い情報を提供できるように、Amazon Q からの特定のコンテンツ (Amazon Q に寄せられる質問やその回答など) をサービスの改善に使用場合があります。

当社が使用するコンテンツとオプトアウト方法については、[Amazon Q 開発者ユーザーガイドの「Amazon Q 開発者サービスの向上」](#)を参照してください。

考慮事項

AWS Glueの Amazon Q データ統合を使用する前に、以下の項目について検討してください。

- 現在、PySpark コード生成はカーネルでのみ機能します。生成されたコードは Python Spark AWS Glue に基づくジョブ用です。
- の Amazon Q データ統合でサポートされるコード生成機能の組み合わせについては AWS Glue、[を参照してくださいサポートされているコード生成機能](#)。

での Amazon Q データ統合のセットアップ AWS Glue

次のセクションでは、AWS Glueの Amazon Q データ統合の設定について説明します。

トピック

- [IAM 許可の設定](#)

IAM 許可の設定

このトピックでは、Amazon Q チャットエクスペリエンスと AWS Glue Studio ノートブックエクスペリエンスに設定する IAM アクセス権限について説明します。

トピック

- [Amazon Q チャット用の IAM アクセス権限の設定](#)
- [Studio ノートブックの IAM 権限の設定 AWS Glue](#)

Amazon Q チャット用の IAM アクセス権限の設定

の Amazon Q データ統合で使用される API にアクセス権限を付与するには、適切な AWS Identity and Access Management (IAM) AWS Glue 権限が必要です。AWS 以下のカスタムポリシーを IAM ID (ユーザー、ロール、グループなど) にアタッチすることでアクセス権限を取得できます。

```
{
```

```
"Version": "2012-10-17",
"Statement": [
  {
    "Effect": "Allow",
    "Action": [
      "glue:StartCompletion",
      "glue:GetCompletion"
    ],
    "Resource": [
      "arn:aws:glue:*:*:completion/*"
    ]
  }
]
```

Studio ノートブックの IAM 権限の設定 AWS Glue

AWS Glue Studio ノートブックで Amazon Q データ統合を有効にするには、ノートブックの IAM ロールに次の権限が付与されていることを確認してください。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "glue:StartCompletion",
        "glue:GetCompletion"
      ],
      "Resource": [
        "arn:aws:glue:*:*:completion/*"
      ]
    },
    {
      "Sid": "CodeWhispererPermissions",
      "Effect": "Allow",
      "Action": [
        "codewhisperer:GenerateRecommendations"
      ],
      "Resource": "*"
    }
  ]
}
```

Note

の Amazon Q データ統合には、AWS SDK を介してプログラムで使用できる API はありません。AWS Glue Amazon Q チャットパネルまたは AWS Glue Studio ノートブックを通じてこのエクスペリエンスを有効にするために、IAM ポリシーでは次の 2 つの API が使用されています: StartCompletion と GetCompletion

権限の割り当て

アクセスを提供するには、ユーザー、グループ、またはロールにアクセス許可を追加します。

- AWS IAM ID センターのユーザーとグループ: 権限セットを作成します。「AWS IAM Identity Center のユーザーガイド」の「[許可セットの作成](#)」の手順に従ってください。
- アイデンティティプロバイダーを介して IAM で管理されているユーザー: ID フェデレーションのロールを作成します。詳細については、「IAM ユーザーガイド」の「[サードパーティー ID プロバイダー \(フェデレーション\) 用のロールの作成](#)」を参照してください。
- IAM ユーザー:
 - ユーザーが担当できるロールを作成します。手順については、「IAM ユーザーガイド」の「[IAM ユーザー用ロールの作成](#)」を参照してください。
 - (お奨めできない方法) ポリシーをユーザーに直接アタッチするか、ユーザーをユーザーグループに追加する。詳細については、「IAM ユーザーガイド」の「[ユーザー \(コンソール\) へのアクセス権限の追加](#)」を参照してください。

サポートされているコード生成機能

AWS Glue の Amazon Q データ統合のコード生成機能の組み合わせは次のとおりです。

ソース	[Transformation] (変換)	Target
以下のフォーマットタイプで使用する S3: json、csv、parquet、hudi、delta	ApplyMapping	以下のフォーマットタイプで使用する S3: json、csv、avro、orc、parquet、hudi、delta
Glue Data Catalog	RenameField	Glue Data Catalog

ソース	[Transformation] (変換)	Target
Amazon Redshift	DropFields	Amazon Redshift
MySQL	SelectFields	MySQL
Postgres	DropNullFields	Postgres
Oracle	フィルター	Oracle
SQL Server	スピゴット	SQL Server
DynamoDB	カスタム SQL コード	DynamoDB
Snowflake	集計	Snowflake
MongoDB	DropDuplicates	MongoDB
カスタム JDBC コネクタ	Join	カスタム JDBC コネクタ
カスタム・スパーク・コネクタ	Union	カスタム・スパーク・コネクタ
グーグル BigQuery		グーグル BigQuery
Teradata		Teradata
Amazon OpenSearch サービス		Amazon OpenSearch サービス
Vertica		Vertica
Azure SQL		アズール DQL
SAP HANA		SAP HANA
アズール・コスモス		アズール・コスモス

インタラクションの例

の Amazon Q AWS Glue データ統合により、Amazon Q パネルに質問を入力できます。AWS Glue が提供するデータ統合機能に関する質問を入力できます。詳細な回答は、参照ドキュメントとともに返されます。

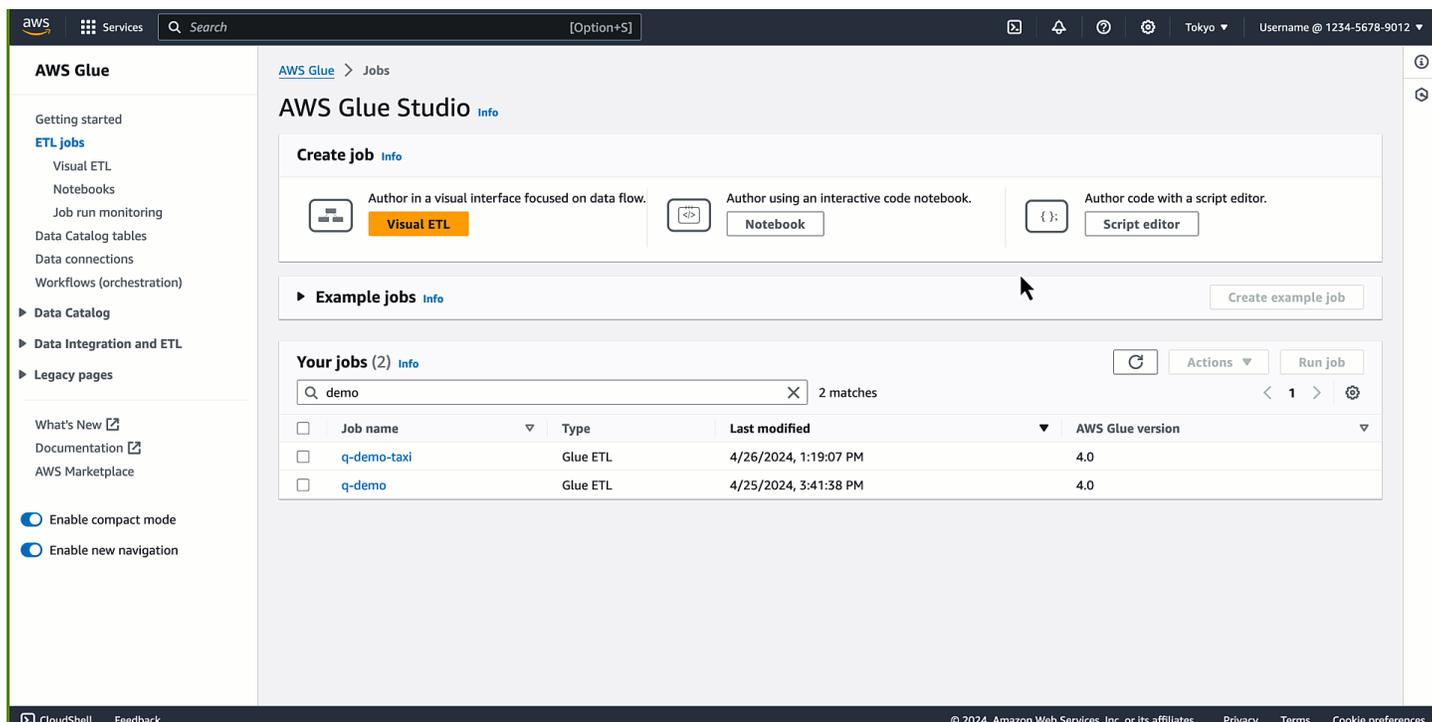
もう 1 つのユースケースは AWS Glue ETL ジョブスクリプトの生成です。データの抽出、変換、読み込みジョブの実行方法について質問できます。PySpark 生成されたスクリプトが返されます。

トピック

- [Amazon Q チャットインタラクション](#)
- [AWS Glue Studio ノートブックのインタラクション](#)

Amazon Q チャットインタラクション

AWS Glue コンソールで新しいジョブの作成を開始し、Amazon Q に「Snowflake から読み取り、フィールドの名前を変更して Redshift に書き込む Glue スクリプトを提供してください」と依頼します。



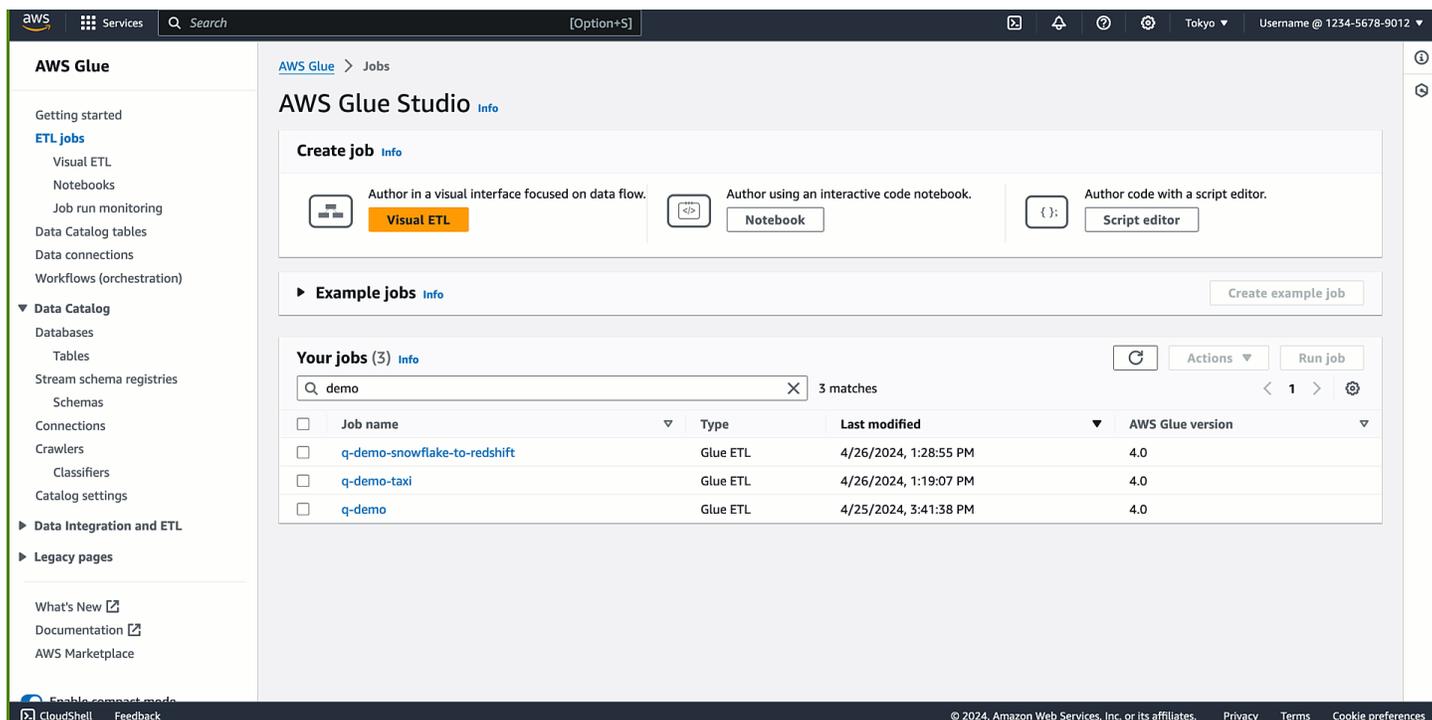
The screenshot shows the AWS Glue Studio console. The left sidebar contains navigation options for AWS Glue, including ETL jobs, Data Catalog, and Legacy pages. The main content area is titled 'AWS Glue Studio' and features a 'Create job' section with three options: 'Visual ETL', 'Notebook', and 'Script editor'. Below this is an 'Example jobs' section with a 'Create example job' button. The 'Your jobs (2)' section displays a table of existing jobs.

<input type="checkbox"/>	Job name	Type	Last modified	AWS Glue version
<input type="checkbox"/>	q-demo-taxi	Glue ETL	4/26/2024, 1:19:07 PM	4.0
<input type="checkbox"/>	q-demo	Glue ETL	4/25/2024, 3:41:38 PM	4.0

コードが生成されていることに気付くでしょう。このレスポンスがあれば、AWS Glue 目的に合ったコードを作成する方法を学び、理解することができます。生成されたコードをスクリプトエディ

ターにコピーして貼り付け、プレースホルダーを設定できます。ジョブに AWS Identity and Access Management (IAM) AWS Glue ロールと接続を設定したら、ジョブを保存して実行します。ジョブが完了すると、Amazon Redshift の Snowflake からエクスポートされたテーブルのクエリを開始できます。

次のプロンプトは、2つの異なるソースからデータを読み取り、それらを個別にフィルタリングして投影し、共通のキーで結合し、出力を3番目のターゲットに書き込みます。Amazon Q に質問: 「S3 から Parquet 形式でデータを読み取り、いくつかのフィールドを選択したい。また、DynamoDB からデータを読み込み、いくつかのフィールドを選択し、いくつかの行をフィルタリングしたいと考えています。これら2つのデータセットを結合して、結果をに書き込みたいと思います。OpenSearch



The screenshot shows the AWS Glue Studio interface. The left sidebar contains navigation options like 'Getting started', 'ETL jobs', 'Data Catalog', and 'Data Integration and ETL'. The main content area is titled 'AWS Glue Studio' and features a 'Create job' section with three options: 'Visual ETL', 'Notebook', and 'Script editor'. Below this is an 'Example jobs' section with a 'Create example job' button. The 'Your jobs (3)' section displays a table of existing jobs:

Job name	Type	Last modified	AWS Glue version
q-demo-snowflake-to-redshift	Glue ETL	4/26/2024, 1:28:55 PM	4.0
q-demo-taxi	Glue ETL	4/26/2024, 1:19:07 PM	4.0
q-demo	Glue ETL	4/25/2024, 3:41:38 PM	4.0

コードが生成されます。ジョブが完了すると、OpenSearch インデックスはダウンストリームのワークロードで使用できるようになり、使用できるようになります。

AWS Glue Studio ノートブックのインタラクション

新しいセルを追加し、達成したいことを説明するコメントを入力します。Tab キーと Enter キーを押すと、推奨コードが表示されます。

最初の目的はデータを抽出することです。「Glue Data Catalog テーブルを読み取るコードをください」、続いて「star_rating>3 でフィルター変換を適用するコードをください」と「フレームを Parquet として S3 に書き込むコードをください」と続きます。

q-nodes [🔗](#) Stop notebook Download Notebook Actions ▼ Save Run

Notebook | **Script** | Job details | Runs | **Data quality - updated** | Schedules | Version Control

```

+ [ ]: [ ]
Worker Type: G.1X
Number of Workers: 5
Session ID: a6846a9a-6489-4599-bf8d-066b59d887da
Applying the following default arguments:
--glue_kernel_version 1.0.4
--enable-glue-datacatalog true
Waiting for session a6846a9a-6489-4599-bf8d-066b59d887da to get into ready status...
Session a6846a9a-6489-4599-bf8d-066b59d887da has been created.
    
```

0 s 1 Initialized (additional servers needed) Glue PySpark | Idle ✓ CodeWhisperer Mode: Edit Ln 1, Col 1 Untitled.ipynb 0

loudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie pre

q-nodes [🔗](#) Stop notebook Download Notebook Actions ▼ Save Run

Notebook | **Script** | Job details | Runs | **Data quality - updated** | Schedules | Version Control

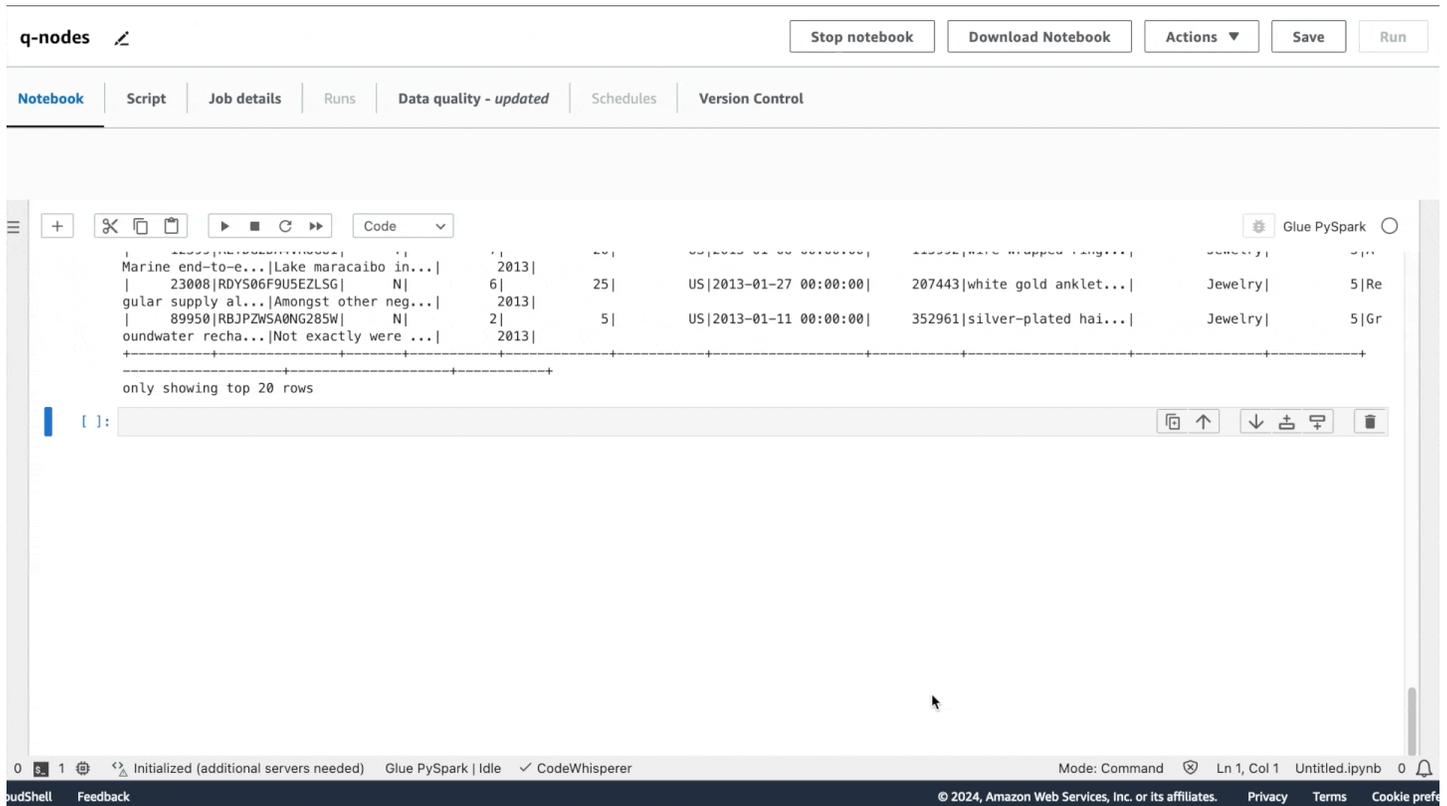
```

+ [ ]: [ ]
|          US| 171306|R00GN0TEQS4ISDM| 90211|white and yellow ...| 3| 5| 5| N|PAP, and regular ...|Words themselves
...|2013-01-23 00:00:00| 2013| Jewelry|
-----
only showing top 20 rows

/opt/amazon/spark/python/lib/pyspark.zip/pyspark/sql/dataframe.py:127: UserWarning: DataFrame constructor is internal. Do not directly use it.
    
```

0 s 1 Initialized (additional servers needed) Glue PySpark | Idle ✓ CodeWhisperer Mode: Edit Ln 1, Col 1 Untitled.ipynb 0

loudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie pre



The screenshot shows the AWS Glue Studio interface. At the top, there are buttons for 'Stop notebook', 'Download Notebook', 'Actions', 'Save', and 'Run'. Below these are tabs for 'Notebook', 'Script', 'Job details', 'Runs', 'Data quality - updated', 'Schedules', and 'Version Control'. The main area displays a data table with columns for product names, years, quantities, and categories. Below the table, there is a code cell with a prompt: 'only showing top 20 rows'. The bottom status bar shows 'Glue PySpark | Idle' and 'CodeWhisperer'.

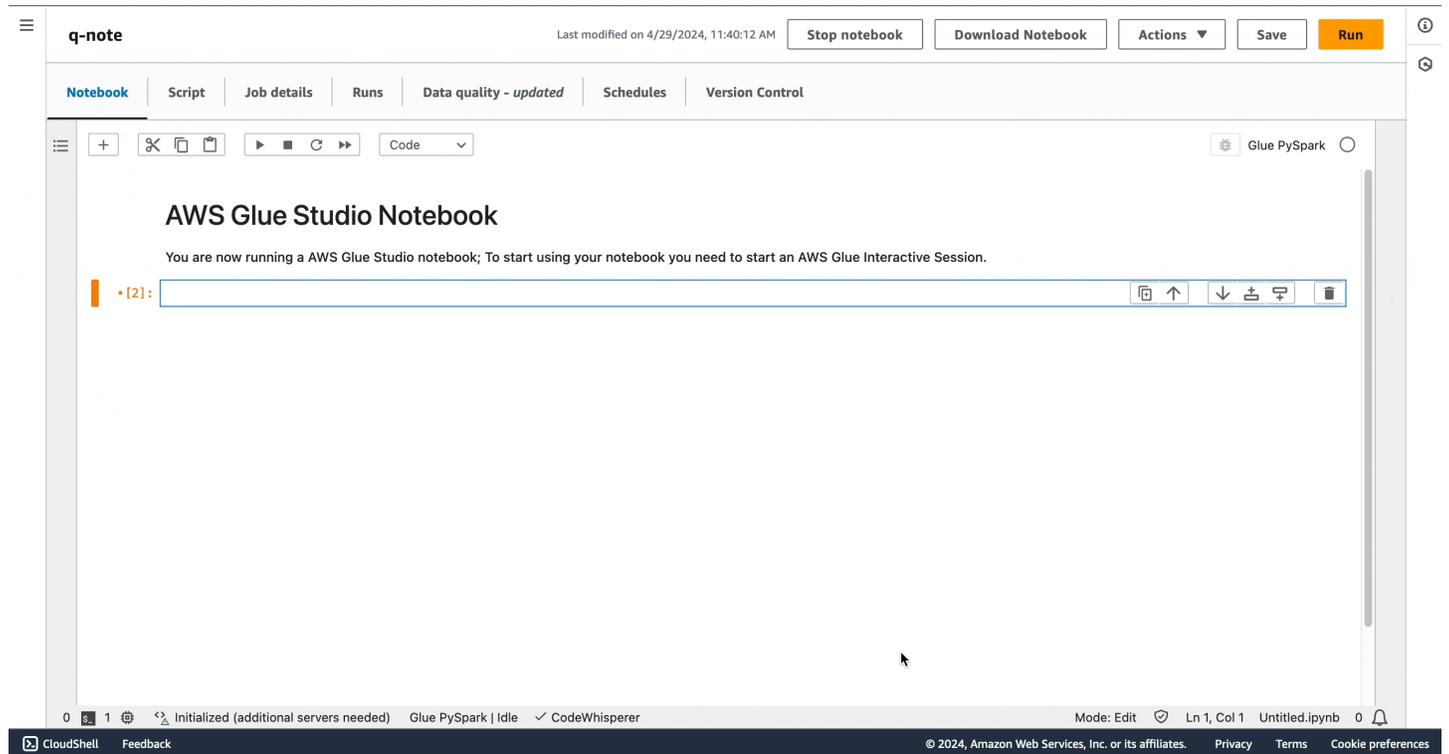
Marine end-to-e... Lake maracaibo in...	2013	25	US 2013-01-27 00:00:00	207443 white gold anklet...	Jewelry	5 Re
23008 RDYS06F9U5EZLSG	N	6				
gular supply al... Amongst other neg...	2013	5	US 2013-01-11 00:00:00	352961 silver-plated hai...	Jewelry	5 Gr
89950 RBJPZWSA0NG285W	N	2				
oundwater recha... Not exactly were ...	2013					

Amazon Q のチャットエクスペリエンスと同様に、このコードは推奨されます。Tab キーを押すと、推奨コードが選択されます。

生成コード内のソースに適したオプションを入力することで、各セルを実行できます。実行中のどの時点でも、`show()`メソッドを使用してデータセットのサンプルをプレビューできます。

複雑なプロンプト

1つの複雑なプロンプトで完全なスクリプトを生成できます。「S3 には JSON データがあり、Oracle には結合が必要なデータがあります。両方のソースから読み取り、結合を行い、結果を Redshift に書き込む Glue スクリプトを提供してください。」



The screenshot displays the AWS Glue Studio Notebook interface. At the top, there's a navigation bar with a hamburger menu, the notebook name 'q-note', and a timestamp 'Last modified on 4/29/2024, 11:40:12 AM'. Action buttons include 'Stop notebook', 'Download Notebook', 'Actions', 'Save', and 'Run'. Below this is a tabbed interface with 'Notebook' selected. The main content area shows the notebook title 'AWS Glue Studio Notebook' and a message: 'You are now running a AWS Glue Studio notebook; To start using your notebook you need to start an AWS Glue Interactive Session.' A code editor below contains a single line of code: '+ [2]:'. The bottom status bar indicates 'Mode: Edit', 'Ln 1, Col 1', and 'Untitled.ipynb'. A footer bar contains 'CloudShell', 'Feedback', and copyright information: '© 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences'.

ノートブックでは、Amazon Q データ統合が Amazon Q AWS Glue チャットで生成されたのと同じコードスニペットを生成したことに気付くかもしれません。

[Run] を選択するか、プログラムで、ノートブックをジョブとして実行できます。

AWS Glue でのオーケストレーション

以下のセクションでは、のジョブのオーケストレーションの情報を提供しますAWS Glue。

トピック

- [トリガーを使用したジョブとクローラの開始](#)
- [AWS Glue で設計図とワークフローを使用して複雑な ETL アクティビティを実行する](#)
- [AWS Glue のブループリントの開発](#)

トリガーを使用したジョブとクローラの開始

AWS Glue では、トリガーと呼ばれる Data Catalog オブジェクトを作成できます。このオブジェクトを使用して、手動または自動で 1 つ以上のクローラを開始したり、抽出、変換、ロード (ETL) ジョブを実行したりできます。トリガーを使用して、依存関係にあるジョブとクローラのチェーンを設計できます。

Note

この同じことは、ワークフローを定義することでも可能です。ワークフローは、複雑なマルチジョブ ETL オペレーションを作成する場合に適しています。詳しくは、「[the section called “設計図とワークフローを使用した複雑な ETL アクティビティの実行”](#)」を参照してください。

トピック

- [AWS Glue トリガー](#)
- [トリガーの追加](#)
- [トリガーの有効化と無効化](#)

AWS Glue トリガー

起動されると、トリガーは指定されたジョブとクローラを開始できます。トリガーは、オンデマンドで、スケジュールに基づいて、またはイベントの組み合わせに基づいて起動します。

Note

1つのトリガーでアクティブにできるクローラは2つだけです。複数のデータストアをクローラする場合は、複数のクローラを同時に実行するのではなく、クローラごとに複数のソースを使用します。

トリガーは作成されると複数の状態のいずれかになります。たとえば、CREATED、ACTIVATED、または DEACTIVATED になります。ACTIVATING などの移行状態もあります。トリガーの起動を一時的に停止するために、トリガーを無効化できます。その後、再度有効化できます。

以下の3種類のトリガーがあります。

予定

cron に基づく時間ベースのトリガー。

スケジュールに基づいて、一連のジョブまたはクローラのトリガーを作成できます。ジョブまたはクローラが実行される頻度、実行される曜日、実行される時間などの制約を指定できます。これらの制約は cron に基づいています。トリガーにスケジュールを設定するときは、cron の機能と制限を考慮してください。たとえば、毎月 31 日にクローラを実行することを選択した場合、いくつかの月には 31 日間はないことに注意してください。cron の詳細については、[ジョブとクローラの時間ベースのスケジュール](#) を参照してください。

条件付き

前のジョブまたはクローラや、複数のジョブまたはクローラが、条件のリストを満たしたときに起動するトリガー。

条件付きトリガーを作成するときは、監視対象のジョブとクローラのリストを指定します。監視対象のジョブまたはクローラごとに、監視するステータス (成功、失敗、タイムアウトなど) を指定します。監視対象のジョブまたはクローラが指定されたステータスで終了すると、トリガーが起動します。監視対象のイベントのいずれかまたはすべてが発生したときに起動するようにトリガーを設定できます。

たとえば、ジョブ J1 とジョブ J2 の両方が正常に完了したときにジョブ J3 を開始するトリガー T1 を設定し、ジョブ J1 またはジョブ J2 のいずれかが失敗した場合にジョブ J4 を開始する別のトリガー T2 を設定できます。

以下の表に、トリガーの監視対象のジョブとクローラの完了状態 (イベント) を示しています。

ジョブの完了状態	クローラの完了状態
<ul style="list-style-type: none">• SUCCEEDED• STOPPED• FAILED• TIMEOUT	<ul style="list-style-type: none">• SUCCEEDED• FAILED• CANCELLED

オンデマンド

有効化されると起動するトリガー。オンデマンドトリガーが ACTIVATED または DEACTIVATED 状態になることはありません。それらのトリガーは常に CREATED 状態のままです。

それらのトリガーは作成されるとすぐに起動可能な状態になるように、作成時に、スケジュールトリガーと条件付きトリガーを有効化するフラグを設定できます。

Important

他のジョブまたはクローラの完了の結果として実行されるジョブまたはクローラは、依存関係にあると言われます。依存関係にあるジョブまたはクローラは、完了したジョブまたはクローラがトリガーによって開始された場合にのみ開始されます。依存関係チェーンのすべてのジョブまたはクローラは、1つのスケジュールまたはオンデマンドトリガーの子であることが必要です。

トリガーを使用してジョブパラメータを渡す

トリガーは、開始するジョブにパラメータを渡すことができます。パラメータとしては、ジョブ引数、タイムアウト値、セキュリティ設定などがあります。トリガーが複数のジョブを開始する場合、パラメータは各ジョブに渡されます。

以下に示しているのは、トリガーによって渡されるジョブ引数のルールです。

- キーと値のペアのキーがデフォルトのジョブ引数と一致する場合、渡された引数はデフォルトの引数を上書きします。キーがデフォルトの引数と一致しない場合、その引数は追加の引数としてジョブに渡されます。
- キーと値のペアのキーが上書不可の引数と一致する場合、渡された引数は無視されます。

詳細については、AWS Glue API の「[the section called “トリガー”](#)」を参照してください。

トリガーの追加

トリガーは、AWS Glue コンソール、AWS Command Line Interface (AWS CLI)、または AWS Glue API を使用して追加できます。

Note

現在、トリガーを使用するとき、AWS Glue コンソールはジョブのみをサポートし、クローラはサポートしません。AWS CLI または AWS Glue API を使用すると、ジョブとクローラの両方でトリガーを設定できます。

トリガーを追加するには (コンソール)

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインの [ETL] で、[トリガー] を選択します。[トリガーを追加] を選択します。
3. 以下のプロパティを指定します。

名前

トリガーに一意の名前を付けます。

トリガータイプ

次のいずれかを指定します。

- [Schedule (スケジュール)]: トリガーは特定の頻度と時間で起動します。
 - [ジョブイベント]: 条件付きトリガー。トリガーは、リスト内のいずれかまたはすべてのジョブが、指定されたステータスと一致すると起動します。トリガーを起動するには、監視されたジョブがトリガーによって開始されている必要があります。どのジョブを選択した場合でも、監視できるジョブイベントは 1 つ (完了ステータス) のみです。
 - [オンデマンド]: トリガーは、有効化されると起動します。
4. トリガーウィザードを完了します。[Review] (確認) ページで [Schedule] (スケジュール) および [Job events] (ジョブイベント) (条件付き) は、[Enable trigger on creation] (作成時にトリガーを有効化する) を選択すると直ちにトリガーを有効化できます。

トリガーを追加するには (AWS CLI)

- 以下のようなコマンドを入力します。

```
aws glue create-trigger --name MyTrigger --type SCHEDULED --schedule "cron(0 12 * * ? *)" --actions CrawlerName=MyCrawler --start-on-creation
```

このコマンドは、MyTrigger という名前のスケジュールトリガーを作成します。このトリガーは、毎日 UTC 午後 12 時に実行され、MyCrawler という名前のクローラを開始します。トリガーは有効化された状態で作成されます。

詳細については、[the section called “AWS Glue トリガー”](#) を参照してください。

ジョブとクローラの時間ベースのスケジュール

AWS Glue では、ジョブとクローラの時間ベースのスケジュールを定義できます。これらのスケジュールの定義は、Unix と同様の [cron](#) 構文を使用します。[協定世界時 \(UTC\)](#) で時間を指定します。スケジュールの最小精度は 5 分です。

スケジュールを使用して実行するようにジョブとクローラを設定する方法の詳細については、「[トリガーを使用したジョブとクローラの開始](#)」を参照してください。

cron 式

cron 式には 6 つの必須フィールドがあり、それらは空白で区切られます。

[Syntax] (構文)

```
cron(Minutes Hours Day-of-month Month Day-of-week Year)
```

フィールド	値	ワイルドカード
分	0 ~ 59	, - * /
時間	0 ~ 23	, - * /
日	1 ~ 31	, - * ? / L W
月	1 ~ 12 または JAN ~ DEC	, - * /

フィールド	値	ワイルドカード
曜日	1~7 または SUN~SAT	, - * ? / L
年	1970~2199	, - * /

ワイルドカード

- ,(カンマ) のワイルドカードには、追加の値が含まれます。Month フィールドの、JAN,FEB,MAR は、1月、2月、3月を含みます。
- -(ダッシュ) のワイルドカードは、範囲を指定します。Day フィールドの、「1-15」は、指定した月の1日から15日を含みます。
- [*] (アスタリスク) のワイルドカードには、フィールドのすべての値が含まれます。Hours フィールドの、* にはすべての時間が含まれています。
- /(スラッシュ) のワイルドカードは、増分を指定します。Minutes フィールドで、「**1/10**」と入力して、その時間の最初の分から始めて、10分毎を指定できます(11分、21分、31分など)。
- [?] (疑問符) のワイルドカードは、任意を意味します。Day-of-month フィールドで7と入力し、7日が何曜日であってもかまわない場合、Day-of-week フィールドに?を入力できます。
- Day-of-month フィールドまたは Day-of-week フィールドにある [L] のワイルドカードは、月または週の最終日を指定します。
- Day-of-month フィールドの、ワイルドカード W は、平日を指定します。Day-of-month フィールドで、3W は月の3番目の平日に最も近い日を指定します。

制限

- Cron 式の Day-of-month フィールドと Day-of-week フィールドを同時に指定することはできません。一方のフィールドに値を指定する場合、もう一方のフィールドで [?] (疑問符) を使用する必要があります。
- 5分より短い間隔を導き出す cron 式はサポートされていません。

例

スケジュールを作成するときは、以下のサンプルの cron 文字列を使用できます。

分	時間	日	月	曜日	年	意味
0	10	*	*	?	*	毎日午前 10:00 (UTC) に実行
15	12	*	*	?	*	毎日午後 12:15 (UTC) に実行
0	18	?	*	MON-FRI	*	毎週月曜日から金曜日まで午後 6:00 (UTC) に実行
0	8	1	*	?	*	毎月 1 日の午前 8:00 (UTC) に実行
0/15	*	*	*	?	*	15 分ごとに実行
0/10	*	?	*	MON-FRI	*	月曜日から金曜日まで 10 分ごとに実行
0/5	8 ~ 17	?	*	MON-FRI	*	毎週月曜日から金曜日まで午前 8:00 から午後 5:55 (UTC) の間

分	時間	日	月	曜日	年	意味
						に 5 分ごとに実行

たとえば、毎日 12:15 UTC のスケジュールで実行するには、次のように指定します。

```
cron(15 12 * * ? *)
```

トリガーの有効化と無効化

AWS Glue コンソール、AWS Command Line Interface (AWS CLI)、または AWS Glue API を使用して、トリガーを有効または無効にできます。

トリガーを有効化または無効化するには (コンソール)

1. にサインイン AWS Management Console し、<https://console.aws.amazon.com/glue/> でAWS Glueコンソールを開きます。
2. ナビゲーションペインの [ETL] で、[トリガー] を選択します。
3. 目的のトリガーの横にあるチェックボックスをオンにし、[Action (アクション)] メニューで [トリガーの有効化] を選択して、トリガーを有効化します。または、[トリガーを無効化] を選択して、トリガーを無効化します。

トリガーを有効化または無効化するには (AWS CLI)

- 以下のいずれかのコマンドを入力します。

```
aws glue start-trigger --name MyTrigger
```

```
aws glue stop-trigger --name MyTrigger
```

トリガーは、開始することで有効化され、停止することで無効化されます。オンデマンドトリガーは、有効化されると、すぐに起動します。

詳細については、[the section called “AWS Glue トリガー”](#) を参照してください。

AWS Glue で設計図とワークフローを使用して複雑な ETL アクティビティを実行する

組織が実行する複雑な抽出、変換、ロード (ETL) プロセスの中には、依存関係のある複数の AWS Glue ジョブとクローラによる実装が最適なものがあります。AWS Glue ワークフローを使用すると、マルチジョブ、マルチクローラの ETL プロセスを設計でき、このプロセスは単一のエンティティとして実行と追跡が可能です。ワークフローを作成し、その中でジョブ、クローラ、トリガーを指定した後に、オンデマンドまたはスケジュールにより、そのワークフローを実行できるようになります。

トピック

- [AWS Glue のワークフローの概要](#)
- [AWS Glue でワークフローを手動により作成および構築する](#)
- [Amazon EventBridge イベントによる AWS Glue ワークフローの開始](#)
- [ワークフローを開始した EventBridge イベントの表示](#)
- [AWS Glue でのワークフローの実行とモニタリング](#)
- [ワークフロー実行の停止](#)
- [ワークフロー実行の修復と再開](#)
- [AWS Glue でのワークフローの実行プロパティの取得と設定](#)
- [AWS Glue APIを使用したワークフローのクエリ](#)
- [AWS Glue での設計図とワークフローの制限](#)
- [AWS Glue での設計図エラーをトラブルシューティングする](#)
- [AWS Glue ブループリントでのペルソナおよびロール用のアクセス許可](#)

AWS Glue のワークフローの概要

AWS Glue では、複数のクローラ、ジョブ、およびトリガーを伴う複雑な ETL (抽出、変換、ロード) アクティビティを作成して可視化できます。各ワークフローは、含まれるすべてのジョブならびにクローラの実行とモニタリングを管理します。ワークフローは、各コンポーネントについて、その実行の進捗状況とステータスを記録します。これにより、タスク全体の概要と各ステップの詳細を把握できます。AWS Glue コンソールは、ワークフローの状態をグラフで表示します。

ワークフローは、AWS Glue 設計図から作成することができます。あるいは、AWS Management Console または AWS Glue API を使用しながら、ワークフローとコンポーネントを毎回手動で作成

することも可能です。設計図の詳細については、「[the section called “設計図の概要”](#)」を参照してください。

ワークフロー内のトリガーは、ジョブおよびクローラの両方を開始でき、同時に、ジョブもしくはクローラの完了によって起動することができます。トリガーを使用することで、相互に依存するジョブとクローラの大規模なチェーンを作成できます。ワークフロー内でジョブとクローラの依存関係を定義するトリガーに加えて、各ワークフローには開始トリガーもあります。開始トリガーには、以下の3種類があります。

- スケジュール – 定義したスケジュールに基づいて、ワークフローが開始されます。スケジュールは、日次、週次、月次その他に設定することも、cron 式に基づきカスタムスケジュールを構成することも可能です。
- オンデマンド – ワークフローは、AWS Glue コンソール、API、あるいは AWS CLI から手動で開始されます。
- EventBridge イベント – ワークフローは、単一の Amazon EventBridge イベントまたは、バッチ化された Amazon EventBridge イベントの発生時に開始されます。このトリガータイプでは、AWS Glue はイベント駆動型アーキテクチャのイベントコンシューマとして機能します。すべての EventBridge イベントタイプがワークフローを開始できます。一般的なユースケースとしては、Amazon S3 バケットに新しいオブジェクトが保存された(S3 の PutObject オペレーション)場合が挙げられます。

イベントのバッチにより開始される場合、ワークフローは、指定された数のイベントを受信するか、指定された時間が経過するまで待機します。EventBridge イベントトリガーを作成するには、オプションでバッチ条件を指定できます。バッチ条件を指定する際には、バッチサイズ (イベント数) を指定する必要があります。また、オプションでバッチウィンドウ (秒数) を指定することもできます。デフォルトで、バッチウィンドウは最大長である 900 秒 (15 分) に設定されています。ワークフローは、最初に満たされたバッチ条件によって開始されます。バッチウィンドウは、最初のイベントが受信された時点で開始されます。トリガーの作成時にバッチ条件を指定しない場合、バッチサイズはデフォルトで 1 に設定されます。

ワークフローが開始されると、バッチ条件はリセットされます。イベントトリガーは、ワークフローを再度開始するために、後続のバッチ条件が満たされるかどうかの監視を開始します。

次の表は、バッチサイズとバッチウィンドウが連携して、ワークフローをトリガーする様子を示しています。

バッチサイズ	バッチウィンドウ	トリガ条件の結果
10		ワークフローは、10 個の EventBridge イベントの受信と、前のイベントから 15 分間の経過のうち、先に発生したいずれかによりトリガーされます。(ウィンドウサイズを指定しない場合のデフォルト設定は 15 分です。)
10	2 分	ワークフローは、10 個の EventBridge イベントの受信と、前のイベントから 2 分間の経過のうち、先に発生したいずれかによりトリガーされます。
1		ワークフローは、1 つ目のイベントの到着時にトリガーされます。ウィンドウサイズの設定は無効です。EventBridge イベントトリガーの作成時にバッチ条件を指定しない場合、バッチサイズはデフォルトで 1 に設定されます。

GetWorkflowRun API オペレーションは、ワークフローをトリガーしたバッチ条件を返します。

ワークフローの作成時は、ワークフローの開始方法に関係なく、ワークフローの同時実行の最大数を指定できます。

イベントまたはイベントのバッチが実行を開始したワークフローが最終的に失敗した場合、対象のイベントまたはイベントのバッチは、以後ワークフロー実行のためには使用されなくなります。新しいワークフローの実行は、次のイベントまたはイベントのバッチを受信した場合にのみ開始されます。

Important

ワークフロー内のジョブ、クローラ、トリガーの総数を、100 以下に制限します。100 を超える値を含めると、ワークフローの実行を再開または停止しようとしたときにエラーが発生する場合があります。

ワークフローに設定された同時実行の上限数を超過している場合は、仮にイベント条件が満たされていても、ワークフローの実行は開始されません。ワークフローの同時実行の制限は、想定されるイ

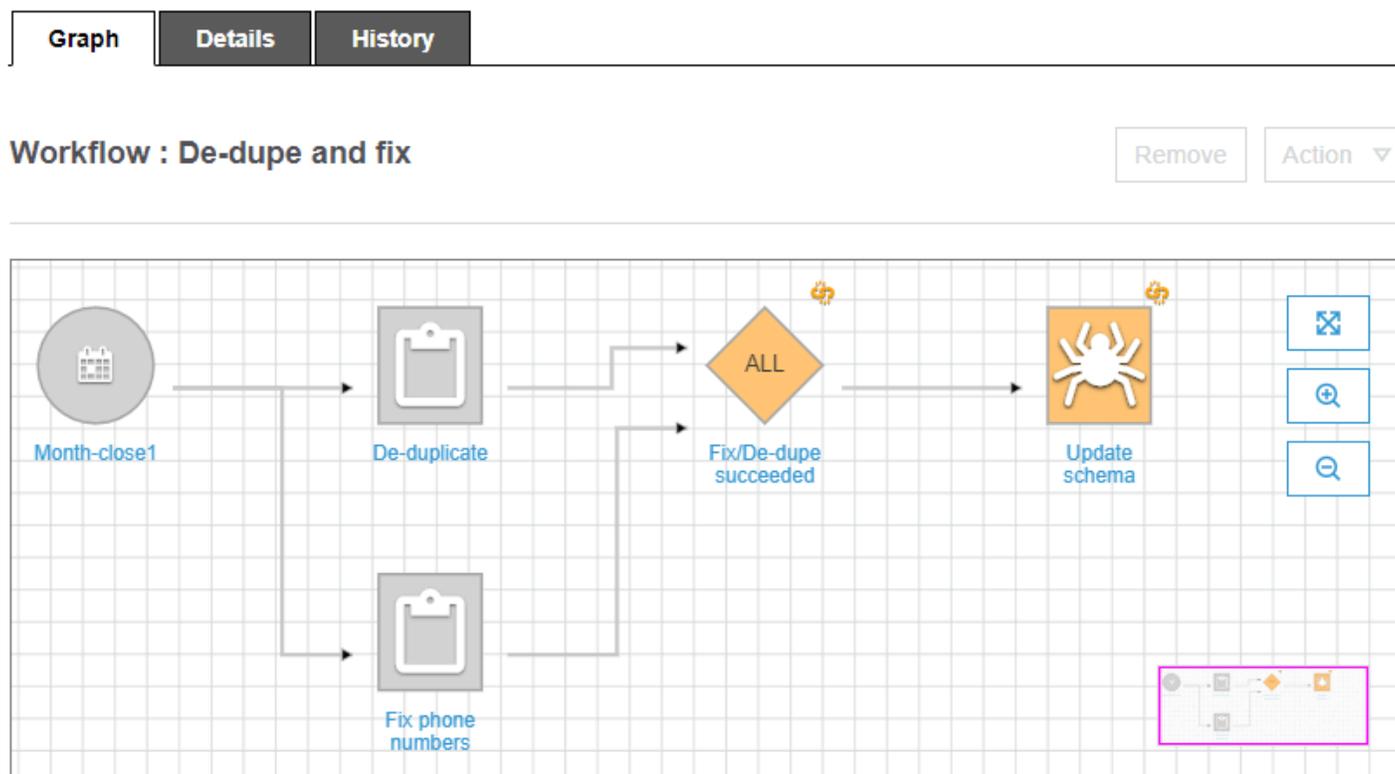
イベントボリュームに基づいて調整することをお勧めします。AWS Glue は、同時実行の制限を超えたために失敗したワークフローの実行を再試行しません。同様に、想定されるイベント量に基づいて、ワークフロー内のジョブとクローラの同時実行の制限を調整することをお勧めします。

ワークフロー実行のプロパティ

ワークフローの実行全体の状態を共有して管理するには、ワークフローのデフォルトの実行プロパティを定義できます。これらのプロパティ (名前と値のペア) は、ワークフローのすべてのジョブで使用できます。AWS Glue API を使用すると、ジョブはワークフローの実行プロパティを取得し、これらのプロパティを変更してワークフローの以降のジョブで使用できます。

ワークフローのグラフ

次の図に、AWS Glue コンソールが表示する、ワークフローの基本的なグラフを示します。ワークフローには、いくつかのコンポーネントが含まれている場合があります。



このワークフローは、スケジュールトリガー Month-close1 により開始されます。このトリガーは、2つのジョブ (De-duplicate および Fix phone numbers) を開始します。これら2つのジョブが正常に完了すると、イベントトリガー Fix/De-dupe succeeded によってクローラ Update schema が開始されます。

ワークフローの静的ビューと動的ビュー

各ワークフローには、静的ビューと動的ビューの表記があります。静的ビューは、ワークフローの設計を示します。動的ビューは、各ジョブや各クローラに関する最新の実行情報を含むランタイムビューです。実行情報には、成功ステータスとエラーの詳細が含まれます。

ワークフローの実行中は、その動的ビューがコンソールに表示されます。また、完了済みのジョブと未実行のジョブがグラフで示されます。実行中のワークフローの動的ビューは、AWS Glue API を使用して取得することもできます。(詳細については、[AWS Glue APIを使用したワークフローのクエリ](#)を参照してください)。

 以下も参照してください。

- [the section called “設計図からのワークフローの作成”](#)
- [the section called “ワークフローの手動による作成と構築”](#)
- [ワークフロー](#) (ワークフロー API 用)

AWS Glue でワークフローを手動により作成および構築する

AWS Glue コンソールにより、ワークフローのノードを一度に 1 つずつ手動で作成し構築できます。

ワークフローは、ジョブ、クローラ、トリガーで構成されます。ワークフローの手動作成を開始する前に、ワークフローに含めるジョブとクローラを作成します。ワークフローのクローラは、オンデマンドで実行するように指定するのが最適です。トリガーは、ワークフローの構築中に新規作成できません。または、既存のトリガーをワークフロー内に複製することもできます。トリガーをクローンすると、トリガーに関連付けられたすべてのカタログオブジェクト (トリガーを起動するジョブまたはクローラ、トリガーにより開始されるジョブまたはクローラ) がワークフローに追加されます。

Important

ワークフロー内のジョブ、クローラ、トリガーの総数を 100 以下に制限します。100 を超える値を含めると、ワークフローの実行を再開または停止しようとしたときにエラーが発生することがあります。

ワークフローを構築するには、ワークフローのグラフにトリガーを追加し、トリガーごとに監視対象のイベントやアクションを定義します。まず、開始トリガーとしてオンデマンドトリガーまたはスケジュールトリガーを追加し、次にイベント (条件付き) トリガーを追加してグラフを完成します。

ステップ 1: ワークフローを作成する

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインの [ETL] で、[ワークフロー] を選択します。
3. [ワークフロー] を選択し、[Add a new ETL workflow (新しい ETL ワークフローの追加)] フォームに入力します。

オプションとして追加したデフォルトの実行プロパティは、ワークフローのすべてのジョブに対する引数として使用できます。詳細については、[AWS Glue でのワークフローの実行プロパティの取得と設定](#) を参照してください。

4. [Add workflow (ワークフローの追加)] を選択します。

新しいワークフローが [ワークフロー] ページのリストに表示されます。

ステップ 2: 開始トリガーを追加する

1. [ワークフロー] ページで、新しいワークフローを選択します。次に、ページの下部で [Graph] (グラフ) タブが選択されていることを確認します。
2. [トリガーを追加] を選択し、[トリガーを追加] ダイアログボックスで、次のいずれかの操作を行います。
 - [Clone existing (既存の複製)] を選択し、複製するトリガーを選択します。その後、[Add] (追加) を選択します。

トリガーがグラフに表示されます。トリガーで監視するジョブやクローラおよびトリガーで開始するジョブやクローラも表示されます。

トリガーを間違えて選択した場合は、そのトリガーをグラフで選択し、[削除] を選択します。

- [Add new (新規追加)] を選択し、[トリガーを追加] フォームに入力します。
 1. [Trigger type] (トリガータイプ) で、[Schedule] (スケジュール)、[On demand] (オンデマンド)、または [EventBridge event] (EventBridge イベント) のいずれかを選択します。

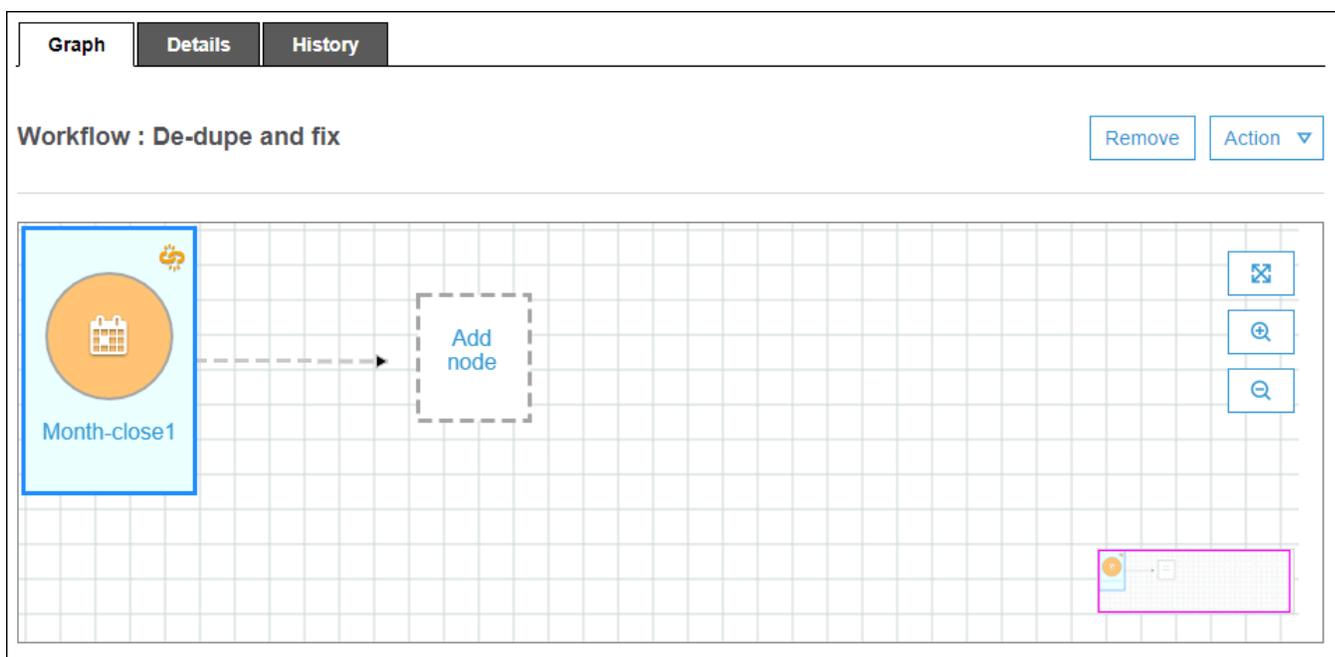
トリガータイプを [Schedule] (スケジュール) にした場合、[Frequency] (頻度) オプションの中から 1 つを選択します。cron 式を入力するには、[Custom] (カスタム) をクリックします。

トリガータイプに [EventBridge event] (EventBridge イベント) を選択した場合は、[Number of events] (イベント数) (バッチサイズ) を入力し、オプションで [Time delay] (遅延時間) (バッチウィンドウ) を入力します。[Time delay] (遅延時間) を省略した場合、デフォルトでバッチウィンドウには 15 分が指定されます。詳細については、[AWS Glue のワークフローの概要](#) を参照してください。

2. 追加] を選択します。

トリガーがグラフ上に表示されます。プレースホルダーノード ([Add node (ノードの追加)] というラベルが付いたノード) も一緒に表示されます。以下の例では、開始トリガーは Month-close1 という名前のスケジュールされたトリガーです。

この段階では、トリガーは保存されていません。



3. 新しいトリガーを追加した場合は、以下のステップを実行します。

a. 次のいずれかを行います。

- プレースホルダーノード ([Add node (ノードの追加)] を選択します。
- 開始トリガーが選択されていることを確認し、[アクション] メニューの [Add jobs/ crawlers to trigger (トリガーにジョブ/クローラを追加)] を選択します。

b. Add job(s) and crawler(s) to trigger (トリガーにジョブおよびクローラを追加) ダイアログボックスで、1 つ以上のジョブまたはクローラを選択し、[追加] を選択します。

トリガーが保存されます。選択したジョブまたはクローラが、トリガーからのコネクタと共にグラフに表示されます。

ジョブやクローラを間違えて追加した場合は、トリガーまたはコネクタを選択して、[削除]を選択できます。

ステップ 3: さらにトリガーを追加する

[Event] (イベント) タイプのトリガーをさらに追加して、ワークフローの構築を続けます。グラフのキャンバスを拡大または縮小するには、グラフの右側にあるアイコンを使用します。追加するトリガーごとに、以下の手順を実行します。

Note

ワークフローの保存のために、行うべきアクションはありません。最後のトリガーを追加し、そのトリガーにアクションを割り当てると、ワークフローが完了し保存されます。後の任意のタイミングでこの作業に戻り、さらにノードを追加することができます。

1. 次のいずれかを行います。

- 既存のトリガーを複製するには、グラフで選択されているノードがないことを確認し、[アクション] メニューの [トリガーを追加] を選択します。
- グラフ上の特定のジョブまたはクローラを監視する新しいトリガーを追加するには、そのジョブまたはクローラのノードを選択し、[トリガーを追加] プレースホルダーノードを選択します。

後のステップで、このトリガーで監視するジョブやクローラをさらに追加できます。

2. [トリガーを追加] ダイアログボックスで、次のいずれかの操作を行います。

- [Add new (新規追加)] を選択し、[トリガーを追加] フォームに入力します。その後、[Add] (追加) を選択します。

トリガーがグラフに表示されます。後のステップでトリガーを完了します。

- [Clone existing (既存の複製)] を選択し、複製するトリガーを選択します。その後、[Add] (追加) を選択します。

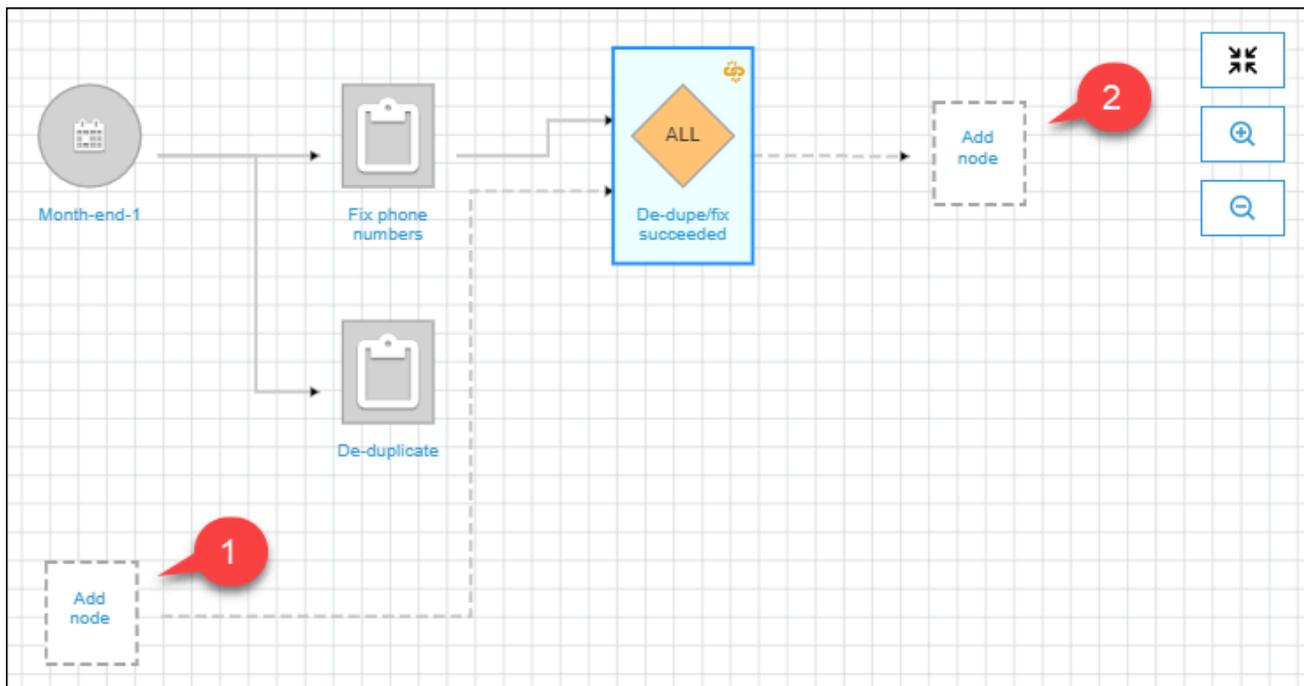
トリガーがグラフに表示されます。トリガーで監視するジョブやクローラおよびトリガーで開始するジョブやクローラも表示されます。

トリガーを間違えて選択した場合は、そのトリガーをグラフで選択し、[削除] を選択します。

3. 新しいトリガーを追加した場合は、以下のステップを実行します。

a. 新しいトリガーを選択します。

次のように、トリガー De-dupe/fix succeeded が選択された状態のグラフが表示され、プレースホルダーノードが (1) イベント用、ならびに (2) アクション用として表示されます。



- b. (トリガーでイベントを既に監視しており、監視対象のジョブやクローラを追加する場合のオプション) 監視対象イベントのプレースホルダーノードを選択し、[Add job(s) and crawler(s) to watch (監視するジョブやクローラの追加)] ダイアログボックスで 1 つ以上のジョブまたはクローラを選択します。監視対象のイベント (成功、失敗など) を選択して、[追加] を選択します。
- c. トリガーが選択されていることを確認し、アクション用のプレースホルダーノードを選択します。
- d. [Add job(s) and crawler(s) to watch (監視対象のジョブやクローラの追加)] ダイアログボックスで 1 つ以上のジョブまたはクローラを選択し、[追加] を選択します。

選択したジョブやクローラがグラフに表示されます。トリガーからのコネクタも表示されません。

Express ワークフローと サービス統合の詳細については、以下を参照してください。

- [AWS Glue のワークフローの概要](#)
- [AWS Glue でのワークフローの実行とモニタリング](#)
- [AWS Glue で設計図からワークフローを作成する](#)

Amazon EventBridge イベントによる AWS Glue ワークフローの開始

Amazon EventBridge (別名 CloudWatch Events) を使用すると、AWS サービスを自動化して、アプリケーションの可用性の問題やリソースの変更などのシステム的なイベントに自動的に対応できます。AWS のサービスからのイベントは、ほぼリアルタイムに EventBridge に提供されます。簡単なルールを記述して、注目するイベントと、イベントがルールに一致した場合に自動的に実行するアクションを指定できます。

EventBridge の機能を使用することで、AWS Glue は、イベント駆動型アーキテクチャにおけるイベントのプロデューサおよびコンシューマとして動作します。ワークフローに対して AWS Glue では、コンシューマとして、あらゆるタイプの EventBridge イベントをサポートしています。このための最も一般的なユースケースとしては、Amazon S3 バケットが新しいオブジェクトを受け取る場合が挙げられます。データが、不規則または定義されていない間隔で到着した際には、可能なかぎり素早くそれを処理できます。

Note

AWS Glue は、EventBridge メッセージの配信を保証していません。EventBridge により重複したメッセージが配信された場合にも、AWS Glue は、重複の排除を行いません。対応すべき処理は、ユースケースに基づいて管理する必要があります。不要なイベントが送信されないように、EventBridge ルール適切に設定します。

開始する前に

Amazon S3 データイベントを使用してワークフローを開始する場合は、対象の S3 バケットのイベントが AWS CloudTrail および EventBridge にログ記録されたことを確認します。そのために

は、CloudTrail の証跡を作成する必要があります。詳細については、「[Creating a trail for your AWS account](#)」を参照してください。

EventBridge イベントを使用してワークフローを開始するには

Note

以下のコマンドで、下記のような置き換えを行います。

- `<workflow-name>` には、ワークフローに割り当てる名前。
- `<trigger-name>` には、トリガーに割り当てる名前。
- `<bucket-name>` には、Amazon S3 バケットの名前。
- `<account-id>` には、有効な AWS アカウント ID。
- `<region>` には、リージョン名 (例: us-east-1)
- `<rule-name>` には、EventBridge ルールに割り当てる名前。

1. EventBridge ルールとターゲットを作成および表示するための、AWS Identity and Access Management (IAM) アクセス許可があることを確認します。以下は、アタッチできるサンプルのポリシーです。オペレーションとリソースを制限するために、このポリシーの範囲を狭めたい場合もあります。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "events:PutRule",
        "events:DisableRule",
        "events>DeleteRule",
        "events:PutTargets",
        "events:RemoveTargets",
        "events:EnableRule",
        "events:List*",
        "events:Describe*"
      ],
      "Resource": "*"
    }
  ]
}
```

```
}
```

2. AWS Glue にイベントが渡される際に、EventBridge サービスが引き受けることができる IAM ロールを作成します。
 - a. IAM コンソールの [Create role] (ロールの作成) ページで、[AWS Service] を選択します。次に、[CloudWatch Events] (CloudWatch イベント) サービスを選択します。
 - b. [Create role] (ロールの作成) ウィザードを完了します。ポリシー `CloudWatchEventsBuiltInTargetExecutionAccess` および `CloudWatchEventsInvocationAccess` が、ウィザードにより自動的にアタッチされます。
 - c. 次のインラインポリシーをロールにアタッチします。このポリシーにより、EventBridge サービスはイベントを AWS Glue に送れるようになります。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "glue:notifyEvent"
      ],
      "Resource": [
        "arn:aws:glue:<region>:<account-id>:workflow/<workflow-name>"
      ]
    }
  ]
}
```

3. 次のコマンドを入力して、ワークフローを作成します。

この他のオプションのコマンドラインパラメータについては、AWS CLI コマンドリファレンスの「[create-workflow](#)」を参照してください。

```
aws glue create-workflow --name <workflow-name>
```

4. 次のコマンドを入力して、ワークフローのための EventBridge イベントトリガーを作成します。これがワークフローの開始トリガーになります。`<actions>`は、実行するアクション (開始するジョブとクローラ) に置き換えます。

引数 actions の記述方法については、AWS CLI コマンドリファレンスの「[create-trigger](#)」を参照してください。

```
aws glue create-trigger --workflow-name <workflow-name> --type EVENT --
name <trigger-name> --actions <actions>
```

単一の EventBridge イベントではなく、バッチされたイベントによってワークフローをトリガーする場合は、代わりに次のコマンドを入力します。

```
aws glue create-trigger --workflow-name <workflow-name> --type EVENT
--name <trigger-name> --event-batching-condition BatchSize=<number-of-
events>,BatchWindow=<seconds> --actions <actions>
```

event-batching-condition 引数の BatchSize は必須であり、BatchWindow はオプションです。BatchWindow が省略された場合、ウィンドウはデフォルトで (最大ウィンドウサイズの) 900 秒に設定されます。

Example

次の例では、3 つの EventBridge イベントを受信した後、または最初のイベントが到着してから 5 分経過後の、いずれか早いタイミングで eventtest ワークフローを開始するトリガーを作成しています。

```
aws glue create-trigger --workflow-name eventtest --type EVENT --name objectArrival
--event-batching-condition BatchSize=3,BatchWindow=300 --actions JobName=test1
```

5. Amazon EventBridge でルールを作成します。

- a. 任意のテキストエディタで JSON オブジェクトを作成し、ルールの詳細を記述します。

次の例では、イベントソースとして Amazon S3 を、イベント名として PutObject を、リクエストパラメータとしてバケット名を、それぞれ指定しています。このルールは、新しいオブジェクトがバケットで受信された場合にワークフローを開始します。

```
{
  "source": [
    "aws.s3"
  ],
  "detail-type": [
```

```

    "AWS API Call via CloudTrail"
  ],
  "detail": {
    "eventSource": [
      "s3.amazonaws.com"
    ],
    "eventName": [
      "PutObject"
    ],
    "requestParameters": {
      "bucketName": [
        "<bucket-name>"
      ]
    }
  }
}

```

新しいオブジェクトがバケット内のフォルダに到着した際にワークフローを開始させるには、次のコードを `requestParameters` に置き換えます。

```

    "requestParameters": {
      "bucketName": [
        "<bucket-name>"
      ]
      "key" : [{ "prefix" : "<folder1>/<folder2>/*"}]}
    }

```

- b. 任意のツールを使用して、ルールを記述した JSON オブジェクトの文字列をエスケープします。

```

{\n  \"source\": [\n    \"aws.s3\"\n  ],\n  \"detail-type\": [\n    \"AWS API Call via CloudTrail\"\n  ],\n  \"detail\": {\n    \"eventSource\": [\n      \"s3.amazonaws.com\"\n    ],\n    \"eventName\": [\n      \"PutObject\"\n    ],\n    \"requestParameters\": {\n      \"bucketName\": [\n        \"<bucket-name>\"\n      ]\n    }\n  }\n}

```

- c. 次のコマンドを実行して、JSON パラメータテンプレートを作成します。このテンプレートを編集して、後続の `put-rule` コマンドの入力パラメータを指定します。出力をファイルに保存します。この例では、`ruleCommand` というファイルに保存しています。

```
aws events put-rule --name <rule-name> --generate-cli-skeleton >ruleCommand
```

--generate-cli-skeleton パラメータの詳細については、AWS Command Line Interface ユーザーガイドの「[Generating AWS CLI skeleton and input parameters from a JSON or YAML input file](#)」を参照してください。

出力ファイルは以下のようになります。

```
{
  "Name": "",
  "ScheduleExpression": "",
  "EventPattern": "",
  "State": "ENABLED",
  "Description": "",
  "RoleArn": "",
  "Tags": [
    {
      "Key": "",
      "Value": ""
    }
  ],
  "EventBusName": ""
}
```

- d. このファイルを編集して、少なくとも Name、EventPattern、および State パラメータを指定しながら、必要に応じて他のパラメータを削除します。EventPattern パラメータでは、前のステップで作成したルール詳細のために、エスケープされた文字列を指定します。

```
{
  "Name": "<rule-name>",
  "EventPattern": "{\n  \"source\": [\n    \"aws.s3\"\n  ],\n  \"detail-type\": [\n    \"AWS API Call via CloudTrail\"\n  ],\n  \"detail\": {\n    \"eventSource\": [\n      \"s3.amazonaws.com\"\n    ],\n    \"eventName\": [\n      \"PutObject\"\n    ],\n    \"requestParameters\": {\n      \"bucketName\": [\n        \"<bucket-name>\"\n      ]\n    }\n  }\n}",
  "State": "DISABLED",
  "Description": "Start an AWS Glue workflow upon new file arrival in an Amazon S3 bucket"
}
```

Note

ワークフローの構築が完了するまで、ルールを無効にしたままにしておくことをお勧めします。

- e. 次の `put-rule` コマンドを入力し、ファイル `ruleCommand` から入力パラメータを読み取ります。

```
aws events put-rule --name <rule-name> --cli-input-json file://ruleCommand
```

次の出力は、正しく処理されたことを示しています。

```
{
  "RuleArn": "<rule-arn>"
}
```

6. 次のコマンドを入力して、ターゲットにルールをアタッチします。ターゲットは、AWS Glue のワークフローです。<rule-name>は、この手順の最初に作成したルールに置き換えます。

```
aws events put-targets --rule <rule-name> --targets
  "Id"="1", "Arn"="arn:aws:glue:<region>:<account-id>:workflow/<workflow-
  name>", "RoleArn"="arn:aws:iam::<account-id>:role/<role-name>" --region <region>
```

次の出力は、正しく処理されたことを示しています。

```
{
  "FailedEntryCount": 0,
  "FailedEntries": []
}
```

7. 次のコマンドを入力して、ルールとターゲットの接続が正常に行われたことを確認します。

```
aws events list-rule-names-by-target --target-arn arn:aws:glue:<region>:<account-
  id>:workflow/<workflow-name>
```

接続に成功した場合は次のような出力になります。<rule-name> は作成したルールの名前です。

```
{
  "RuleNames": [
    "<rule-name>"
  ]
}
```

8. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
9. ワークフローを選択し、ワークフローグラフで、開始トリガーとそのアクション (ワークフローが開始するジョブまたはクローラ) の表示を確認します。その後、[ステップ 3: さらにトリガーを追加する](#) の手順に進みます。または、AWS Glue API か AWS Command Line Interface を使用して、ワークフローに他のコンポーネントを追加します。
10. ワークフローの指定を完了したら、ルールを有効にします。

```
aws events enable-rule --name <rule-name>
```

これで、EventBridge イベントまたはイベントのバッチによって、ワークフローを開始する準備が整いました。

i 以下も参照してください。

- [Amazon EventBridge ユーザーガイド](#)
- [AWS Glue のワークフローの概要](#)
- [AWS Glue でワークフローを手動により作成および構築する](#)

ワークフローを開始した EventBridge イベントの表示

ワークフローを開始した Amazon EventBridge イベントのイベント ID を表示できます。ワークフローがバッチされたイベントによって開始された場合は、そのバッチ内のすべてのイベントの、イベント ID を表示することが可能です。

バッチサイズが 1 より大きいワークフローでは、どのバッチ条件 (バッチサイズでのイベントの到着数、またはバッチウィンドウの有効期限) によってワークフローが開始されたかを確認することもできます。。

ワークフローを開始した EventBridge イベントを表示するには (コンソール)

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインで、[Workflows] (ワークフロー) をクリックします。
3. ワークフローを選択します。下部で [History] (履歴) タブを選択します。
4. ワークフロー実行を選択し、[View run details] (実行の詳細を表示する) をクリックします。
5. 実行の詳細ページにある [Run properties] (実行プロパティ) フィールドで、[aws:eventIds] キーを検索します。

このキーの値は、EventBridge イベント ID のリストです。

ワークフローを開始した EventBridge イベントを表示するには (AWS API)

- Python スクリプトに次のコードを含めます。

```
workflow_params =  
    glue_client.get_workflow_run_properties(Name=workflow_name, RunId=workflow_run_id)  
batched_events = workflow_params['aws:eventIds']
```

batched_events は、それぞれがイベント ID である文字列のリストです。

 以下も参照してください。

- [Amazon EventBridge ユーザーガイド](#)
- [the section called “ワークフローの概要”](#)

AWS Glue でのワークフローの実行とモニタリング

ワークフローの開始トリガーがオンデマンドトリガーである場合は、AWS Glue コンソールからワークフローを開始できます。ワークフローの実行とモニタリングを行うには、以下の手順を完了します。ワークフローが失敗した場合は、実行グラフを表示し、どのノードで失敗したのかを判断します。ワークフローが設計図から作成された場合は、その設計図の実行を表示して、ワークフローの作成に使用された設計図パラメータ値を確認することで、トラブルシューティングの助けになります。(詳細については、[the section called “設計図の実行の表示”](#) を参照してください)。

ワークフローを実行およびモニタリングするには、AWS Glue コンソール、API、あるいは AWS Command Line Interface (AWS CLI) を使用します

ワークフローを実行およびモニタリングするには (コンソール)

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。
2. ナビゲーションペインの [ETL] で、[ワークフロー] を選択します。
3. ワークフローを選択します。[アクション] メニューの [実行] を選択します。
4. ワークフローのリストで、[Last run status] (最終実行ステータス) 列を確認します。進行中のワークフローのステータスを表示するには、[Refresh] (更新) ボタンをクリックします。
5. ワークフローの実行中、もしくはワークフローが完了 (または失敗) した後、以下の手順を実行して実行の詳細を表示します。
 - a. ワークフローが選択されていることを確認し、[History] (履歴) タブを選択します。
 - b. 現在の、または最も新しいワークフロー実行を選択し、[View run details] (実行の詳細を表示する) をクリックします。

ワークフローの実行グラフが、現在の実行ステータスを表示します。

- c. グラフ内の任意のノードを選択して、ノードの詳細とステータスを表示します。

The screenshot displays the AWS Glue console interface. On the left, a workflow graph is shown with three nodes: a green circle labeled 'myDemoBPWorkflow1_starti...', a red square labeled 'myDemoBPWorkflow1_etl_j...', and a grey diamond labeled 'myDemoBPWorkflow1_myD...'. The red square node is highlighted with a blue border and a red 'X' icon, indicating it has failed. A 'Resume run' button is visible in the top right of the graph area. Below the graph is a legend: 'Legend: ✓ Completed, 🔄 Running, ✗ Failed, ⚠ Warning, 🚫 Error'. On the right side, the 'Job details' panel is open, showing 'Selected run' information: 'Tue, 21 Jul 2020 19:55:10 GMT - FAILED'. Below this, a table lists job details: Name (myDemoBPWorkflow1_etl_jo), Description (-), Job run id (jr_8e74182b093deea6bf63d), Status (Failed), Resume (checkbox), Retry attempt (-), Job run error (Error: Invalid argument type), Execution time (28), Start time (Tue, 21 Jul 2020 19:55:10 G), and End time (Tue, 21 Jul 2020 20:21:17 G).

ワークフローを実行およびモニタリングするには (AWS CLI)

1. 次のコマンドを入力します。<workflow-name> は、実行するワークフローに置き換えます。

```
aws glue start-workflow-run --name <workflow-name>
```

ワークフローが正常に開始されると、コマンドから実行 ID が返されます。

2. `get-workflow-run` コマンドを実行して、ワークフロー実行の状態を表示します。ワークフロー名と実行 ID を指定します。

```
aws glue get-workflow-run --name myWorkflow --run-id
wr_d2af14217e8eae775ba7b1fc6fc7a42c795aed3cbcd8763f9415452e2dbc8705
```

以下に、このコマンドの出力例を示します。

```
{
  "Run": {
    "Name": "myWorkflow",
    "WorkflowRunId":
"wr_d2af14217e8eae775ba7b1fc6fc7a42c795aed3cbcd8763f9415452e2dbc8705",
    "WorkflowRunProperties": {
      "run_state": "COMPLETED",
      "unique_id": "fee63f30-c512-4742-a9b1-7c8183bdaae2"
    },
    "StartedOn": 1578556843.049,
    "CompletedOn": 1578558649.928,
    "Status": "COMPLETED",
    "Statistics": {
      "TotalActions": 11,
      "TimeoutActions": 0,
      "FailedActions": 0,
      "StoppedActions": 0,
      "SucceededActions": 9,
      "RunningActions": 0,
      "ErroredActions": 0
    }
  }
}
```

 以下も参照してください。

- [the section called “ワークフローの概要”](#)
- [the section called “設計図の概要”](#)

ワークフロー実行の停止

ワークフロー実行を停止するには、AWS Glue コンソール、AWS Command Line Interface (AWS CLI)、または AWS Glue API を使用できます。ワークフロー実行を停止すると、すべての実行中のジョブとクローラーは即座に終了され、まだ開始されていないジョブとクローラーは開始されることがありません。すべての実行中のジョブとクローラーが停止するまでに最大 1 分かかる場合があります。ワークフロー実行のステータスは [実行中] から [停止中] に変わり、ワークフロー実行が完全に停止すると、ステータスは [停止済み] になります。

ワークフロー実行の停止後に実行グラフを表示し、完了したジョブとクローラー、および開始しなかったジョブとクローラーを確認できます。これにより、データの整合性を確保するためのステップを実行する必要があるかどうかを判断できます。ワークフロー実行を停止すると、自動ロールバックオペレーションは実行されません。

ワークフロー実行を停止するには (コンソール)

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。
2. ナビゲーションペインの [ETL] で、[ワークフロー] を選択します。
3. 実行中のワークフローを選択し、[履歴] タブを選択します。
4. ワークフロー実行を選択し、[実行の停止] を選択します。

実行ステータスが [停止中] に変わります。

5. (オプション) ワークフロー実行を選択し、[実行の詳細を表示する] を選択して、実行グラフを確認します。

ワークフロー実行を停止するには (AWS CLI)

- 次のコマンドを入力します。<workflow-name> はワークフローの名前に置き換え、<run-id> は停止するワークフロー実行の実行 ID に置き換えます。

```
aws glue stop-workflow-run --name <workflow-name> --run-id <run-id>
```

次に stop-workflow-run コマンドの例を示します。

```
aws glue stop-workflow-run --name my-workflow --run-id  
wr_137b88917411d128081069901e4a80595d97f719282094b7f271d09576770354
```

ワークフロー実行の修復と再開

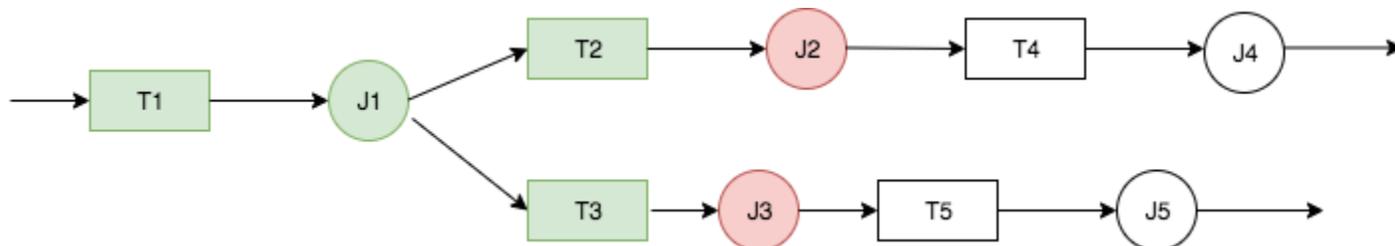
ワークフロー内の 1 つ以上のノード (ジョブまたはクローラ) が正常に完了しない場合、ワークフローは部分的にしか実行されなかったことを意味します。この根本原因を見つけて修正した上で、ワークフローの実行を再開するノードを 1 つ以上選択し、対象のワークフローを再開します。選択したノードと、それらのノードの下流にあるすべてのノードが実行されます。

トピック

- [ワークフロー実行の再開の仕組み](#)
- [ワークフロー実行を再開する](#)
- [ワークフロー実行の再開に関する注意と制限事項](#)

ワークフロー実行の再開の仕組み

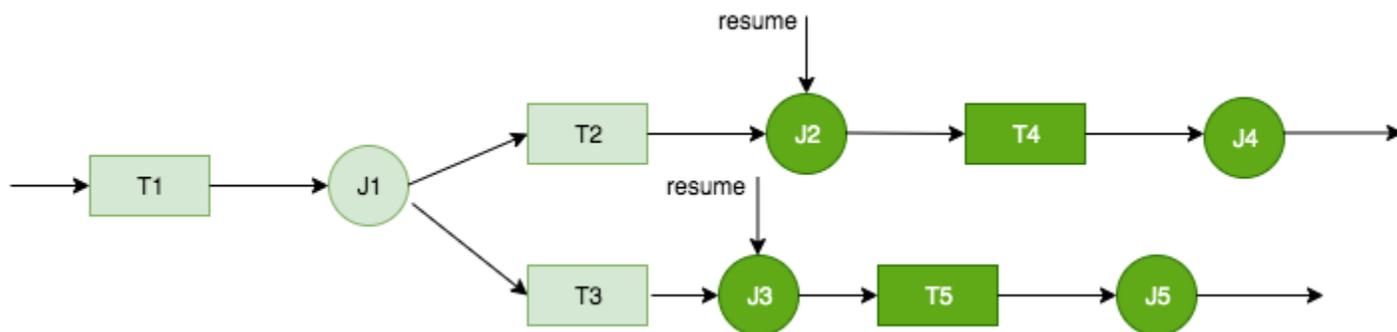
次の図にある、ワークフロー W1 について考えてみます。



ワークフロー実行は以下のように進行します。

1. トリガー T1 が、ジョブ J1 を開始します。
2. J1 が正常に完了すると、T2 と T3 がトリガーされ、それぞれジョブ J2 と J3 が実行されます。
3. ジョブ J2 と J3 が失敗します。
4. トリガー T4 と T5 は、J2 と J3 の正常な完了に依存するため、これらのトリガーは起動せずジョブ J4 と J5 も実行されません。ワークフロー W1 は部分的にのみ実行されます。

ここで、J2 と J3 の失敗の原因となった問題が修正されたとします。ワークフローの実行を再開する開始点として、J2 および J3 が選択されます。



ワークフローの実行は以下のように再開されます。

1. ジョブ J2 および J3 が正常に実行されます。
2. T4 と T5 がトリガーされます。
3. ジョブ J4 および J5 が正常に実行されます。

実行が再開されたワークフローは、新しい実行 ID を持つ別のワークフローとして追跡されます。ワークフロー履歴を表示することで、以前に行われたすべてのワークフロー実行の実行 ID を表示できます。次のスクリーンショットの例では、実行 ID `wr_c7a22...` (2 行目) を持つワークフローには、完了できなかったノードがあります。ユーザーが問題を解決し、ワークフローの実行を再開しました。その結果が、実行 ID `wr_a07e55...` (1 行目) として示されています。

Run ID	Previous run ID	Run status	Execution time
<code>wr_a07e55f2087afdd415a404403f644a4265278...</code>	<code>wr_c7a2219a8dc412f1366a5b30df3c58be30b9...</code>	Completed	17 Minutes
<code>wr_c7a2219a8dc412f1366a5b30df3c58be30b9...</code>	-	Completed	8 Minutes

Note

この説明の残りの部分では、「再開されたワークフロー実行」という表現は、以前のワークフロー実行が再開された際に作成されたワークフロー実行のことを指します。「元のワークフロー実行」とは、部分的にしか実行されず、再開する必要があったワークフロー実行を指します。

再開されたワークフロー実行のグラフ

再開されたワークフロー実行では、ノードのサブセットのみが実行されますが、実行のグラフは完全な形で表示されます。つまり、再開されたワークフローで実行されなかったノードは、元のワークフロー実行の実行グラフからコピーされています。元のワークフロー実行で実行され、グラフにコピーされたジョブおよびクローラーのノードには、実行の詳細が含まれます。

前の図のワークフロー W1 について、もう一度考えてみます。ワークフロー実行が J2 および J3 から再開されると、再開されたワークフロー実行の実行グラフには、すべてのジョブ (J1 ~ J5)、およびすべてのトリガー (T1 ~ T5) が表示されます。J1 の実行に関する詳細が、元のワークフロー実行からコピーされます。

ワークフロー実行のスナップショット

ワークフローの実行が開始されると、AWS Glue は、その時点でのワークフロー設計グラフのスナップショットを作成します。スナップショットの使用期間は、ワークフローの実行中です。実行の開始後にトリガーを変更した場合、その変更は、実行中のワークフローには影響しません。スナップショットを使用すると、ワークフローの実行を一貫した方法で行うことができます。

スナップショットでは、トリガーのみが不変になります。ワークフローの実行中に下流のジョブとクローラーに加えた変更は、その時点の実行に対しても有効になります。

ワークフロー実行を再開する

ワークフローの実行を再開するには、以下の手順に従います。ワークフロー実行は、AWS Glue コンソール、API、または AWS Command Line Interface (AWS CLI) により再開させることができます。

ワークフロー実行を再開するには (コンソール)

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。

ワークフローを表示し、その実行を再開するためのアクセス許可を持つユーザーとして、サインインします。

Note

ワークフローの実行を再開するには、`glue:ResumeWorkflowRun` AWS Identity and Access Management(IAM) のアクセス許可が必要です。

2. ナビゲーションペインで、[Workflows] (ワークフロー) をクリックします。
3. ワークフローを選択し、[History] (履歴) タブを選択します。

- 部分的にしか実行されていないワークフロー実行を選択し、[View run details] (実行の詳細を表示する) をクリックします。
- 実行グラフで、再起動しワークフロー実行の再開ポイントとする、最初の (または唯一の) ノードを選択します。
- グラフの右側にある詳細ペインで、[Resume] (再開) チェックボックスをオンにします。

The screenshot shows the AWS Glue console interface. On the left, a workflow graph is displayed with three nodes: a green circle (Completed), a red square with a clipboard icon (Failed), and a grey diamond (ALL). The failed node is highlighted with a blue border. A legend at the top indicates the status of each node. On the right, the 'Job details' panel shows the selected run as 'Tue, 21 Jul 2020 19:55:10 GMT - FAILED'. The status is 'Failed' and the error message is 'Error: Invalid argument type'. A 'Resume run' button is visible at the top right of the graph area.

ノードの色が変わり、右上に小さな [Resume] (再開) アイコンが表示されます。

The screenshot shows the AWS Glue console interface after the resume action. The failed node is now highlighted with a purple border and a small 'C' icon in the top right corner. The 'Job details' panel on the right shows the selected run as 'Tue, 21 Jul 2020 19:55:10 GMT - RESUME'. The status is 'Resume' and the 'Resume' checkbox is checked. The 'Resume run' button is still visible at the top right of the graph area.

- 再起動する追加のノードについて、前述の 2 つのステップを完了します。
- [Resume run] (実行を再開) をクリックします。

ワークフロー実行を再開するには (AWS CLI)

- glue:ResumeWorkflowRun IAM アクセス許可が付与されていることを確認します。
- 再起動するノードのノード ID を取得します。

- a. 元のワークフロー実行のために、`get-workflow-run` コマンドを実行します。ワークフロー名と実行 ID を指定し、次の例に示すように、`--include-graph` オプションを追加します。コンソールの [History] (履歴) タブから実行 ID を取得します。もしくは `get-workflow` コマンドを実行し取得します。

```
aws glue get-workflow-run --name cloudtrailtest1 --run-id
wr_a07e55f2087afdd415a404403f644a4265278f68b13ba3da08c71924e5e3c3a8 --include-
graph
```

このコマンドは、グラフのノードとエッジを大きな JSON オブジェクトとして返します。

- b. ノードオブジェクトの `Type` および `Name` プロパティごとに、対象としているノードを見つけます。

次に、出力されるノードオブジェクトの例を示します。

```
{
  "Type": "JOB",
  "Name": "test1_post_failure_4592978",
  "UniqueId":
"wnode_d1b2563c503078b153142ee76ce545fe5ceef66e053628a786ddd74a05da86fd",
  "JobDetails": {
    "JobRuns": [
      {
        "Id":
"jr_690b9f7fc5cb399204bc542c6c956f39934496a5d665a42de891e5b01f59e613",
        "Attempt": 0,
        "TriggerName": "test1_aggregate_failure_649b2432",
        "JobName": "test1_post_failure_4592978",
        "StartedOn": 1595358275.375,
        "LastModifiedOn": 1595358298.785,
        "CompletedOn": 1595358298.785,
        "JobRunState": "FAILED",
        "PredecessorRuns": [],
        "AllocatedCapacity": 0,
        "ExecutionTime": 16,
        "Timeout": 2880,
        "MaxCapacity": 0.0625,
        "LogGroupName": "/aws-glue/python-jobs"
      }
    ]
  }
}
```

```
}
```

- c. ノードオブジェクトの UniqueId プロパティから、ノード ID を取得します。
3. `resume-workflow-run` コマンドを実行します。次の例に示すように、ワークフロー名、実行 ID、およびノード ID のスペース区切りのリストを指定します。

```
aws glue resume-workflow-run --name cloudtrailtest1 --run-id
wr_a07e55f2087afdd415a404403f644a4265278f68b13ba3da08c71924ebe3c3a8 --node-
ids wnode_ca1f63e918fb855e063aed2f42ec5762ccf71b80082ae2eb5daeb8052442f2f3
wnode_d1b2563c503078b153142ee76ce545fe5ceef66e053628a786ddd74a05da86fd
```

このコマンドは、再開された (新しい) ワークフローの実行の ID と、開始されるノードのリストを出力します。

```
{
  "RunId": "wr_2ada0d3209a262fc1156e4291134b3bd643491bcfb0ceead30bd3e4efac24de9",
  "NodeIds": [
    "wnode_ca1f63e918fb855e063aed2f42ec5762ccf71b80082ae2eb5daeb8052442f2f3"
  ]
}
```

この例での `resume-workflow-run` コマンドでは再起動するノード 2 つをリストしていますが、出力例では 1 つのノードのみが再起動されると表示していることに注意してください。これは、1 つのノードが他のノードの下流にあり、いずれにせよ、この下流のノードはワークフローの通常のフローによって再開されるためです。

ワークフロー実行の再開に関する注意と制限事項

ワークフロー実行を再開する場合には、次の注意事項と制限事項に留意してください。

- 実行を再開できるのは、COMPLETED 状態のワークフローのみです。

Note

ワークフロー実行の中で 1 つ以上のノードが未完了の場合でも、ワークフロー実行の状態は COMPLETED となります。実行グラフをチェックして、正常に完了しなかったノードを見つける必要があります。

- 元のワークフロー実行により実行が試みられた任意のジョブまたはクローラのノードから、ワークフロー実行を再開することができます。トリガーノードからは、ワークフロー実行を再開することはできません。
- ノードを再起動しても、その状態はリセットされません。処理が部分的に完了しているデータのロールバックは行われません。
- 同じワークフロー実行を、複数回再開することが可能です。再開されたワークフローの実行が部分的にしか実行されない場合は、その問題に対処した上で実行を再開します。
- 再起動のために2つのノードを選択し、それらが互いに依存している場合には、上流にあるノードが下流ノードより先に実行されます。実際、下流にあるノードはワークフローの通常のフローによって実行されるため、このノードを選択することは冗長的となります。

AWS Glue でのワークフローの実行プロパティの取得と設定

ワークフローの実行プロパティを使用して、AWS Glue ワークフローのジョブ間の状態を共有して管理します。デフォルトの実行プロパティは、ワークフローの作成時に設定できます。次に、ジョブの実行時に、ジョブは実行プロパティの値を取得して必要に応じて変更し、ワークフローの以降のジョブに対する入力として使用できます。ジョブが実行プロパティを変更した場合、新しい値が有効なのはワークフローの実行中に限られます。デフォルトの実行プロパティは影響を受けません。

AWS Glue ジョブがワークフローに含まれていない場合、これらのプロパティは設定されません。

次の Python コード例は抽出、変換、ロード (ETL) ジョブの一部で、ワークフローの実行プロパティを取得する方法を示しています。

```
import sys
import boto3
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from awsglue.context import GlueContext
from pyspark.context import SparkContext

glue_client = boto3.client("glue")
args = getResolvedOptions(sys.argv, ['JOB_NAME', 'WORKFLOW_NAME', 'WORKFLOW_RUN_ID'])
workflow_name = args['WORKFLOW_NAME']
workflow_run_id = args['WORKFLOW_RUN_ID']
workflow_params = glue_client.get_workflow_run_properties(Name=workflow_name,
                                                         RunId=workflow_run_id)["RunProperties"]

target_database = workflow_params['target_database']
```

```
target_s3_location = workflow_params['target_s3_location']
```

次のコードは、`target_format` 実行プロパティを 'csv' に設定することで続行します。

```
workflow_params['target_format'] = 'csv'  
glue_client.put_workflow_run_properties(Name=workflow_name, RunId=workflow_run_id,  
RunProperties=workflow_params)
```

詳細については、次を参照してください:

- [GetWorkflowRunProperties アクション \(Python: get_workflow_run_properties\)](#)
- [PutWorkflowRunProperties アクション \(Python: put_workflow_run_properties\)](#)

AWS Glue APIを使用したワークフローのクエリ

AWS Glue には、ワークフローを管理するための高度な API が用意されています。AWS Glue API を使用してワークフローの静的ビューまたは実行中のワークフローの動的ビューを取得できます。詳細については、[ワークフロー](#)を参照してください。

トピック

- [静的ビューのクエリ](#)
- [動的ビューのクエリ](#)

静的ビューのクエリ

GetWorkflow API オペレーションを使用して、ワークフローの設計を示す静的ビューを取得します。このオペレーションは、ノードおよびエッジで構成される有向グラフを返します。ノードは、トリガー、ジョブ、またはクローラを表します。エッジは、ノード間の関係を定義します。AWS Glue コンソールのグラフでは、エッジがコネクタ (矢印) で表されます。

このオペレーションは、NetworkX、igraph、JGraphT、Java Universal Network/Graph (JUNG) Framework などの一般的なグラフ処理ライブラリでも使用できます。これらのライブラリが表すグラフはすべてが類似しているため、必要な変換は最小限です。

この API から返される静的ビューは、ワークフローに関連付けられているトリガーの最新の定義に基づく最新のビューです。

グラフ定義

ワークフローのグラフ G は順序付きペア (N, E) で、 N はノードのセットを表し、 E はエッジのセットを表します。ノードはグラフ内の頂点で、一意の数値で識別されます。ノードのタイプは、トリガー、ジョブ、またはクローラです。例: $\{name:T1, type:Trigger, uniqueId:1\}$, $\{name:J1, type:Job, uniqueId:2\}$ 。

エッジは、2 タプルの形式 $(src, dest)$ です。 src と $dest$ はノードで、 src から $dest$ への有向エッジがあります。

静的ビューのクエリの例を示します。

ジョブ $J1$ の完了時にジョブ $J2$ をトリガーする条件付きのトリガー T を考えます。

```
J1 ----> T ----> J2
```

ノード: $J1$ 、 T 、 $J2$

エッジ: $(J1, T)$ 、 $(T, J2)$

動的ビューのクエリ

`GetWorkflowRun` API オペレーションを使用して、実行中のワークフローの動的ビューを取得します。このオペレーションは、グラフの同じ静的ビューを、ワークフローの実行に関連するメタデータと共に返します。

実行の際、`GetWorkflowRun` 呼び出しのジョブを表すノードでは、最新のワークフロー実行の一部として開始されたジョブ実行のリストを含んでいます。このリストを使用して、各ジョブの実行ステータスをグラフ自体に表示できます。まだ実行されていないダウンストリームの依存関係の場合、このフィールドは `null` に設定されます。グラフ化された情報により、任意の時点における任意のワークフローの現在の状態を確認できます。

この API から返される動的ビューは、ワークフローの実行が開始された時点における静的ビューに基づいています。

ランタイムノードの例: $\{name:T1, type: Trigger, uniqueId:1\}$ 、 $\{name:J1, type:Job, uniqueId:2, jobDetails:\{jobRuns\}\}$ 、 $\{name:C1, type:Crawler, uniqueId:3, crawlerDetails:\{crawls\}\}$

例 1: 動的ビュー

次の例は、単純な 2 トリガーワークフローを示しています。

- ノード: t1、j1、t2、j2
- エッジ: (t1, j1)、(j1, t2)、(t2, j2)

GetWorkflow レスポンスの内容は次のとおりです。

```
{
  Nodes : [
    {
      "type" : Trigger,
      "name" : "t1",
      "uniqueId" : 1
    },
    {
      "type" : Job,
      "name" : "j1",
      "uniqueId" : 2
    },
    {
      "type" : Trigger,
      "name" : "t2",
      "uniqueId" : 3
    },
    {
      "type" : Job,
      "name" : "j2",
      "uniqueId" : 4
    }
  ],
  Edges : [
    {
      "sourceId" : 1,
      "destinationId" : 2
    },
    {
      "sourceId" : 2,
      "destinationId" : 3
    },
    {
      "sourceId" : 3,
      "destinationId" : 4
    }
  ]
}
```

GetWorkflowRun レスポンスの内容は次のとおりです。

```
{
  Nodes : [
    {
      "type" : Trigger,
      "name" : "t1",
      "uniqueId" : 1,
      "jobDetails" : null,
      "crawlerDetails" : null
    },
    {
      "type" : Job,
      "name" : "j1",
      "uniqueId" : 2,
      "jobDetails" : [
        {
          "id" : "jr_12334",
          "jobRunState" : "SUCCEEDED",
          "errorMessage" : "error string"
        }
      ],
      "crawlerDetails" : null
    },
    {
      "type" : Trigger,
      "name" : "t2",
      "uniqueId" : 3,
      "jobDetails" : null,
      "crawlerDetails" : null
    },
    {
      "type" : Job,
      "name" : "j2",
      "uniqueId" : 4,
      "jobDetails" : [
        {
          "id" : "jr_1233sdf4",
          "jobRunState" : "SUCCEEDED",
          "errorMessage" : "error string"
        }
      ],
      "crawlerDetails" : null
    }
  ]
}
```

```

    ],
    Edges : [
      {
        "sourceId" : 1,
        "destinationId" : 2
      },
      {
        "sourceId" : 2,
        "destinationId" : 3
      },
      {
        "sourceId" : 3,
        "destinationId" : 4
      }
    ]
  }

```

例 2: 複数のジョブと 1 つの条件付きのトリガー

次の例は、複数のジョブと 1 つの条件付きのトリガー (t3) を持つワークフローを示しています。

Consider Flow:

```

T(t1) ----> J(j1) ----> T(t2) ----> J(j2)
      |                               |
      |                               |
      >+-----> T(t3) <-----+
              |
              |
              J(j3)

```

Graph generated:

Nodes: t1, t2, t3, j1, j2, j3

Edges: (t1, j1), (j1, t2), (t2, j2), (j1, t3), (j2, t3), (t3, j3)

AWS Glue での設計図とワークフローの制限

設計図とワークフローについての制限を次に示します。

設計図での制限

設計図については、以下の制限に注意してください。

- 設計図は、Amazon S3 バケットが置かれているのと同じ AWS リージョンに登録される必要があります。

- 設計図を AWS アカウント間で共有する場合は、Amazon S3 内で設計図の ZIP アーカイブに対する読み取りを許可する必要があります。設計図の ZIP アーカイブに対する読み取り許可を持つユーザーは、自分の AWS アカウントに設計図を登録し、それを使用することができます。
- 設計図のパラメータセットは、単一の JSON オブジェクトとして格納されます。このオブジェクトの最大長は 128 KB です。
- 設計図の ZIP アーカイブにおける非圧縮状態での最大サイズは 5 MB です。圧縮状態での最大サイズは 1 MB です。
- ワークフロー内のジョブ、クローラ、トリガーの総数を、100 以下に制限します。100 を超える値を含めると、ワークフローの実行を再開または停止しようとした場合にエラーが発生することがあります。

ワークフローの制限

ワークフローについては、以下の制限に注意してください。これらのコメントの中には、ワークフローを手動で作成しているユーザーに直接関係するものがあります。

- Amazon EventBridge のイベントトリガーの最大バッチサイズは 100 です。最大ウィンドウサイズは 900 秒 (15 分) です。
- トリガーは 1 つのワークフローにのみ関連付けることができます。
- 1 つの開始トリガー (オンデマンドまたはスケジュール) のみが許可されます。
- ワークフロー内のジョブまたはクローラがワークフロー外のトリガーによって開始された場合、ジョブまたはクローラの完了 (成功またはその他) に依存するワークフロー内のトリガーは起動されません。
- 同様に、ワークフロー内のジョブまたはクローラに、ワークフロー内とワークフロー外の両方でジョブまたはクローラの完了 (成功またはそれ以外) に依存するトリガーがあり、さらにジョブまたはクローラがワークフロー内から開始されている場合には、ジョブまたはクローラの完了時、ワークフロー内のトリガーのみが起動されます。

AWS Glue での設計図エラーをトラブルシューティングする

AWS Glue の設計図を使用する際にエラーが発生した場合は、次の解決策を参照しながら問題の原因を突き止め、それを修正してください。

トピック

- [エラー: PySpark モジュールが見つからない](#)

- [エラー: 設計図の設定ファイルが見つからない](#)
- [エラー:インポートされたファイルが見つからない](#)
- [エラー:リソースでの iamPassRole の実行が承認されない](#)
- [エラー: cron スケジュールが無効](#)
- [エラー: 既存のものと同じ名前を使用してのトリガー](#)
- [エラー:名前:foo という名前のワークフローが既に存在する](#)
- [エラー: 指定された layoutGenerator のパスにモジュールが見つからない](#)
- [エラー: Connections フィールドの検証エラー](#)

エラー: PySpark モジュールが見つからない

AWS Glue が、「Unknown error executing layout generator function ModuleNotFoundError: No module named 'pyspark'」というエラーを返します。

設計図アーカイブを解凍すると、以下のいずれかのようにになります。

```
$ unzip compaction.zip
Archive:  compaction.zip
  creating:  compaction/
  inflating:  compaction/blueprint.cfg
  inflating:  compaction/layout.py
  inflating:  compaction/README.md
  inflating:  compaction/compaction.py

$ unzip compaction.zip
Archive:  compaction.zip
  inflating:  blueprint.cfg
  inflating:  compaction.py
  inflating:  layout.py
  inflating:  README.md
```

1 つ目のケースでは、設計図に関連するすべてのファイルが compaction という名前のフォルダの下に配置され、その後で compaction.zip という名前の zip ファイルに変換されています。

2 つ目のケースでは、設計図に必要なすべてのファイルをフォルダに配置せずに、zip ファイル compaction.zip に対しルートファイルとして追加しています。

上記の形式のどちらを使用しても、ファイルを作成することができます。ただし、blueprint.cfg の中には、レイアウトを生成するスクリプト内の関数名を指す正しいパスが必要です。

例

ケース 1 の場合、`blueprint.cfg` には以下のように `layoutGenerator` が必要です

```
layoutGenerator": "compaction.layout.generate_layout"
```

ケース 2 の場合、`blueprint.cfg` には以下のように `layoutGenerator` が必要です

```
layoutGenerator": "layout.generate_layout"
```

このパスが正しく指定されていない場合は、次のようなエラーが表示されることがあります。例えば、フォルダをケース 2 で説明した構造としながら、ケース 1 で示した `layoutGenerator` が使用されている場合、上記のようなエラーが発生する場合があります。

エラー: 設計図の設定ファイルが見つからない

AWS Glue が、「Unknown error executing layout generator function FileNotFoundError: [Errno 2] No such file or directory: '/tmp/compaction/blueprint.cfg」というエラーを返します。

`blueprint.cfg` ファイルは、ZIP アーカイブのルートレベルか、ZIP アーカイブと同じ名前のフォルダ内に配置する必要があります。

設計図の ZIP アーカイブを抽出すると、通常、`blueprint.cfg` は、下記のいずれかのパスで見つかります。以下のパスのいずれかに配置されていない場合には、上記のエラーが表示されます。

```
$ unzip compaction.zip
Archive:  compaction.zip
   creating:  compaction/
  inflating:  compaction/blueprint.cfg

$ unzip compaction.zip
Archive:  compaction.zip
  inflating:  blueprint.cfg
```

エラー:インポートされたファイルが見つからない

AWS Glue が、「Unknown error executing layout generator function FileNotFoundError: [Errno 2] No such file or directory: '* *demo-project/foo.py」というエラーを返します。

レイアウト生成スクリプトに他のファイルを読み込む機能がある場合は、インポートするファイルの完全なパスを指定する必要があります。例えば、Conversion.py スクリプトは Layout.py から参照されることがあります。詳細については、「[Sample Blueprint Project](#)」を参照してください。

エラー: リソースでの iamPassRole の実行が承認されない

AWS Glue が、「User: arn:aws:sts::123456789012:assumed-role/AWSGlueServiceRole/ GlueSession is not authorized to perform: iam:PassRole on resource: arn:aws:iam::123456789012:role/AWSGlueServiceRole」というエラーを返します。

ワークフロー内のジョブとクローラが、設計図からワークフローを作成するために渡されたロールと同じロールを引き受ける場合、設計図のロールには、それ自体に対する iam:PassRole アクセス許可を含める必要があります。

ワークフロー内のジョブとクローラが、設計図からワークフローのエンティティを作成するために渡されたロール以外のロールを引き受ける場合、設計図のロールには、それ自体のロールではなく引き受けたロールに対する、iam:PassRole アクセス許可を含める必要があります。

詳細については、「[Permissions for Blueprint Roles](#)」を参照してください。

エラー: cron スケジュールが無効

AWS Glue が、「The schedule cron(0 0 * * * *) is invalid.」というエラーを返します。

有効な [cron](#) 式を指定します。詳細については、「[Time-Based Schedules for Jobs and Crawlers](#)」を参照してください。

エラー: 既存のものと同じ名前を使用してのトリガー

AWS Glue が、「Trigger with name 'foo_starting_trigger' already submitted with different configuration」というエラーを返します。

設計図でワークフローを作成するために、レイアウトスクリプトの中でトリガーを定義する必要はありません。トリガーの作成は、2つのアクション間に定義された依存関係に基づいて、設計図ライブラリによって管理されます。

トリガーの名前は、次のように決定されます。

- ワークフローの開始トリガーの場合は、<workflow_name>_starting_trigger のような名前になります。

- ワークフロー内のノード (ジョブ/クローラ) で、1 つまたは複数の上流ノードの完了に依存している場合、AWS Glue はトリガーを、<workflow_name>_<node_name>_trigger のような名前で定義します

このエラーは、既に存在するものと同じ名前で、トリガーが発生したことを意味します。既存のトリガーを削除すれば、ワークフローの作成を再実行できます。

Note

ワークフローを削除しても、ワークフロー内のノードは削除されません。ワークフローが削除されても、トリガーが残されている場合があります。このため、作成したワークフローを削除した後に、同じ名前で同じ設計図から再作成しようとする、ワークフローが既に存在していることを知らせるエラーは表示されずに、トリガーが既に存在している、というエラーが出力されることがあります。

エラー:名前:foo という名前のワークフローが既に存在する

ワークフロー名は一意である必要があります。別の名前で再試行してください。

エラー: 指定された layoutGenerator のパスにモジュールが見つからない

AWS Glue が、「Unknown error executing layout generator function ModuleNotFoundError: No module named 'crawl_s3_locations'」というエラーを返します。

```
layoutGenerator": "crawl_s3_locations.layout.generate_layout"
```

例えば、上記の LayoutGenerator パスを指定しながら設計図アーカイブを解凍する際には、次のようにする必要があります。

```
$ unzip crawl_s3_locations.zip
Archive:  crawl_s3_locations.zip
  creating:  crawl_s3_locations/
  inflating:  crawl_s3_locations/blueprint.cfg
  inflating:  crawl_s3_locations/layout.py
  inflating:  crawl_s3_locations/README.md
```

アーカイブの解凍時に、設計図アーカイブが次のようになっている場合は、上記のエラーが発生することがあります。

```
$ unzip crawl_s3_locations.zip
Archive:  crawl_s3_locations.zip
  inflating: blueprint.cfg
  inflating: layout.py
  inflating: README.md
```

`crawl_s3_locations` という名前のフォルダがなく、`layoutGenerator` パスがモジュール `crawl_s3_locations` を介してレイアウトファイルを参照している場合に、上記のエラーが発生する可能性があります。

エラー: Connections フィールドの検証エラー

AWS Glue が、「Unknown error executing layout generator function TypeError: Value ['foo'] for key Connections should be of type <class 'dict'>!」というエラーを返します。

これは検証に関するエラーです。Job クラス内の Connections フィールドはディクショナリの使用を想定しており、代わりに値のリストが提供された場合には、エラーの原因となります。

```
User input was list of values
Connections= ['string']

Should be a dict like the following
Connections*={'Connections': ['string']}
```

設計図からワークフローを作成する際に発生する、これらの実行時エラーを回避するには、「[Testing a Blueprint](#)」に概略されているとおりに、ワークフロー、ジョブ、およびクローラの定義を検証します。

レイアウトスクリプト内での AWS Glue ジョブ、クローラ、ワークフローの定義に関しては、「[AWS Glue Blueprint Classes Reference](#)」にある構文を参照してください。

AWS Glue ブループリントでのペルソナおよびロール用のアクセス許可

以下は、AWS Identity and Access Management ブループリントのための典型的なペルソナと、ペルソナおよびロール用に提案されている AWS Glue(IAM) アクセス許可ポリシーです。

トピック

- [設計図のペルソナ](#)
- [設計図のペルソナ用のアクセス許可](#)
- [設計図のロール用のアクセス許可](#)

設計図のペルソナ

以下は、一般的に、AWS Glue ブループリントのライフサイクルに關与するペルソナです

ペルソナ	説明
AWS Glue 開発者	設計図を開発、テスト、公開します。
AWS Glue 管理者	設計図を登録および保守し、アクセス許可を付与します。
データアナリスト	設計図を実行してワークフローを作成します。

(詳細については、[the section called “設計図の概要”](#) を参照してください)。

設計図のペルソナ用のアクセス許可

設計図のペルソナごとに、提案されているアクセス許可を以下に示します。

ブループリントのための AWS Glue デベロッパー権限

AWS Glue デベロッパーには設計図の公開に使用される Amazon S3 バケットに対する、書き込みのアクセス許可が必要です。多くの場合、デベロッパーはアップロードした後に設計図を登録します。その場合、デベロッパーには [the section called “ブループリントのための AWS Glue 管理者のアクセス許可”](#) に一覧表示されているアクセス許可が必要です。さらに、デベロッパーが登録後に設計図をテストする場合は、[the section called “設計図のためのデータアナリストの権限”](#) に一覧表示されているアクセス許可も必要となります。

ブループリントのための AWS Glue 管理者のアクセス許可

次のポリシーでは、AWS Glue ブループリントを登録、表示、保守するためのアクセス許可が付与されます。

Important

次のポリシーでは、`<s3-bucket-name>` および `<prefix>` を、登録のためにアップロードされた設計図の ZIP アーカイブを指す Amazon S3 パスに置き換えます。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "glue:CreateBlueprint",
        "glue:UpdateBlueprint",
        "glue>DeleteBlueprint",
        "glue:GetBlueprint",
        "glue:ListBlueprints",
        "glue:BatchGetBlueprints"
      ],
      "Resource": "*"
    },
    {
      "Effect": "Allow",
      "Action": [
        "s3:GetObject"
      ],
      "Resource": "arn:aws:s3:::<s3-bucket-name>/<prefix>/*"
    }
  ]
}
```

設計図のためのデータアナリストの権限

次のポリシーは、設計図を実行し、結果のワークフローとワークフローコンポーネントを表示するアクセス許可を付与します。また、ワークフローおよびワークフローコンポーネントを作成するために PassRole が引き受ける、AWS Glue というロールを付与します。

このポリシーでは、任意のリソースに対するアクセス許可が付与されます。個々の設計図へのきめ細かなアクセスを構成する場合は、設計図の ARN に次の形式を使用します。

```
arn:aws:glue:<region>:<account-id>:blueprint/<blueprint-name>
```

Important

次のポリシーでは、<account-id> を有効な AWS アカウントに置き換え、<role-name> を設計図の実行に使用するロールの名前に置き換えます。このロールに必要なアクセス許

可については、「[the section called “設計図のロール用のアクセス許可”](#)」を参照してください。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "glue:ListBlueprints",
        "glue:GetBlueprint",
        "glue:StartBlueprintRun",
        "glue:GetBlueprintRun",
        "glue:GetBlueprintRuns",
        "glue:GetCrawler",
        "glue:ListTriggers",
        "glue:ListJobs",
        "glue:BatchGetCrawlers",
        "glue:GetTrigger",
        "glue:BatchGetWorkflows",
        "glue:BatchGetTriggers",
        "glue:BatchGetJobs",
        "glue:BatchGetBlueprints",
        "glue:GetWorkflowRun",
        "glue:GetWorkflowRuns",
        "glue:ListCrawlers",
        "glue:ListWorkflows",
        "glue:GetJob",
        "glue:GetWorkflow",
        "glue:StartWorkflowRun"
      ],
      "Resource": "*"
    },
    {
      "Effect": "Allow",
      "Action": "iam:PassRole",
      "Resource": "arn:aws:iam::<account-id>:role/<role-name>"
    }
  ]
}
```

設計図のロール用のアクセス許可

設計図からワークフローを作成するために、IAM ロール用として提案されているアクセス許可を次に示します。このロールには `glue.amazonaws.com` との信頼関係が必要です。

⚠ Important

次のポリシーでは、`<account-id>` を有効な AWS アカウントに置き換え、`<role-name>` をロールの名前に置き換えます。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "glue:CreateJob",
        "glue:GetCrawler",
        "glue:GetTrigger",
        "glue>DeleteCrawler",
        "glue:CreateTrigger",
        "glue>DeleteTrigger",
        "glue>DeleteJob",
        "glue:CreateWorkflow",
        "glue>DeleteWorkflow",
        "glue:GetJob",
        "glue:GetWorkflow",
        "glue:CreateCrawler"
      ],
      "Resource": "*"
    },
    {
      "Effect": "Allow",
      "Action": "iam:PassRole",
      "Resource": "arn:aws:iam::<account-id>:role/<role-name>"
    }
  ]
}
```

Note

ワークフロー内のジョブとクローラが、このロール以外のロールを引き受ける場合、このポリシーには、設計図のロールのものではなく引き受けたロールでの `iam:PassRole` アクセス許可を付与する必要があります。

AWS Glue のブループリントの開発

組織が、類似した一連の ETL ユースケースを実行する場合、それらすべてをパラメータ化した単一のワークフローで処理することでメリットが得られます。このようなニーズに応えるため、AWS Glue では、設計図を定義して、ワークフローの生成に使用できるようになってます。設計図にはパラメータを渡すことができるので、データアナリストは単一の設計図からさまざまなワークフローを作成して、類似の ETL ユースケースを処理できます。作成した設計図は、別の部門、チーム、プロジェクトでの再利用が可能です。

トピック

- [AWS Glue のブループリント概要](#)
- [AWS Glue のブループリントの開発](#)
- [AWS Glue でブループリントを登録する](#)
- [AWS Glue でブループリントを表示する](#)
- [AWS Glue でブループリントを更新する](#)
- [AWS Glue で設計図からワークフローを作成する](#)
- [AWS Glue でブループリントの実行を表示する](#)

AWS Glue のブループリント概要

Note

ブループリント機能は、現在、AWS Glue コンソールの次のリージョンではご利用いただけません。アジアパシフィック (ジャカルタ)、中東 (UAE)。

AWS Glue ブループリントにより、AWS Glue ワークフローを作成ならびに共有する手段が提供されます。複雑な ETL プロセスで、使用できるユースケースが類似している場合には、ユースケースごとに AWS Glue ワークフローを作成する代わりにブループリントを 1 つ作成します。

設計図では、ワークフローに含めるジョブとクローラを指定し、ワークフローを作成するために設計図が実行される際にワークフローのユーザーから提供すべきパラメータも指定します。パラメータを使用することで、単一のブループリントにより、類似したさまざまなユースケースのワークフローを生成できるようになります。ワークフローの詳細については、「[AWS Glue のワークフローの概要](#)」を参照してください。

設計図に関するユースケースの例を次に示します。

- 既存のデータセットを分割する場合。Amazon Simple Storage Service (Amazon S3) のソースならびにターゲットへのパス、およびパーティション列のリストが、設計図への入力パラメータとなります。
- Amazon DynamoDB テーブルのスナップショットを、Amazon Redshift のような SQL データストア内に作成する場合。ブループリントへの入力パラメータは、DynamoDB テーブル名と AWS Glue の接続となります。このパラメータにより、Amazon Redshift クラスターと送信先のデータベースを指定します。
- 複数の Amazon S3 パスの CSV データを Parquet に変換する場合。AWS Glue ワークフローに、パスごとの個別のクローラとジョブを含める必要があります。AWS Glue Data Catalog 内の送信先データベースと、Amazon S3 パスのコンマ区切りのリストが入力パラメータとなります。この場合、ワークフローが作成するクローラおよびジョブの数は可変であることに注意してください。

設計図のコンポーネント

設計図は、以下のコンポーネントを含む ZIP アーカイブです。

- Python レイアウトジェネレータースクリプト

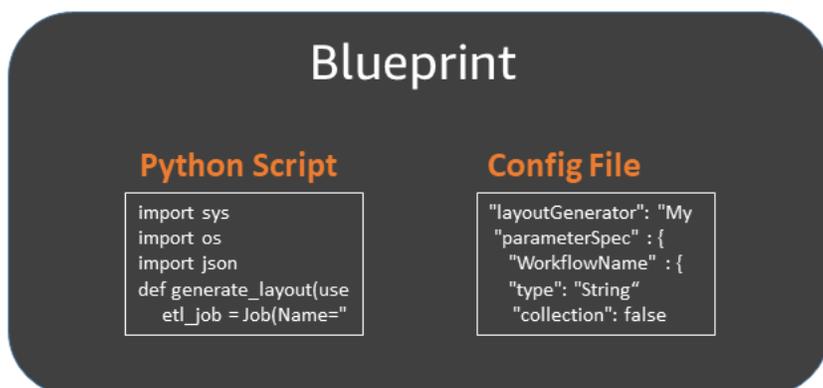
ワークフローのレイアウト (ワークフロー用に作成するジョブとクローラ、ジョブとクローラのプロパティ、およびジョブとクローラの依存関係) を指定する関数が含まれています。この関数はブループリントのパラメータを受け取り、AWS Glue がワークフロー生成のために使用するワークフロー構造 (JSON オブジェクト) を返します。ワークフローの生成には Python スクリプトを使用するので、ユースケースに適した独自のロジックを追加できます。

- 設定ファイル

ワークフローのレイアウトを生成する Python 関数の完全な修飾名を指定します。スクリプトで使用されるすべての設計図パラメータの名前、データタイプ、およびその他のプロパティも指定します。

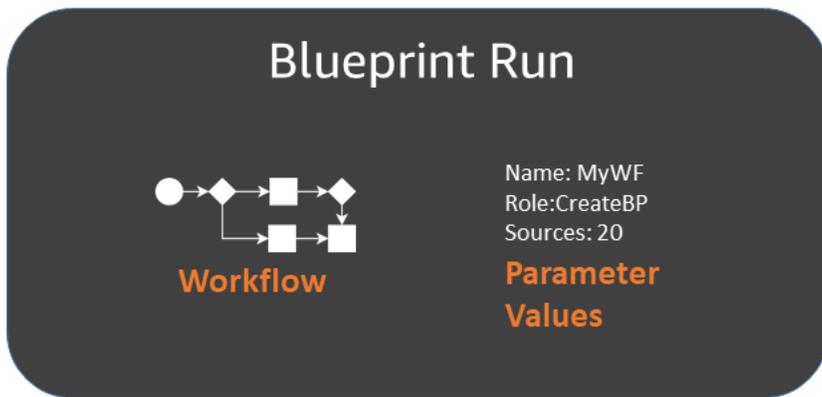
- (オプション) ETL スクリプトとサポートファイル

高度なユースケースとして、ジョブで使用する ETL スクリプトの場所をパラメータ化できます。ZIP アーカイブにジョブスクリプトファイルを含め、スクリプトのコピー先となる Amazon S3 ロケーションを設計図パラメータとして指定できます。レイアウトジェネレータスクリプトは、指定された場所に ETL スクリプトをコピーし、そのコピー先をジョブスクリプトの場所プロパティとして指定します。スクリプトで処理するための、ライブラリやその他のサポートファイルを含めることもできます。



設計図の実行

ブループリントからワークフローを作成すると、AWS Glue でブループリントが実行され、ワークフローがカプセル化するジョブ、クローラー、およびトリガーを作成する非同期プロセスが開始します。AWS Glue でブループリントの実行を使用して、ワークフローとそのコンポーネントの作成をオーケストレーションします。設計図の実行ステータスを表示すると、作成プロセスのステータスを知ることができます。設計図の実行では、設計図パラメータに指定した値の保存も行われます。



ブループリントの実行は、AWS Glue コンソールまたは AWS Command Line Interface (AWS CLI) を使用して表示できます。ワークフローを表示またはトラブルシューティングする場合には、いつでも設計図の実行に戻り、ワークフローの作成に使用された設計図のパラメータ値を確認できます。

設計図のライフサイクル

ブループリントは、AWS Glue により開発、テスト、登録が行われ、これを実行することでワークフローが作成されます。設計図のライフサイクルには、通常 3 つのペルソナが関与します。

ペルソナ	タスク
AWS Glue デベロッパー	<ul style="list-style-type: none"> ワークフローレイアウトのスクリプトを記述し、設定ファイルを作成します。 AWS Glue によって提供されるライブラリを使用して、ブループリントをローカルでテストします。 スクリプト、設定ファイル、およびサポートファイルを含む ZIP アーカイブを作成し、Amazon S3 内のロケーションでそのアーカイブを公開します。 バケットオブジェクトに対する読み取りアクセス許可を、AWS Glue 管理者の AWS アカウントに付与するバケットポリシーを、Amazon S3 バケットに追加します。 Amazon S3 にある ZIP アーカイブに対する、読み取りのための IAM アクセス許可を、AWS Glue 管理者に付与します。

ペルソナ	タスク
AWS Glue 管理者	<ul style="list-style-type: none">• AWS Glue にブループリントを登録します。AWS Glue で予約した Amazon S3 の場所に、ZIP アーカイブのコピーを作成します。• 設計図に対する IAM アクセス許可をデータアナリストに付与します。
データアナリスト	<ul style="list-style-type: none">• 設計図を実行したワークフローの作成と、設計図パラメータ値の指定を行います。設計図の実行ステータスをチェックして、ワークフローとワークフローコンポーネントが正常に生成されたことを確認します。• ワークフローを実行し、トラブルシューティングします。ワークフローを実行する前に、ワークフローのデザイングラフを AWS Glue コンソールに表示し、そのワークフローを検証します。

 以下も参照してください。

- [AWS Glue のブループリントの開発](#)
- [AWS Glue で設計図からワークフローを作成する](#)
- [AWS Glue ブループリントでのペルソナおよびロール用のアクセス許可](#)

AWS Glue のブループリントの開発

AWS Glue デベロッパーは、データアナリストがワークフローを生成するための設計図を作成および公開できます。

トピック

- [設計図の開発の概要](#)
- [設計図の開発のための前提条件](#)
- [設計図コードを記述する](#)
- [設計図プロジェクト例](#)

- [設計図をテストする](#)
- [設計図を公開する](#)
- [AWS Glue ブループリントクラスリファレンス](#)
- [設計図の例](#)

i 以下も参照してください。

- [AWS Glue のブループリント概要](#)

設計図の開発の概要

開発プロセスの最初のステップは、設計図からのメリットが活かせるような、共通点のあるユースケースを特定することです。典型的なユースケースとしては、一般化した方法での解決が想定される、反復的な ETL の問題が考えられます。次に、その一般化されたユースケースを実装するために設計図を構成します。一般化されたユースケースを使用して特定のユースケースを定義できるように、設計図の入力パラメータを定義します。

設計図は、その設計図のパラメータ構成ファイルを含むプロジェクトと、生成するワークフローのレイアウトを定義するスクリプトにより構成されます。レイアウトは、作成するジョブとクローラ (または設計図の用語におけるエンティティ) を定義します。

レイアウトスクリプトでは、トリガーを直接指定しません。その代わりに、スクリプトが作成するジョブやクローラー間の依存関係を指定するコードを書き込みます。AWS Glue は依存関係の指定に基づき、トリガーを生成します。レイアウトスクリプトからは、すべてのワークフローエンティティの仕様を含む、ワークフローオブジェクトが出力されます。

ワークフローオブジェクトの構築には、以下に示す AWS Glue のブループリントライブラリを使用します。

- `aws glue . blueprint . base_resource` – ライブラリで使用される基本リソースのライブラリ。
- `aws glue . blueprint . workflow` – Workflow クラスを定義するためのライブラリ。
- `aws glue . blueprint . job` – Job クラスを定義するためのライブラリ。
- `aws glue . blueprint . crawler` – Crawler クラスを定義するためのライブラリ。

上記の他には、Python シェルで使用できるライブラリのみが、レイアウト生成用のライブラリとしてサポートされています。

設計図ライブラリで定義されたメソッドを使用することで、公開前の設計図をローカルでテストできます。

設計図をデータアナリストに提供する準備ができたなら、スクリプト、パラメータ設定ファイル、および追加のスクリプトやライブラリなどのサポートファイルを、単一のデプロイ可能なアセットの中にパッケージ化します。次に、このアセットを Amazon S3 にアップロードし、AWS Glue への登録を管理者に依頼します。

その他の設計図プロジェクトの例については、「[設計図プロジェクト例](#)」および「[設計図の例](#)」を参照してください。

設計図の開発のための前提条件

ブループリントを開発するには、AWS Glue と Apache Spark ETL ジョブや Python シェルジョブのスクリプト記述に関する知識が必要です。また、以下のセットアップタスクも完了する必要があります。

- 設計図レイアウトスクリプトで使用する 4 個の AWS Python ライブラリのダウンロード
- AWS SDK のセットアップ。
- AWS CLI のセットアップ。

Python ライブラリをダウンロードする

以下のライブラリを GitHub からダウンロードし、自分のプロジェクトにインストールします。

- https://github.com/aws-labs/aws-glue-blueprint-libs/tree/master/awsglue/blueprint/base_resource.py
- <https://github.com/aws-labs/aws-glue-blueprint-libs/tree/master/awsglue/blueprint/workflow.py>
- <https://github.com/aws-labs/aws-glue-blueprint-libs/tree/master/awsglue/blueprint/crawler.py>
- <https://github.com/aws-labs/aws-glue-blueprint-libs/tree/master/awsglue/blueprint/job.py>

AWS Java SDK のセットアップ

AWS Java SDK に対しては、設計図の API を含む jar ファイルを追加する必要があります。

1. まだ追加していない場合には、AWS SDK for Java のセットアップを行います。

- Java 1.x の場合は、「AWS SDK for Java デベロッパーガイド」の「[AWS SDK for Java の入手方法](#)」にある手順に従ってください。
 - Java 2.x の場合は、「AWS SDK for Java 2.x デベロッパーガイド」の「[AWS SDK for Java 2.x のセットアップ](#)」にある手順に従ってください。
2. 設計図の API へのアクセス許可が付与されている、クライアント jar ファイルをダウンロードします。
 - Java 1.x 用ファイル: `s3://awsglue-custom-blueprints-preview-artifacts/awsglue-java-sdk-preview/AWSGlueJavaClient-1.11.x.jar`
 - Java 2.x 用ファイル: `s3://awsglue-custom-blueprints-preview-artifacts/awsglue-java-sdk-v2-preview/AwsJavaSdk-Glue-2.0.jar`
 3. Java クラスパスの先頭にクライアント jar を追加して、AWS Java SDK で提供された AWS Glue クライアントをオーバーライドします。

```
export CLASSPATH=<path-to-preview-client-jar>:$CLASSPATH
```

4. (オプション) 次の Java アプリケーションを使用して SDK をテストします。このアプリケーションは空のリストを出力するように設定されています。

`accessKey` および `secretKey` は自分の認証情報に置き換え、`us-east-1` は使用しているリージョンに置き換えます。

```
import com.amazonaws.auth.AWSCredentials;
import com.amazonaws.auth.AWSCredentialsProvider;
import com.amazonaws.auth.AWSStaticCredentialsProvider;
import com.amazonaws.auth.BasicAWSCredentials;
import com.amazonaws.services.glue.AWSGlue;
import com.amazonaws.services.glue.AWSGlueClientBuilder;
import com.amazonaws.services.glue.model.ListBlueprintsRequest;

public class App{
    public static void main(String[] args) {
        AWSCredentials credentials = new BasicAWSCredentials("accessKey",
"secretKey");
        AWSCredentialsProvider provider = new
AWSStaticCredentialsProvider(credentials);
        AWSGlue glue = AWSGlueClientBuilder.standard().withCredentials(provider)
.withRegion("us-east-1").build();
```

```
ListBlueprintsRequest request = new
ListBlueprintsRequest().withMaxResults(2);
System.out.println(glue.listBlueprints(request));
}
}
```

AWS Python SDK をセットアップする

次の手順では、Python バージョン 2.7 以降、またはバージョン 3.6 以降が、コンピュータ上にインストールされていることを前提としています。

1. 次の boto3 wheel ファイルをダウンロードします。ファイルを開くか保存することを促された場合は、次のファイルを保存します。s3://awsglue-custom-blueprints-preview-artifacts/aws-python-sdk-preview/boto3-1.17.31-py2.py3-none-any.whl
2. 次の botocore wheel ファイルをダウンロードします。s3://awsglue-custom-blueprints-preview-artifacts/aws-python-sdk-preview/botocore-1.20.31-py2.py3-none-any.whl
3. 使用している Python のバージョンを確認します。

```
python --version
```

4. Python のバージョンに応じて、以下のコマンドを入力します (Linux の場合)。

- Python 2.7 以降の場合。

```
python3 -m pip install --user virtualenv
source env/bin/activate
```

- Python 3.6 以降の場合。

```
python3 -m venv python-sdk-test
source python-sdk-test/bin/activate
```

5. botocore wheel ファイルをインストールします。

```
python3 -m pip install <download-directory>/botocore-1.20.31-py2.py3-none-any.whl
```

6. boto3 wheel ファイルをインストールします。

```
python3 -m pip install <download-directory>/boto3-1.17.31-py2.py3-none-any.whl
```

7. `~/.aws/credentials` および `~/.aws/config` ファイル内で、認証情報およびデフォルトのリージョンを指定します。詳細については、「[AWS CLI ユーザーガイド](#)」の「AWS Command Line Interface の設定。」を参照してください。
8. (オプション) セットアップをテストします。次のコマンドは、空のリストを返すように設定されています。

`us-east-1` は、実際のリージョンに置き換えます。

```
$ python
>>> import boto3
>>> glue = boto3.client('glue', 'us-east-1')
>>> glue.list_blueprints()
```

プレビュー AWS CLI をセットアップする

1. まだの場合は、コンピュータ上で AWS Command Line Interface (AWS CLI) のインストールおよび (または) 更新を行ってください。この作業は、Python インストーラユーティリティの `pip` を使用することで、最も簡単に完了できます。

```
pip install awscli --upgrade --user
```

AWS CLI を使用したインストールに関する完全な手順は、「[Installing the AWS Command Line Interface](#)」でご確認いただけます。

2. 次の AWS CLI wheel ファイルをダウンロードします。 `s3://awsglue-custom-blueprints-preview-artifacts/awscli-preview-build/awscli-1.19.31-py2.py3-none-any.whl`
3. AWS CLI wheel ファイルをインストールする。

```
python3 -m pip install awscli-1.19.31-py2.py3-none-any.whl
```

4. `aws configure` コマンドを実行します。AWS 認証情報 (アクセスキー、シークレットキーを含む) と AWS リージョンを設定します。AWS CLI の設定に関する情報は、「[Configuring the AWS CLI](#)」でご確認ください。
5. AWS CLI をテストします。次のコマンドは、空のリストを返すように設定されています。

`us-east-1` は、実際のリージョンに置き換えます。

```
aws glue list-blueprints --region us-east-1
```

設計図コードを記述する

作成する各設計図プロジェクトには、少なくとも以下のファイルを含む必要があります。

- ワークフローを定義する Python レイアウトスクリプト。このスクリプトには、ワークフロー内のエンティティ (ジョブとクローラ)、および、それらの間の依存関係を定義する関数が含まれています。
- 以下を定義する設定ファイル `blueprint.cfg`。
 - ワークフローレイアウト定義関数を指す完全なパス。
 - 設計図が受け取るパラメータ。

トピック

- [設計図レイアウトスクリプトを作成する](#)
- [設定ファイルを作成する](#)
- [設計図パラメータを指定する](#)

設計図レイアウトスクリプトを作成する

設計図のレイアウトスクリプトには、ワークフロー内のエンティティを生成するための関数を含む必要があります。この関数には好きな名前を付けることができます。AWS Glue は構成ファイルを使用して、関数の完全修飾名を決定します。

レイアウト関数は以下を実行します。

- (オプション) Job オブジェクトを作成する Job クラスをインスタンス化し、Command および Role の引数を渡します。これらは、AWS Glue コンソールまたは API を使用しながらジョブを作成する際に指定する、ジョブのプロパティです。
- (オプション) Crawler オブジェクトを作成する Crawler クラスをインスタンス化し、名前、ロール、ターゲット引数を渡します。
- オブジェクト (ワークフローエンティティ) 間の依存関係を示すために、追加の引数 `DependsOn` および `WaitForDependencies` を `Job()` および `Crawler()` に渡します。これらの引数については、このセクションの後半で説明します。
- Workflow に返すためのワークフローオブジェクトを作成する AWS Glue クラスをインスタンス化し、Name 引数、Entities 引数、およびオプションの `OnSchedule` 引数を渡します。引数 `Entities` では、ワークフローに含めるすべてのジョブとクローラを指定します。Entities オ

プロジェクトを構築する方法については、このセクションで後述するプロジェクト例をご覧ください。

- Workflow オブジェクトを返します。

Job、Crawler、および Workflow クラスの定義については、「[AWS Glue ブループリントクラス リファレンス](#)」を参照してください。

レイアウト関数は、以下の引数を受け取る必要があります。

引数	説明
user_params	設計図パラメータの名前と値に関する Python ディクショナリ。詳細については、 設計図パラメータを指定する を参照してください。
system_params	2 つのプロパティ (region および accountId) を含む Python ディクショナリ。

次に、Layout.py ファイルに記述された、レイアウトジェネレータスクリプトの例を示します。

```
import argparse
import sys
import os
import json
from awsglue.blueprint.workflow import *
from awsglue.blueprint.job import *
from awsglue.blueprint.crawler import *

def generate_layout(user_params, system_params):

    etl_job = Job(Name="{}_etl_job".format(user_params['WorkflowName']),
                  Command={
                      "Name": "glueetl",
                      "ScriptLocation": user_params['ScriptLocation'],
                      "PythonVersion": "2"
                  },
                  Role=user_params['PassRole'])
    post_process_job = Job(Name="{}_post_process".format(user_params['WorkflowName']),
                           Command={
                               "Name": "pythonshell",
```

```
        "ScriptLocation": user_params['ScriptLocation'],
        "PythonVersion": "2"
    },
    Role=user_params['PassRole'],
    DependsOn={
        etl_job: "SUCCEEDED"
    },
    WaitForDependencies="AND")
sample_workflow = Workflow(Name=user_params['WorkflowName'],
    Entities=Entities(Jobs=[etl_job, post_process_job]))
return sample_workflow
```

このサンプルスクリプトは、必要な設計図ライブラリをインポートし、2つのジョブを含むワークフローを生成する `generate_layout` 関数をインクルードします。これは非常にシンプルなスクリプトです。より複雑なスクリプトでは、追加のロジックとパラメータを使用して、多数のジョブとクローラを含むワークフローを生成できます。さらに、ジョブとクローラの数を変更可変にすることも可能です。

DependsOn 引数を使用する

引数 `DependsOn` は、このエンティティがワークフロー内の他のエンティティに対し維持する依存関係を、ディクショナリ向けに表現したものです。これには、以下の形式が使用されます。

```
DependsOn = {dependency1 : state, dependency2 : state, ...}
```

このディクショナリーのキーはエンティティの名前ではなく、オブジェクトリファレンスを表し、値は監視すべき状態に対応する文字列です。AWS Glue 適切なトリガーを推測します。有効な状態については、「[Condition Structure](#)」を参照してください。

例えば、ジョブは、クローラの正常な完了に依存する場合があります。crawler2 という名前のクローラオブジェクトを以下のように定義したとします。

```
crawler2 = Crawler(Name="my_crawler", ...)
```

この時、crawler2 に依存するオブジェクトには、次のようなコンストラクタ引数が含まれます。

```
DependsOn = {crawler2 : "SUCCEEDED"}
```

例:

```
job1 = Job(Name="Job1", ..., DependsOn = {crawler2 : "SUCCEEDED", ...})
```

仮にエンティティで `DependsOn` が省略されている場合には、そのエンティティはワークフローの開始トリガーに依存します。

WaitForDependencies 引数を使用する

`WaitForDependencies` 引数は、ジョブまたはクローラエンティティが、依存するすべてのエンティティが完了するまで待機するのか、あるいは、いずれかのエンティティが完了するまで待機するのかを決定します。

許容される値は「AND」または「ANY」です。

OnSchedule 引数を使用する

`Workflow` クラスコンストラクタの引数 `OnSchedule` は、ワークフローの開始トリガーを定義する cron 式です。

この引数を指定した場合、AWS Glue はスケジュールと対応付けながら、スケジュールトリガーを作成します。特に指定されていない場合は、ワークフローの開始トリガーはオンデマンドトリガーになります。

設定ファイルを作成する

設計図の設定ファイルは必須のファイルで、ワークフローを生成するためのスクリプトエントリーポイントと、設計図が受け取るパラメータを定義します。ファイル名は、`blueprint.cfg` とする必要があります。

設定ファイルの例を次に示します。

```
{
  "layoutGenerator": "DemoBlueprintProject.Layout.generate_layout",
  "parameterSpec" : {
    "WorkflowName" : {
      "type": "String",
      "collection": false
    },
    "WorkerType" : {
      "type": "String",
      "collection": false,
      "allowedValues": ["G1.X", "G2.X"],
      "defaultValue": "G1.X"
    }
  }
}
```

```
    },
    "Dpu" : {
      "type" : "Integer",
      "allowedValues" : [2, 4, 6],
      "defaultValue" : 2
    },
    "DynamoDBTableName": {
      "type": "String",
      "collection" : false
    },
    "ScriptLocation" : {
      "type": "String",
      "collection": false
    }
  }
}
```

layoutGenerator プロパティでは、レイアウトを生成するスクリプト内の関数について、完全な修飾名を指定します。

parameterSpec プロパティでは、この設計図が受け取るパラメータを指定します。詳細については、[設計図パラメータを指定する](#) を参照してください。

Important

設定ファイルに設計図パラメータとしてワークフロー名を含めるか、あるいはレイアウトスクリプトにより一意のワークフロー名を生成する必要があります。

設計図パラメータを指定する

設定ファイルには、parameterSpec JSON オブジェクトの設計図パラメータの仕様が含まれます。parameterSpec には、1 つ以上のパラメータオブジェクトが含まれます。

```
"parameterSpec": {
  "<parameter_name>": {
    "type": "<parameter-type>",
    "collection": true|false,
    "description": "<parameter-description>",
    "defaultValue": "<default value for the parameter if value not specified>"
    "allowedValues": "<list of allowed values>"
  },
}
```

```

    "<parameter_name>": {
      ...
    }
  }
}

```

以下は、各パラメータオブジェクトでのコーディング規則です。

- パラメータ名と type は必須です。その他のプロパティはすべてオプションです。
- defaultValue プロパティを指定する場合のパラメータはオプションです。それ以外の場合、パラメータは必須であり、この設計図からワークフローを作成するデータアナリストは、その値を提供する必要があります。
- collection プロパティに true を設定した場合、このパラメータは値のコレクションを参照できます。コレクションは、任意のデータ型にすることができます。
- allowedValues を指定すると、ブループリントからワークフローを作成する際にデータアナリストが選択する、値のドロップダウンリストが AWS Glue コンソールに表示されます。

type で使用可能な値を以下に示します。

パラメータのデータ型	メモ
String	-
Integer	-
Double	-
Boolean	指定できる値は true および false です。AWS Glue コンソールの [Create a workflow from <blueprint>] (<blueprint> からのワークフローの作成) ページで、チェックボックスを生成します。
S3Uri	s3:// で開始する Amazon S3 パスを完成させます。[Create a workflow from <blueprint>] (<blueprint> からのワークフローの作成) ページで、テキストフィールドおよび [Browse] (参照) ボタンを生成します。
S3Bucket	Amazon S3 バケット名のみ。[Create a workflow from <blueprint>] (<blueprint> からのワークフローの作成) ページで、バケットピッカーを生成します。

パラメータのデータ型	メモ
IAMRoleArn	AWS Identity and Access Management (IAM ロール) の Amazon リソースネーム (ARN)。[Create a workflow from <blueprint>] (<blueprint> からのワークフローの作成) ページで、ロールピッカーを生成します。
IAMRoleName	IAM ロールの名前。[Create a workflow from <blueprint>] (<blueprint> からのワークフローの作成) ページで、ロールピッカーを生成します。

設計図プロジェクト例

データ形式の変換は、抽出、変換、ロード (ETL) で頻繁に行われるユースケースです。一般的な分析ワークロードでは、CSV や JSON などのテキスト形式よりも、Parquet や ORC などの列ベースのファイル形式の使用が好まれます。このサンプル設計図を使用すると、Amazon S3 上のファイルで、CSV/JSON/その他のデータを Parquet に変換できます。

この設計図では、設計図パラメータで定義された S3 パスのリストを取得し、そのデータを Parquet 形式に変換した上で、別の設計図パラメータで指定された S3 のロケーションに書き込んでいます。このレイアウトスクリプトは、パスごとにクローラとジョブを作成します。また、このレイアウトスクリプトでは、Conversion.py 内の ETL スクリプトを、別の設計図パラメータで指定された S3 バケットにアップロードしています。その後、このレイアウトスクリプトは、アップロードされたスクリプトを各ジョブの ETL スクリプトとして指定します。プロジェクト用の ZIP アーカイブには、レイアウトスクリプト、ETL スクリプト、および設計図の設定ファイルが含まれています。

その他の設計図プロジェクトの例については、「[設計図の例](#)」を参照してください。

以下は、Layout.py ファイル内にあるレイアウトスクリプトです。

```
from awsglue.blueprint.workflow import *
from awsglue.blueprint.job import *
from awsglue.blueprint.crawler import *
import boto3

s3_client = boto3.client('s3')

# Ingesting all the S3 paths as Glue table in parquet format
def generate_layout(user_params, system_params):
    #Always give the full path for the file
```

```

with open("ConversionBlueprint/Conversion.py", "rb") as f:
    s3_client.upload_fileobj(f, user_params['ScriptsBucket'], "Conversion.py")
etlScriptLocation = "s3://{}/Conversion.py".format(user_params['ScriptsBucket'])

crawlers = []
jobs = []
workflowName = user_params['WorkflowName']
for path in user_params['S3Paths']:
    tablePrefix = "source_"
    crawler = Crawler(Name="{}_crawler".format(workflowName),
                      Role=user_params['PassRole'],
                      DatabaseName=user_params['TargetDatabase'],
                      TablePrefix=tablePrefix,
                      Targets= {"S3Targets": [{"Path": path}]})
    crawlers.append(crawler)
    transform_job = Job(Name="{}_transform_job".format(workflowName),
                       Command={"Name": "glueetl",
                                "ScriptLocation": etlScriptLocation,
                                "PythonVersion": "3"},
                       Role=user_params['PassRole'],
                       DefaultArguments={"--database_name":
user_params['TargetDatabase'],
                                        "--table_prefix": tablePrefix,
                                        "--region_name": system_params['region'],
                                        "--output_path":
user_params['TargetS3Location']},
                       DependsOn={crawler: "SUCCEEDED"},
                       WaitForDependencies="AND")
    jobs.append(transform_job)
conversion_workflow = Workflow(Name=workflowName, Entities=Entities(Jobs=jobs,
Crawlers=crawlers))
return conversion_workflow

```

以下は、対応する設計図の設定ファイル blueprint.cfg です。

```

{
  "layoutGenerator": "ConversionBlueprint.Layout.generate_layout",
  "parameterSpec" : {
    "WorkflowName" : {
      "type": "String",
      "collection": false,
      "description": "Name for the workflow."
    },
  },

```

```
"S3Paths" : {
  "type": "S3Uri",
  "collection": true,
  "description": "List of Amazon S3 paths for data ingestion."
},
"PassRole" : {
  "type": "IAMRoleName",
  "collection": false,
  "description": "Choose an IAM role to be used in running the job/crawler"
},
"TargetDatabase": {
  "type": "String",
  "collection" : false,
  "description": "Choose a database in the Data Catalog."
},
"TargetS3Location": {
  "type": "S3Uri",
  "collection" : false,
  "description": "Choose an Amazon S3 output path: ex:s3://<target_path>/."
},
"ScriptsBucket": {
  "type": "S3Bucket",
  "collection": false,
  "description": "Provide an S3 bucket name(in the same AWS Region) to store
the scripts."
}
}
```

以下に示す、ファイル Conversion.py 内のスクリプトは、アップロードされた ETL スクリプトです。ここでは、変換中にパーティション分割スキームが保持されることに注意してください。

```
import sys
from pyspark.sql.functions import *
from pyspark.context import SparkContext
from awsglue.transforms import *
from awsglue.context import GlueContext
from awsglue.job import Job
from awsglue.utils import getResolvedOptions
import boto3

args = getResolvedOptions(sys.argv, [
  'JOB_NAME',
```

```
'region_name',
'database_name',
'table_prefix',
'output_path'])
databaseName = args['database_name']
tablePrefix = args['table_prefix']
outputPath = args['output_path']

glue = boto3.client('glue', region_name=args['region_name'])

glue_context = GlueContext(SparkContext.getOrCreate())
spark = glue_context.spark_session
job = Job(glue_context)
job.init(args['JOB_NAME'], args)

def get_tables(database_name, table_prefix):
    tables = []
    paginator = glue.get_paginator('get_tables')
    for page in paginator.paginate(DatabaseName=database_name, Expression=table_prefix
+"""):
        tables.extend(page['TableList'])
    return tables

for table in get_tables(databaseName, tablePrefix):
    tableName = table['Name']
    partitionList = table['PartitionKeys']
    partitionKeys = []
    for partition in partitionList:
        partitionKeys.append(partition['Name'])

# Create DynamicFrame from Catalog
dyf = glue_context.create_dynamic_frame.from_catalog(
    name_space=databaseName,
    table_name=tableName,
    additional_options={
        'useS3ListImplementation': True
    },
    transformation_ctx='dyf'
)

# Resolve choice type with make_struct
dyf = ResolveChoice.apply(
    frame=dyf,
    choice='make_struct',
```

```
        transformation_ctx='resolvechoice_' + tableName
    )

    # Drop null fields
    dyf = DropNullFields.apply(
        frame=dyf,
        transformation_ctx="dropnullfields_" + tableName
    )

    # Write DynamicFrame to S3 in glueparquet
    sink = glue_context.getSink(
        connection_type="s3",
        path=outputPath,
        enableUpdateCatalog=True,
        partitionKeys=partitionKeys
    )
    sink.setFormat("glueparquet")

    sink.setCatalogInfo(
        catalogDatabase=databaseName,
        catalogTableName=tableName[len(tablePrefix):]
    )
    sink.writeFrame(dyf)

job.commit()
```

Note

このサンプルの設計図に対し、入力として提供できる Amazon S3 パスは 2 つだけです。これは、AWS Glue トリガーが呼び出すことができるクローラアクションが、2 つだけに制限されているためです。

設計図をテストする

コードの開発時には、ローカルなテストを実行して、ワークフローのレイアウトに誤りがないことを確認する必要があります。

ローカルテストでは AWS Glue ジョブ、クローラ、トリガーは生成されません。代わりに、レイアウトスクリプトをローカルで実行し、`to_json()` および `validate()` メソッドによりオブジェクトを画面表示し、エラーを発見します。これらのメソッドは、ライブラリで定義されている 3 つのクラスすべてで使用できます。

`user_params` がレイアウト関数に渡す、引数 `system_params` および AWS Glue を処理するには 2 つの方法があります。テストベンチコードでは、サンプル設計図パラメータ値のディクショナリを作成し、それを引数 `user_params` としてレイアウト関数に渡すことができます。または、`user_params` への参照を削除し、その部分をハードコーディングのための文字列に置き換えます。

コードの引数 `system_params` の中で、`region` および `accountId` プロパティを利用している場合には、`system_params` を独自の辞書に含めて渡すことができます。

設計図をテストするには

1. Python インタプリタをライブラリのあるディレクトリで起動するか、設計図ファイルならびに提供されたライブラリを、希望の統合開発環境 (IDE) にロードします。
2. コードによって提供されたライブラリがインポートされていることを確認します。
3. 任意のエンティティまたは Workflow オブジェクトで `validate()` または `to_json()` 呼び出すためのコードを、レイアウト関数に追加します。例えば、コードが `mycrawler` という名前の Crawler オブジェクトを作成する場合には、以下のように `validate()` を呼び出せます。

```
mycrawler.validate()
```

`mycrawler` は以下のように表示できます。

```
print(mycrawler.to_json())
```

オブジェクトで `to_json` を呼び出す場合、別途 `validate()` を呼び出す必要はありません。`validate()` は `to_json()` により呼び出されます。

これにより、ワークフローオブジェクトでのこれらのメソッドの呼び出しが、最も効率的に行えます。スクリプトで、ワークフローオブジェクトに `my_workflow` と名前を付けているのであれば、次のようにワークフローオブジェクトを検証して画面表示します。

```
print(my_workflow.to_json())
```

`to_json()` と `validate()` の詳細については、「[クラスメソッド](#)」を参照してください。

また、このセクションの後半の例に示すように、`pprint` をインポートし、ワークフローオブジェクトを書式を設定しながら表示することも可能です。

4. コードを実行し、エラーを修正した上で、最後に `validate()` または `to_json()` の呼び出しを削除します。

Example

次の例では、サンプルの設計図パラメータのディクショナリを構築し、それをレイアウト関数 `generate_compaction_workflow` の引数 `user_params` に渡す方法を示しています。また、生成されたワークフローオブジェクトを、書式を指定しながら印刷するための方法も示しています。

```
from pprint import pprint
from awsglue.blueprint.workflow import *
from awsglue.blueprint.job import *
from awsglue.blueprint.crawler import *

USER_PARAMS = {"WorkflowName": "compaction_workflow",
               "ScriptLocation": "s3://awsexamplebucket1/scripts/threaded-
compaction.py",
               "PassRole": "arn:aws:iam::111122223333:role/GlueRole-ETL",
               "DatabaseName": "cloudtrial",
               "TableName": "ct_cloudtrail",
               "CoalesceFactor": 4,
               "MaxThreadWorkers": 200}

def generate_compaction_workflow(user_params: dict, system_params: dict) -> Workflow:
    compaction_job = Job(Name=f"{user_params['WorkflowName']}_etl_job",
                        Command={"Name": "glueetl",
                                "ScriptLocation": user_params['ScriptLocation'],
                                "PythonVersion": "3"},
                        Role="arn:aws:iam::111122223333:role/
AWSGlueServiceRoleDefault",
                        DefaultArguments={"DatabaseName": user_params['DatabaseName'],
                                         "TableName": user_params['TableName'],
                                         "CoalesceFactor":
user_params['CoalesceFactor'],
                                         "max_thread_workers":
user_params['MaxThreadWorkers']})

    catalog_target = {"CatalogTargets": [{"DatabaseName": user_params['DatabaseName'],
                                         "Tables": [user_params['TableName']]}]}

    compacted_files_crawler = Crawler(Name=f"{user_params['WorkflowName']}_post_crawl",
```

```
Targets = catalog_target,
Role=user_params['PassRole'],
DependsOn={compaction_job: "SUCCEEDED"},
WaitForDependencies="AND",
SchemaChangePolicy={"DeleteBehavior": "LOG"})

compaction_workflow = Workflow(Name=user_params['WorkflowName'],
                               Entities=Entities(Jobs=[compaction_job],
Crawlers=[compacted_files_crawler]))
return compaction_workflow

generated = generate_compaction_workflow(user_params=USER_PARAMS, system_params={})
gen_dict = generated.to_json()

pprint(gen_dict)
```

設計図を公開する

作成した設計図は、Amazon S3 にアップロードする必要があります。これには、設計図を公開するために使用する Amazon S3 バケットに対する、書き込みアクセス許可が必要です。また、ブループリントの登録を行う AWS Glue 管理者には、Amazon S3 バケットからの読み取りアクセスが許可されている必要があります。AWS Identity and Access Management 設計図でのペルソナおよびロールのために提案されている、AWS Glue (IAM) アクセス許可ポリシーについては、「[AWS Glue ブループリントでのペルソナおよびロール用のアクセス許可](#)」を参照してください。

設計図を公開するには

1. 必要なスクリプト、リソース、および設計図の設定ファイルを作成します。
2. すべてのファイルを ZIP アーカイブに追加し、その ZIP ファイルを Amazon S3 にアップロードします。設計図を登録し実行するリージョンと、同じリージョンにある S3 バケットを使用します。

ZIP ファイルは、次のコマンドを使用して、コマンドラインから作成することができます。

```
zip -r folder.zip folder
```

3. AWS での希望するアカウントに読み取りアクセス許可を付与する、バケットポリシーを追加します。次に、そのポリシー例を示します。

```
{
```

```
"Version": "2012-10-17",
"Statement": [
  {
    "Effect": "Allow",
    "Principal": {
      "AWS": "arn:aws:iam::111122223333:root"
    },
    "Action": "s3:GetObject",
    "Resource": "arn:aws:s3:::my-blueprints/*"
  }
]
```

4. Amazon S3 バケットに対する IAM の `s3:GetObject` の許可を、AWS Glue 管理者またはブループリントを登録する任意のユーザーに付与します。管理者に付与するポリシーの例については、「[ブループリントのための AWS Glue 管理者のアクセス許可](#)」を参照してください。

ブループリントに関するローカルでのテストが完了したら、AWS Glue 上でもブループリントをテストします。ブループリントを AWS Glue 上でテストするには、登録が完了している必要があります。登録された設計図を表示できるユーザーを制限するには、IAM 認証を使用するか、テスト用に分離されたアカウントを使用します。

 以下も参照してください。

- [AWS Glue でブループリントを登録する](#)

AWS Glue ブループリントクラスリファレンス

AWS Glue ブループリント用のライブラリでは、ワークフローレイアウトスクリプトで使用する 3 つのクラス (Job、Crawler、および Workflow) を定義します。

トピック

- [Job クラス](#)
- [Crawler クラス](#)
- [Workflow クラス](#)
- [クラスメソッド](#)

Job クラス

Job クラスは、AWS Glue での ETL ジョブを表します。

必須のコンストラクター引数

Job クラスのコンストラクタに必須な引数を、以下に示します。

引数名	型	説明
Name	str	ジョブに割り当てる名前。AWS Glue では他のブループリントの実行で作成されたジョブとの区別するため、名前にはランダムに生成されるサフィックスを付けています。
Role	str	ジョブが実行中に引き受けるロールの Amazon リソースネーム (ARN)。
Command	dict	API ドキュメントの JobCommand 構造 に掲載されている Job コマンド。

オプションのコンストラクター引数

Job クラスのコンストラクタで、オプションとなっている引数を以下に示します。

引数名	型	説明
DependsOn	dict	ジョブが依存するワークフローエンティティの一覧。詳細については、 DependsOn 引数を使用する を参照してください。
WaitForDependencies	str	実行前のジョブが、依存するすべてのエンティティが完了するまで待機するのか、あるいは任意のエンティティが完了することを待つのかを示します。詳細については、 WaitForDependencies 引数を使用する を参照してください。ジョブが 1 つのエンティティのみに依存する場合は、この設定は省略します。

引数名	型	説明
(ジョブプロパティ)	-	Job 構造 API ドキュメント (AWS Glue と CreatedOn 以外)LastModifiedOn の一覧に掲載されている、いずれかのジョブのプロパティ。

Crawler クラス

Crawler クラスは、AWS Glue でのクローラーを表します。

必須のコンストラクター引数

Crawler クラスのコンストラクタに必須な引数を、以下に示します。

引数名	型	説明
Name	str	AWS Glue は、ランダムに生成されたサフィックスをクローラーの名前に付加することで、異なるブループリントの実行によって作成されたクローラを区別します。
Role	str	クローラが実行中に引き受ける必要のあるロールの ARN。
Targets	dict	クロールするターゲットのコレクション。 Targets クラスコンストラクターの引数は、API ドキュメントの CrawlerTargets 構造 に定義されています。Targets コンストラクタの引数はすべてオプションです。ただし、少なくとも 1 つを渡す必要があります。

オプションのコンストラクター引数

Crawler クラスのコンストラクタで、オプションとなっている引数を以下に示します。

引数名	型	説明
DependsOn	dict	クローラが依存するワークフローエンティティのリスト。詳細については、 DependsOn 引数を使用する を参照してください。
WaitForDependencies	str	実行前のクローラが、依存するすべてのエンティティが完了するまで待機するのか、任意のエンティティが完了することを待つのかを示します。詳細については、 WaitForDependencies 引数を使用する を参照してください。クローラが 1 つのエンティティのみに依存する場合は、この設定は省略します。
(クローラのプロパティ)	-	<p>Crawler 構造 API ドキュメントの AWS Glue の一覧に掲載されている、いずれかのクローラのプロパティ (以下のものを除く)。</p> <ul style="list-style-type: none"> • State • CrawlElapsedTime • CreationTime • LastUpdated • LastCrawl • Version

Workflow クラス

Workflow クラスは、AWS Glue のワークフローを表します。ワークフローレイアウトスクリプトは、Workflow オブジェクト。AWS Glue でこのオブジェクトに基づいてワークフローを作成します。

必須のコンストラクター引数

Workflow クラスのコンストラクタに必須な引数を、以下に示します。

引数名	型	説明
Name	str	ワークフローに割り当てる名前。
Entities	Entities	ワークフローに含めるエンティティ (ジョブおよびクローラ) のコレクション。Entities クラスのコンストラクタは、引数 Jobs (Job オブジェクトのリスト)、および、Crawlers 引数 (Crawler オブジェクトのリスト) を受け取ります。

オプションのコンストラクター引数

Workflow クラスのコンストラクタで、オプションとなっている引数を以下に示します。

引数名	型	説明
Description	str	「」を参照してください Workflow 構造
DefaultRunProperties	dict	「」を参照してください Workflow 構造
OnSchedule	str	cron 式

クラスメソッド

上記の 3 つのクラスには、以下のメソッドが含まれています。

validate()

オブジェクトのプロパティを検証し、エラーが見つかった場合はメッセージを出力して終了します。エラーが発見されない場合は、出力を生成しません。Workflow クラスでは、ワークフロー内のすべてのエンティティで、自分自身を呼び出します。

to_json()

オブジェクトを JSON としてシリアル化します。また、validate() を呼び出します。Workflow クラスでは、この JSON オブジェクトにはジョブとクローラのリストと、ジョブとクローラの依存関係の仕様によって生成されたトリガーのリストが含まれます。

設計図の例

[AWS Glue ブループリントの Githubリポジトリ](#)には、設計図プロジェクトの例が多数提供されています。これらのサンプルは参考用であり、本番での使用は意図されていません。

これらのサンプルプロジェクトのタイトルを以下に示します。

- **Compaction:** この設計図では、必要なファイルサイズに基づいて、複数の入力ファイルをより大きなチャンクの中に圧縮するジョブを作成します。
- **Conversion:** この設計図では、さまざまな標準ファイル形式の入力ファイルを、分析ワークロード向けに最適化された Apache Parquet 形式に変換します。
- **Crawling Amazon S3 locations:** この設計図は、複数の Amazon S3 ロケーションをクロールして、Data Catalog にメタデータテーブルを追加します。
- **Custom connection to Data Catalog:** この設計図では、AWS Glue のカスタムコネクタを使用しデータストアにアクセスします。レコードを読み取り、レコードスキーマに基づいて AWS Glue データカタログのテーブル定義を設定します。
- **Encoding:** この設計図では、UTF 以外のファイルを UTF エンコードされたファイルに変換します。
- **Partitioning:** この設計図では、パーティショニングジョブを作成します。特定のパーティションキーに基づいて出力ファイルをパーティションに配置します。
- **Importing Amazon S3 data into a DynamoDB table:** この設計図では、Amazon S3 から DynamoDB テーブルにデータをインポートします。
- **Standard table to governed:** この設計図では、Lake Formation テーブルに AWS Glue データカタログテーブルをインポートします。

AWS Glue でブループリントを登録する

AWS Glue デベロッパーが設計図をコーディングして、Amazon Simple Storage Service (Amazon S3) に ZIP アーカイブをアップロードした後、AWS Glue 管理者はその設計図を登録する必要があります。設計図は、登録されることで使用が可能になります。

ブループリントを登録すると、AWS Glue はブループリントアーカイブを予約済みの Amazon S3 ロケーションにコピーします。その後、アーカイブは、アップロードした場所から削除できます。

設計図を登録するには、アップロードされたアーカイブが保存される Amazon S3 ロケーションに対する、読み取りアクセス許可が必要です。また、AWS Identity and Access Management (IAM) のア

クセス許可 `glue:CreateBlueprint` も必要です。ブループリントの登録、表示、および保守を行う必要がある AWS Glue 管理者のために、提案されているアクセス許可については、[ブループリントのための AWS Glue 管理者のアクセス許可](#) を参照してください。

ブループリントは AWS Glue コンソール、AWS Glue API、または AWS Command Line Interface (AWS CLI) を使用して登録できます。

設計図を登録するには (コンソール)

1. Amazon S3 にある設計図の ZIP アーカイブに対する、読み取りアクセス許可 (`s3:GetObject`) が付与されていることを確認します。
2. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。

設計図を登録するための、アクセス許可が付与されたユーザーとしてサインインします。設計図の ZIP アーカイブを含む Amazon S3 バケットと同じ、AWS リージョンに切り替えます。

3. ナビゲーションペインで [Blueprints] (ブループリント) を選択します。次に、[Blueprints] (設計図) ページで [Add blueprint] (設計図を追加) をクリックします。
4. 設計図の名前と、必要に応じて説明を入力します。
5. [ZIP archive location (S3)] (ZIP アーカイブの場所 (S3)) に、アップロードされた設計図の ZIP アーカイブがある Amazon S3 パスを入力します。パスの先頭は `s3://` で開始し、またアーカイブのファイル名を含めます。
6. (オプション) 1 つ以上のタグを追加します。
7. [Add blueprint] (設計図を追加) をクリックします。

[Blueprints] (設計図) ページが再び開き、設計図のステータスが「CREATING」となっていることを表示します。ステータスが「ACTIVE」または「FAILED」に変わるまで、[Refresh] (更新) ボタンをクリックします。

8. ステータスが「FAILED」となった場合は、設計図を選択した上で、[Actions] (アクション) メニューで [View] (表示) をクリックします。

詳細ページに、失敗の理由が表示されます。エラーメッセージが、「Unable to access object at location...」または「Access denied on object at location...」となっている場合は、以下の要件を確認してください。

- サインインしているユーザーには、Amazon S3 内にある設計図の ZIP アーカイブに対する読み取りアクセス許可が必要です。

- ZIP アーカイブを含む Amazon S3 バケットには、オブジェクトに対する読み取りアクセス許可を AWS アカウント ID に付与するためのバケットポリシーが必要です。(詳細については、[AWS Glue のブループリントの開発](#) を参照してください)。
 - 使用している Amazon S3 バケットは、コンソールでサインインしているリージョンと同じリージョンにある必要があります。
9. データアナリストに、設計図へのアクセス許可が付与されていることを確認します。

データアナリストのために提案された IAM ポリシーは、「[設計図のためのデータアナリストの権限](#)」で確認できます。このポリシーでは、任意のリソースに `glue:GetBlueprint` を付与します。リソースレベルでポリシーをより細かく設定している場合には、その新しく作成したリソースに対するアクセス許可をデータアナリストに付与します。

設計図を登録するには (AWS CLI)

1. 次のコマンドを入力します。

```
aws glue create-blueprint --name <blueprint-name> [--description <description>] --blueprint-location s3://<s3-path>/<archive-filename>
```

2. 設計図のステータスを確認するには、次のコマンドを入力します。ステータスが ACTIVE または FAILED に変化するまで、このコマンドを繰り返します。

```
aws glue get-blueprint --name <blueprint-name>
```

ステータスが FAILED となり、エラーメッセージに「Unable to access object at location...」または「Access denied on object at location...」が出力された場合は、以下の要件を確認してください。

- サインインしているユーザーには、Amazon S3 内にある設計図の ZIP アーカイブに対する読み取りアクセス許可が必要です。
- ZIP アーカイブを含む Amazon S3 バケットには、オブジェクトに対する読み取りアクセス許可を AWS アカウント ID に付与するためのバケットポリシーが必要です。詳細については、[設計図を公開する](#) を参照してください。
- 使用している Amazon S3 バケットは、コンソールでサインインしているリージョンと同じリージョンにある必要があります。

i 以下も参照してください。

- [AWS Glue のブループリント概要](#)

AWS Glue でブループリントを表示する

設計図を表示して、その説明とステータスおよびパラメータ仕様を確認し、設計図の ZIP アーカイブをダウンロードします。

ブループリントを表示するには、AWS Glue コンソール、AWS Glue API、または AWS Command Line Interface (AWS CLI) を使用します。

設計図を表示するには (コンソール)

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。
2. ナビゲーションペインで [Blueprints] (ブループリント) を選択します。
3. [Blueprints] (設計図) ページで、設計図を選択します。次に、[Actions] (アクション) メニューで [View] (表示) をクリックします。

設計図を表示するには (AWS CLI)

- 設計図の名前、説明、ステータスのみを表示するには、次のコマンドを入力します。<code>blueprint-name</code>は、表示する設計図の名前に置き換えます。

```
aws glue get-blueprint --name <blueprint-name>
```

次のような出力が表示されます。

```
{
  "Blueprint": {
    "Name": "myDemoBP",
    "CreatedOn": 1587414516.92,
    "LastModifiedOn": 1587428838.671,
    "BlueprintLocation": "s3://awsexamplebucket1/demo/
DemoBlueprintProject.zip",
    "Status": "ACTIVE"
  }
}
```

```
}
```

以下のコマンドを入力することで、パラメータの仕様も表示できます。

```
aws glue get-blueprint --name <blueprint-name> --include-parameter-spec
```

次のような出力が表示されます。

```
{
  "Blueprint": {
    "Name": "myDemoBP",
    "CreatedOn": 1587414516.92,
    "LastModifiedOn": 1587428838.671,
    "ParameterSpec": "{\"WorkflowName\":{\"type\":\"String\",\"collection\":false,\"description\":null,\"defaultValue\":null,\"allowedValues\":null},\"PassRole\":{\"type\":\"String\",\"collection\":false,\"description\":null,\"defaultValue\":null,\"allowedValues\":null},\"DynamoDBTableName\":{\"type\":\"String\",\"collection\":false,\"description\":null,\"defaultValue\":null,\"allowedValues\":null},\"ScriptLocation\":{\"type\":\"String\",\"collection\":false,\"description\":null,\"defaultValue\":null,\"allowedValues\":null}}",
    "BlueprintLocation": "s3://awsexamplebucket1/demo/DemoBlueprintProject.zip",
    "Status": "ACTIVE"
  }
}
```

引数 `--include-blueprint` を使用して、出力に URL を含めることもできます。この URL をブラウザに貼り付けて、AWS Glue に保存されているブループリントの zip アーカイブをダウンロードできます。

 以下も参照してください。

- [AWS Glue のブループリント概要](#)

AWS Glue でブループリントを更新する

レイアウトスクリプトの改訂、設計図のパラメータセットの改訂、またはサポートファイルの改訂が行われた場合は、設計図を更新できます。更新により、新しいバージョンの設計図が作成されます。

設計図を更新した場合でも、その設計図から作成された既存のワークフローに影響は与えません。

設計図の更新には、AWS Glue コンソール、AWS Glue API、または AWS Command Line Interface (AWS CLI) を使用します。

以下に示す手順では、AWS Glue デベロッパーが、設計図の ZIP アーカイブを Amazon S3 で作成およびアップロードし、その更新を行うと想定しています。

設計図を更新するには (コンソール)

1. Amazon S3 にある設計図の ZIP アーカイブに対する、読み取りアクセス許可 (s3:GetObject) が付与されていることを確認します。
2. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。

設計図を更新するためのアクセス許可が付与されたユーザーとしてサインインします。設計図の ZIP アーカイブを含む Amazon S3 バケットと同じ、AWS リージョンに切り替えます。

3. ナビゲーションペインで [Blueprints] (ブループリント) を選択します。
4. [Blueprints] (設計図) ページで設計図を選択し、[Actions] (アクション) メニューで、[Edit] (編集) をクリックします。
5. [Edit a blueprint] (設計図の編集) ページで、設計図の [Description] (説明)、もしくは [ZIP archive location (S3)] (ZIP アーカイブの場所 (S3)) を更新します。パスにはアーカイブのファイル名を含める必要があります。
6. [Save] を選択します。

[Blueprints] (設計図) ページが再び開き、設計図のステータスが「UPDATING」となっていることを表示します。ステータスが「ACTIVE」または「FAILED」に変わるまで、[Refresh] (更新) ボタンをクリックします。

7. ステータスが「FAILED」となった場合は、設計図を選択した上で、[Actions] (アクション) メニューで [View] (表示) をクリックします。

詳細ページに、失敗の理由が表示されます。エラーメッセージが、「Unable to access object at location...」または「Access denied on object at location...」となっている場合は、以下の要件を確認してください。

- サインインしているユーザーには、Amazon S3 内にある設計図の ZIP アーカイブに対する読み取りアクセス許可が必要です。

- ZIP アーカイブを含む Amazon S3 バケットには、オブジェクトに対する読み取りアクセス許可を AWS アカウント ID に付与するためのバケットポリシーが必要です。詳細については、[設計図を公開する](#) を参照してください。
- 使用している Amazon S3 バケットは、コンソールでサインインしているリージョンと同じリージョンにある必要があります。

Note

更新が失敗した場合、次回の設計図実行では、正常に登録または更新された設計図の最新バージョンが使用されます。

設計図を更新するには (AWS CLI)

1. 次のコマンドを入力します。

```
aws glue update-blueprint --name <blueprint-name> [--description <description>] --  
blueprint-location s3://<s3-path>/<archive-filename>
```

2. 設計図のステータスを確認するには、次のコマンドを入力します。ステータスが ACTIVE または FAILED に変化するまで、このコマンドを繰り返します。

```
aws glue get-blueprint --name <blueprint-name>
```

ステータスが FAILED となり、エラーメッセージに「Unable to access object at location...」または「Access denied on object at location...」が出力された場合は、以下の要件を確認してください。

- サインインしているユーザーには、Amazon S3 内にある設計図の ZIP アーカイブに対する読み取りアクセス許可が必要です。
- ZIP アーカイブを含む Amazon S3 バケットには、オブジェクトに対する読み取りアクセス許可を AWS アカウント ID に付与するためのバケットポリシーが必要です。詳細については、[設計図を公開する](#) を参照してください。
- 使用している Amazon S3 バケットは、コンソールでサインインしているリージョンと同じリージョンにある必要があります。

 以下も参照してください。

- [AWS Glue のブループリント概要](#)

AWS Glue で設計図からワークフローを作成する

AWS Glue ワークフローはコンポーネントを 1 つずつ手動で作成したり、AWS Glue [ブループリント](#) からワークフローを作成することもできます。AWS Glue には一般的なユースケースのブループリントが含まれます。AWS Glue デベロッパーが、追加の設計図を作成することも可能です。

Important

ワークフロー内のジョブ、クローラ、トリガーの総数を 100 以下に制限します。100 を超える値を含めると、ワークフローの実行を再開または停止しようとしたときにエラーが発生することがあります。

設計図を使用すると、設計図で定義されている一般化されたユースケースをベースにしながら、特定のユースケース向けのワークフローをすばやく生成できます。設計図パラメータの値を指定することで、特定のユースケースを定義します。例えば、データセットをパーティション化する設計図では、パラメータとして Amazon S3 のソースとターゲットのパスを使用します。

AWS Glue は、ブループリントを実行することで、そのブループリントからワークフローを作成します。設計図の実行では、指定したパラメータ値が保存されます。このパラメータは、ワークフローとそのコンポーネントの作成に関する進行状況と結果を追跡するためにも使用されます。ワークフローをトラブルシューティングする場合は、設計図の実行を表示して、ワークフローの作成に使用された設計図のパラメータ値を確認します。

ワークフローを作成および表示するには、特定の IAM アクセス許可が必要です。提案されている IAM ポリシーについては、「[設計図のためのデータアナリストの権限](#)」を参照してください。

ブループリントからワークフローを作成するには、AWS Glue コンソール、AWS Glue API、または AWS Command Line Interface (AWS CLI) を使用します。

設計図からワークフローを作成するには (コンソール)

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。

- ワークフローを作成するための、アクセス許可が付与されたユーザーとしてサインインします。
- ナビゲーションペインで [Blueprints] (ブループリント) を選択します。
 - 設計図を選択し、[Actions] (アクション) メニューで、[Create workflow] (ワークフローの作成) をクリックします。
 - [Create a workflow from <blueprint-name>] (<blueprint-name> からのワークフローの作成) ページで、以下の情報を入力します。

設計図パラメータ

これらのパラメータは、設計図の構成によって異なります。パラメータについては、デベロッパーにお問い合わせください。設計図には、通常、ワークフロー名がパラメータとして含まれています。

IAM ロール

ワークフローとそのコンポーネントを作成するために、AWS Glue が引き受けるロールです。ロールには、ワークフロー、ジョブ、クローラ、トリガーを作成および削除するための、アクセス許可が付与されている必要があります。提案されているロールのポリシーについては、「[設計図のロール用のアクセス許可](#)」を参照してください。

- 送信 を選択します。

[Blueprint Details] (設計図の詳細) ページが表示され、その下部に設計図の実行が一覧表示されます。

- 設計図の実行の一覧で、最上位にある設計図の実行から、ワークフローの作成ステータスを確認します。

初期ステータスは、RUNNING です。ステータスが SUCCEEDED または FAILED に遷移するまで、[Refresh] (更新) ボタンをクリックします。

- 次のいずれかを行います。
 - 完了ステータスが SUCCEEDED であれば、[Workflows] (ワークフロー) ページを開き、新しく作成したワークフローを選択して実行できます。ワークフローを実行する前に、設計グラフを確認することができます。
 - 完了ステータスが FAILED の場合には、設計図実行を選択した上で [Actions] (アクション) メニューで [View] (表示) をクリックし、エラーメッセージを表示します。

Express ワークフローと サービス統合の詳細については、以下を参照してください。

- [AWS Glue のワークフローの概要](#)
- [AWS Glue でブループリントを更新する](#)
- [AWS Glue でワークフローを手動により作成および構築する](#)

AWS Glue でブループリントの実行を表示する

設計図の実行を表示して、次の情報を表示します。

- 作成されたワークフローの名前。
- ワークフローの作成に使用された設計図のパラメータ値。
- ワークフロー作成オペレーションのステータス。

ブループリントを表示するには、AWS Glue コンソール、AWS Glue API、または AWS Command Line Interface (AWS CLI) を使用します。

設計図の実行を表示するには (コンソール)

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。
2. ナビゲーションペインで [Blueprints] (ブループリント) を選択します。
3. [Blueprints] (設計図) ページで、設計図を選択します。次に、[Actions] (アクション) メニューで [View] (表示) をクリックします。
4. [Blueprint Details] (設計図の詳細) ページの最下部で、設計図の実行を選択し、[Actions] (アクション) メニューから [View] 表示を選択します。

設計図の実行を表示するには (AWS CLI)

- 次のコマンドを入力します。<blueprint-name> を設計図の名前に置き換えます。<blueprint-run-id> を設計図の実行の ID に置き換えます。

```
aws glue get-blueprint-run --blueprint-name <blueprint-name> --run-id <blueprint-run-id>
```

 以下も参照してください。

- [AWS Glue のブループリント概要](#)

AWS Glue の場合は AWS CloudFormation

AWS CloudFormation は、多くの AWS リソースを作成できるサービスです。AWS Glue には、AWS Glue Data Catalog でオブジェクトを作成するための API オペレーションが用意されています。ただし、AWS Glue オブジェクトや他の関連する AWS リソースオブジェクトを AWS CloudFormation テンプレートファイルで定義して作成するほうが便利な場合があります。この場合、オブジェクトの作成プロセスを自動化できます。

AWS CloudFormation では、JSON (JavaScript Object Notation) または YAML (YAML Ain't Markup Language) のどちらかの簡略化された構文を使用し、AWS リソースの作成を記述します。AWS CloudFormation テンプレートを使用して、データベース、テーブル、パーティション、クローラー、分類子、接続などのデータカタログオブジェクトを定義できます。ジョブ、トリガー、開発エンドポイントなどの ETL オブジェクトを定義することもできます。必要なすべての AWS リソースを記述するテンプレートを作成すると、これらのリソースが AWS CloudFormation で自動的にプロビジョニングおよび設定されます。

詳細については、AWS CloudFormation ユーザーガイドの「[What Is AWS CloudFormation?](#)」および「[Working with AWS CloudFormation Templates](#)」を参照してください。

管理者として AWS Glue と互換性がある AWS CloudFormation テンプレートを使用する場合は、依存する AWS CloudFormation および AWS のサービスとアクションにアクセス権を付与する必要があります。AWS CloudFormation リソースを作成するアクセス権限を付与するには、次のポリシーを、AWS CloudFormation を使用するユーザーにアタッチします。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "cloudformation:*"
      ],
      "Resource": "*"
    }
  ]
}
```

次の表は、AWS CloudFormation テンプレートで自動的に実行できるアクションの一覧です。AWS CloudFormation テンプレートに追加できる AWS リソースタイプやプロパティタイプに関する情報へのリンクが含まれています。

AWS Glue リソース:	AWS CloudFormation テンプレート	AWS Glue サンプル
分類子	AWS::Glue::Classifier	Grok 分類子 、 JSON 分類子 、 XML 分類子
Connection	AWS::Glue::Connection	MySQL 接続
Crawler	AWS::Glue::Crawler	Amazon S3 クローラー 、 MySQL クローラー
データベース	AWS::Glue::Database	空のデータベース 、 テーブルを含むデータベース
開発エンドポイント	AWS::Glue::DevEndpoint	開発エンドポイント
ジョブ	AWS::Glue::Job	Amazon S3 ジョブ 、 JDBC ジョブ
機械学習変換	AWS::Glue::MLTransform	機械学習変換
データ品質ルールセット	AWS::Glue::DataQualityRuleset	データ品質ルールセット 、 EventBridge スケジューラを使用したデータ品質ルールセット
パーティション	AWS::Glue::Partition	テーブルのパーティション
テーブル	AWS::Glue::Table	データベース内のテーブル
Trigger (トリガー)	AWS::Glue::Trigger	オンデマンドのトリガー 、 スケジュールされたトリガー 、 条件付きトリガー

使用を開始するには、以下のサンプルテンプレートを独自のメタデータを使用してカスタマイズします。次に AWS CloudFormation コンソールを使用して AWS CloudFormation スタックを作成し、AWS Glue および関連サービスにオブジェクトを追加します。AWS Glue オブジェクトの多くのフィールドはオプションです。これらのテンプレートは、AWS Glue オブジェクトの使用や機能に必須または必要なフィールドを示しています。

AWS CloudFormation テンプレートは JSON 形式または YAML 形式のいずれかで使用できます。以下の例では、読みやすい YAML を使用しています。各例には、テンプレートで定義されている値を説明するコメント (#) が含まれています。

AWS CloudFormation テンプレートには Parameters セクションを含めることができます。このセクションは、サンプルテキストを編集して変更できます。または、YAML ファイルを AWS CloudFormation コンソールに送信してスタックを作成するときに変更できます。テンプレートの Resources セクションには、AWS Glue および関連オブジェクトの定義が含まれています。AWS CloudFormation テンプレートの構文定義には、詳細なプロパティ構文を含むプロパティが含まれている場合があります。AWS Glue オブジェクトを作成するために、すべてのプロパティが必要になるわけではありません。これらのサンプルは、AWS Glue オブジェクトを作成するための一般的なプロパティの値の例です。

AWS Glue データベース用のサンプル AWS CloudFormation テンプレート

Data Catalog の AWS Glue データベースにはメタデータテーブルが含まれています。このデータベースは、構成するプロパティが非常に少なく、AWS CloudFormation テンプレートを使用して Data Catalog に作成できます。次のサンプルテンプレートでは、使用を開始する方法と、AWS Glue での AWS CloudFormation スタックの使い方を示します。このサンプルテンプレートで作成されるリソースは `cfn-mysampledatabase` というデータベースのみです。このデータベースは、サンプルのテキストを編集するか、YAML の送信時に AWS CloudFormation コンソールで値を変更することで、変更できます。

次に示すのは、AWS Glue データベースを作成するための一般的なプロパティの値の例です。AWS CloudFormation 用の AWS Glue データベーステンプレートの詳細については、「[AWS::Glue::Database](#)」を参照してください。

```
---  
AWSTemplateFormatVersion: '2010-09-09'
```

```
# Sample CloudFormation template in YAML to demonstrate creating a database named
mysampledatabase
# The metadata created in the Data Catalog points to the flights public S3 bucket
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:
  CFNDatabaseName:
    Type: String
    Default: cfn-mysampledatabse

# Resources section defines metadata for the Data Catalog
Resources:
# Create an AWS Glue database
CFNDatabaseFlights:
  Type: AWS::Glue::Database
  Properties:
    # The database is created in the Data Catalog for your account
    CatalogId: !Ref AWS::AccountId
    DatabaseInput:
      # The name of the database is defined in the Parameters section above
      Name: !Ref CFNDatabaseName
      Description: Database to hold tables for flights data
      LocationUri: s3://crawler-public-us-east-1/flight/2016/csv/
      #Parameters: Leave AWS database parameters blank
```

AWS Glue データベース、テーブル、およびパーティション用のサンプル AWS CloudFormation テンプレート

AWS Glue テーブルには、ETL スクリプトで処理するデータの構造と場所を定義するメタデータが含まれています。テーブル内に、データを並列処理するためのパーティションを定義できます。パーティションは、キーを使用して定義したデータのチャンクです。たとえば、キーとして月を使用すると、1月のすべてのデータが同じパーティションに含まれます。AWS Glue では、データベースにテーブルを含め、テーブルにパーティションを含めることができます。

次のサンプルでは、AWS CloudFormation テンプレートを使用して、データベース、テーブル、およびパーティションを事前設定する方法を示します。元のデータ形式は csv であり、カンマ (,) で区切られています。テーブルを作成するには事前にデータベースが必要であり、パーティションを作成するには事前にテーブルが必要であるため、テンプレートでは DependsOn ステートメントを使用して、これらのオブジェクトの作成時に相互の依存関係を定義します。

次のサンプルの値では、一般に利用可能な Amazon S3 バケットのフライトデータを含むテーブルを定義します。わかりやすくするために、データのいくつかの列と 1 つのパーティションキーのみが定義されています。4 つのパーティションも Data Catalog に定義されています。基本データのストレージを記述するいくつかのフィールドも StorageDescriptor フィールドに示されています。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CloudFormation template in YAML to demonstrate creating a database, a table,
  and partitions
# The metadata created in the Data Catalog points to the flights public S3 bucket
#
# Parameters substituted in the Resources section
# These parameters are names of the resources created in the Data Catalog
Parameters:
  CFNDatabaseName:
    Type: String
    Default: cfn-database-flights-1
  CFNTableName1:
    Type: String
    Default: cfn-manual-table-flights-1
# Resources to create metadata in the Data Catalog
Resources:
###
# Create an AWS Glue database
CFNDatabaseFlights:
  Type: AWS::Glue::Database
  Properties:
    CatalogId: !Ref AWS::AccountId
    DatabaseInput:
      Name: !Ref CFNDatabaseName
      Description: Database to hold tables for flights data
###
# Create an AWS Glue table
CFNTableFlights:
  # Creating the table waits for the database to be created
  DependsOn: CFNDatabaseFlights
  Type: AWS::Glue::Table
  Properties:
    CatalogId: !Ref AWS::AccountId
    DatabaseName: !Ref CFNDatabaseName
    TableInput:
      Name: !Ref CFNTableName1
```

```

    Description: Define the first few columns of the flights table
    TableType: EXTERNAL_TABLE
    Parameters: {
"classification": "csv"
}
#
    ViewExpandedText: String
    PartitionKeys:
    # Data is partitioned by month
    - Name: mon
      Type: bigint
    StorageDescriptor:
    OutputFormat: org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat
    Columns:
    - Name: year
      Type: bigint
    - Name: quarter
      Type: bigint
    - Name: month
      Type: bigint
    - Name: day_of_month
      Type: bigint
    InputFormat: org.apache.hadoop.mapred.TextInputFormat
    Location: s3://crawler-public-us-east-1/flight/2016/csv/
    SerdeInfo:
    Parameters:
    field.delim: ","
    SerializationLibrary: org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
# Partition 1
# Create an AWS Glue partition
CFNPartitionMon1:
    DependsOn: CFNTableFlights
    Type: AWS::Glue::Partition
    Properties:
    CatalogId: !Ref AWS::AccountId
    DatabaseName: !Ref CFNDatabaseName
    TableName: !Ref CFNTableName1
    PartitionInput:
    Values:
    - 1
    StorageDescriptor:
    OutputFormat: org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat
    Columns:
    - Name: mon
      Type: bigint

```

```
    InputFormat: org.apache.hadoop.mapred.TextInputFormat
    Location: s3://crawler-public-us-east-1/flight/2016/csv/mon=1/
    SerdeInfo:
      Parameters:
        field.delim: ","
      SerializationLibrary: org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
# Partition 2
# Create an AWS Glue partition
CFNPartitionMon2:
  DependsOn: CFNTableFlights
  Type: AWS::Glue::Partition
  Properties:
    CatalogId: !Ref AWS::AccountId
    DatabaseName: !Ref CFNDatabaseName
    TableName: !Ref CFNTableName1
    PartitionInput:
      Values:
        - 2
      StorageDescriptor:
        OutputFormat: org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat
        Columns:
          - Name: mon
            Type: bigint
        InputFormat: org.apache.hadoop.mapred.TextInputFormat
        Location: s3://crawler-public-us-east-1/flight/2016/csv/mon=2/
        SerdeInfo:
          Parameters:
            field.delim: ","
          SerializationLibrary: org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
# Partition 3
# Create an AWS Glue partition
CFNPartitionMon3:
  DependsOn: CFNTableFlights
  Type: AWS::Glue::Partition
  Properties:
    CatalogId: !Ref AWS::AccountId
    DatabaseName: !Ref CFNDatabaseName
    TableName: !Ref CFNTableName1
    PartitionInput:
      Values:
        - 3
      StorageDescriptor:
        OutputFormat: org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat
        Columns:
```

```
- Name: mon
  Type: bigint
  InputFormat: org.apache.hadoop.mapred.TextInputFormat
  Location: s3://crawler-public-us-east-1/flight/2016/csv/mon=3/
  SerdeInfo:
    Parameters:
      field.delim: ","
    SerializationLibrary: org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
# Partition 4
# Create an AWS Glue partition
CFNPartitionMon4:
  DependsOn: CFNTableFlights
  Type: AWS::Glue::Partition
  Properties:
    CatalogId: !Ref AWS::AccountId
    DatabaseName: !Ref CFNDatabaseName
    TableName: !Ref CFNTableName1
    PartitionInput:
      Values:
        - 4
    StorageDescriptor:
      OutputFormat: org.apache.hadoop.hive ql.io.HiveIgnoreKeyTextOutputFormat
      Columns:
        - Name: mon
          Type: bigint
      InputFormat: org.apache.hadoop.mapred.TextInputFormat
      Location: s3://crawler-public-us-east-1/flight/2016/csv/mon=4/
      SerdeInfo:
        Parameters:
          field.delim: ","
        SerializationLibrary: org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe
```

AWS Glue Grok 分類子用のサンプル AWS CloudFormation テンプレート

AWS Glue 分類子はデータのスキーマを決定します。1つのタイプのカスタム分類子では、grok パターンを使用してデータをマッチングします。パターンがマッチすると、カスタム分類子ではテーブルのスキーマを作成し、分類子の定義に設定された値に classification を設定します。

このサンプルで作成する分類子では、message という列が 1 つあるスキーマを作成し、分類を greedy に設定します。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a classifier
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:

# The name of the classifier to be created
CFNClassifierName:
  Type: String
  Default: cfn-classifier-grok-one-column-1

#
#
# Resources section defines metadata for the Data Catalog
Resources:
# Create classifier that uses grok pattern to put all data in one column and classifies
it as "greedy".
CFNClassifierFlights:
  Type: AWS::Glue::Classifier
  Properties:
    GrokClassifier:
      #Grok classifier that puts all data in one column
      Name: !Ref CFNClassifierName
      Classification: greedy

      GrokPattern: "%{GREEDYDATA:message}"
      #CustomPatterns: none
```

AWS Glue JSON 分類子用のサンプル AWS CloudFormation テンプレート

AWS Glue 分類子はデータのスキーマを決定します。1 つのタイプのカスタム分類子では、分類のための分類子の JSON データを定義する JsonPath 文字列を使用します。「[Writing JsonPath Custom Classifiers](#)」に記載されているように、AWS Glue では JsonPath の演算子のサブセットがサポートされています。

パターンが一致すると、カスタム識別子を使用してテーブルのスキーマが作成されます。

このサンプルでは、オブジェクトの `Records3` 配列で各スキーマを使用してスキーマを作成する識別子を作成します。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a JSON classifier
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:

# The name of the classifier to be created
CFNClassifierName:
  Type: String
  Default: cfn-classifier-json-one-column-1

#
#
# Resources section defines metadata for the Data Catalog
Resources:
# Create classifier that uses a JSON pattern.
CFNClassifierFlights:
  Type: AWS::Glue::Classifier
  Properties:
    JSONClassifier:
      #JSON classifier
      Name: !Ref CFNClassifierName
      JsonPath: $.Records3[*]
```

AWS Glue XML 分類子用のサンプル AWS CloudFormation テンプレート

AWS Glue 分類子はデータのスキーマを決定します。1つのタイプのカスタム分類子は、XML タグを指定して、解析中の XML ドキュメントの各レコードを含む要素を指定します。パターンがマッチすると、カスタム分類子ではテーブルのスキーマを作成し、分類子の定義に設定された値に `classification` を設定します。

このサンプルで作成する分類子では、Record タグの各レコードでスキーマを作成し、分類を XML に設定します。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating an XML classifier
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:

# The name of the classifier to be created
CFNClassifierName:
  Type: String
  Default: cfn-classifier-xml-one-column-1

#
#
# Resources section defines metadata for the Data Catalog
Resources:
# Create classifier that uses the XML pattern and classifies it as "XML".
CFNClassifierFlights:
  Type: AWS::Glue::Classifier
  Properties:
    XMLClassifier:
      #XML classifier
      Name: !Ref CFNClassifierName
      Classification: XML
      RowTag: <Records>
```

Amazon S3 の AWS Glue クローラー用のサンプル AWS CloudFormation テンプレート

AWS Glue クローラーでは、データに対応するメタデータテーブルをデータカタログに作成します。次に、これらのテーブル定義を ETL ジョブのソースおよびターゲットとして使用できます。

このサンプルでは、クローラー、必要な IAM ロール、および AWS Glue データベースをデータカタログに作成します。このクローラーを実行すると、クローラーは IAM ロールを取得し、公開フライトデータ用のテーブルをデータベースに作成します。テーブルは、プレフィックス `cfn_sample_1_` を使用して作成されます。このテンプレートで作成された IAM ロールでは、グローバルのアクセス許可が付与されるので、カスタムロールを作成する必要がある場合もあります。

この分類子で定義されるカスタム分類子はありません。AWS Glue の組み込み分類子がデフォルトで使用されます。

このサンプルを AWS CloudFormation コンソールに送信する場合は、IAM ロールを作成することを確認する必要があります。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a crawler
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:

# The name of the crawler to be created
CFNCrawlerName:
  Type: String
  Default: cfn-crawler-flights-1
CFNDatabaseName:
  Type: String
  Default: cfn-database-flights-1
CFNTablePrefixName:
  Type: String
  Default: cfn_sample_1_

#
#
# Resources section defines metadata for the Data Catalog
Resources:
#Create IAM Role assumed by the crawler. For demonstration, this role is given all
permissions.
CFNRoleFlights:
  Type: AWS::IAM::Role
  Properties:
    AssumeRolePolicyDocument:
      Version: "2012-10-17"
      Statement:
        -
          Effect: "Allow"
          Principal:
            Service:
              - "glue.amazonaws.com"
          Action:
```

```

    - "sts:AssumeRole"
  Path: "/"
  Policies:
    -
      PolicyName: "root"
      PolicyDocument:
        Version: "2012-10-17"
        Statement:
          -
            Effect: "Allow"
            Action: "*"
            Resource: "*"
# Create a database to contain tables created by the crawler
CFNDatabaseFlights:
  Type: AWS::Glue::Database
  Properties:
    CatalogId: !Ref AWS::AccountId
    DatabaseInput:
      Name: !Ref CFNDatabaseName
      Description: "AWS Glue container to hold metadata tables for the flights
crawler"
#Create a crawler to crawl the flights data on a public S3 bucket
CFNCrawlerFlights:
  Type: AWS::Glue::Crawler
  Properties:
    Name: !Ref CFNCrawlerName
    Role: !GetAtt CFNRoleFlights.Arn
    #Classifiers: none, use the default classifier
    Description: AWS Glue crawler to crawl flights data
    #Schedule: none, use default run-on-demand
    DatabaseName: !Ref CFNDatabaseName
    Targets:
      S3Targets:
        # Public S3 bucket with the flights data
        - Path: "s3://crawler-public-us-east-1/flight/2016/csv"
    TablePrefix: !Ref CFNTablePrefixName
    SchemaChangePolicy:
      UpdateBehavior: "UPDATE_IN_DATABASE"
      DeleteBehavior: "LOG"
    Configuration: "{\"Version\":1.0,\"CrawlerOutput\":{\"Partitions\":
{\"AddOrUpdateBehavior\":\"InheritFromTable\"},\"Tables\":{\"AddOrUpdateBehavior\":
\"MergeNewColumns\"}}}"

```

AWS Glue 接続用のサンプル AWS CloudFormation テンプレート

Data Catalog の AWS Glue 接続には、JDBC データベースに接続するために必要な JDBC およびネットワーク情報が含まれています。この情報は、JDBC データベースに接続して ETL ジョブをクローलまたは実行するときに使用されます。

このサンプルでは、Amazon RDS MySQL データベース (devdb) への接続を作成します。この接続を使用する場合は、IAM ロール、データベース認証情報、およびネットワーク接続の値も指定する必要があります。テンプレートの必須フィールドの詳細を参照してください。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a connection
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:

# The name of the connection to be created
  CFNConnectionName:
    Type: String
    Default: cfn-connection-mysql-flights-1
  CFNJDBCString:
    Type: String
    Default: "jdbc:mysql://xxx-mysql.yyyyyyyyyyyyyyy.us-east-1.rds.amazonaws.com:3306/
devdb"
  CFNJDBCUser:
    Type: String
    Default: "master"
  CFNJDBCPassword:
    Type: String
    Default: "12345678"
    NoEcho: true
#
#
# Resources section defines metadata for the Data Catalog
Resources:
  CFNConnectionMySQL:
    Type: AWS::Glue::Connection
    Properties:
      CatalogId: !Ref AWS::AccountId
```

```
ConnectionInput:
  Description: "Connect to MySQL database."
  ConnectionType: "JDBC"
  #MatchCriteria: none
  PhysicalConnectionRequirements:
    AvailabilityZone: "us-east-1d"
    SecurityGroupIdList:
      - "sg-7d52b812"
    SubnetId: "subnet-84f326ee"
  ConnectionProperties: {
    "JDBC_CONNECTION_URL": !Ref CFNJDBCString,
    "USERNAME": !Ref CFNJDBCUser,
    "PASSWORD": !Ref CFNJDBCPassword
  }
  Name: !Ref CFNConnectionName
```

JDBC の AWS Glue クローラー用のサンプル AWS CloudFormation テンプレート

AWS Glue クローラーでは、データに対応するメタデータテーブルをデータカタログに作成します。次に、これらのテーブル定義を ETL ジョブのソースおよびターゲットとして使用できます。

このサンプルでは、クローラー、必要な IAM ロール、および AWS Glue データベースをデータカタログに作成します。このクローラーを実行すると、クローラーは IAM ロールを取得し、MySQL データベースに保存されている公開フライトデータ用のテーブルをデータベースに作成します。テーブルは、プレフィックス `cfn_jdbc_1_` を使用して作成されます。このテンプレートで作成された IAM ロールでは、グローバルのアクセス許可が付与されるので、カスタムロールを作成する必要がある場合もあります。JDBC データに対してはカスタム分類子を定義できません。AWS Glue の組み込み分類子がデフォルトで使用されます。

このサンプルを AWS CloudFormation コンソールに送信する場合は、IAM ロールを作成することを確認する必要があります。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a crawler
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
```

```
Parameters:

# The name of the crawler to be created
CFNCrawlerName:
  Type: String
  Default: cfn-crawler-jdbc-flights-1
# The name of the database to be created to contain tables
CFNDatabaseName:
  Type: String
  Default: cfn-database-jdbc-flights-1
# The prefix for all tables crawled and created
CFNTableNamePrefix:
  Type: String
  Default: cfn_jdbc_1_
# The name of the existing connection to the MySQL database
CFNConnectionName:
  Type: String
  Default: cfn-connection-mysql-flights-1
# The name of the JDBC path (database/schema/table) with wildcard (%) to crawl
CFNJDBCPath:
  Type: String
  Default: saldev/%

#
#
# Resources section defines metadata for the Data Catalog
Resources:
#Create IAM Role assumed by the crawler. For demonstration, this role is given all
permissions.
CFNRoleFlights:
  Type: AWS::IAM::Role
  Properties:
    AssumeRolePolicyDocument:
      Version: "2012-10-17"
      Statement:
        -
          Effect: "Allow"
          Principal:
            Service:
              - "glue.amazonaws.com"
          Action:
            - "sts:AssumeRole"
    Path: "/"
    Policies:
      -
```

```

    PolicyName: "root"
    PolicyDocument:
      Version: "2012-10-17"
      Statement:
        -
          Effect: "Allow"
          Action: "*"
          Resource: "*"
# Create a database to contain tables created by the crawler
CFNDatabaseFlights:
  Type: AWS::Glue::Database
  Properties:
    CatalogId: !Ref AWS::AccountId
    DatabaseInput:
      Name: !Ref CFNDatabaseName
      Description: "AWS Glue container to hold metadata tables for the flights
crawler"
#Create a crawler to crawl the flights data in MySQL database
CFNCrawlerFlights:
  Type: AWS::Glue::Crawler
  Properties:
    Name: !Ref CFNCrawlerName
    Role: !GetAtt CFNRoleFlights.Arn
    #Classifiers: none, use the default classifier
    Description: AWS Glue crawler to crawl flights data
    #Schedule: none, use default run-on-demand
    DatabaseName: !Ref CFNDatabaseName
    Targets:
      JdbcTargets:
        # JDBC MySQL database with the flights data
        - ConnectionName: !Ref CFNConnectionName
          Path: !Ref CFNJDBCPath
        #Exclusions: none
    TablePrefix: !Ref CFNTablePrefixName
    SchemaChangePolicy:
      UpdateBehavior: "UPDATE_IN_DATABASE"
      DeleteBehavior: "LOG"
    Configuration: "{\"Version\":1.0,\"CrawlerOutput\":{\"Partitions\":
{\"AddOrUpdateBehavior\": \"InheritFromTable\"},\"Tables\":{\"AddOrUpdateBehavior\":
\"MergeNewColumns\"}}}"

```

Amazon S3 から Amazon S3 への AWS Glue ジョブのサンプル AWS CloudFormation テンプレート

Data Catalog の AWS Glue ジョブには、AWS Glue でスクリプトを実行するために必要なパラメータ値が含まれています。

このサンプルで作成するジョブでは、Amazon S3 バケットのフライトデータを csv 形式で読み取り、Amazon S3 の Parquet ファイルに書き込みます。このジョブで実行するスクリプトは既存している必要があります。環境に応じた ETL スクリプトを AWS Glue コンソールで生成できます。このジョブ実行時に、適切なアクセス許可が設定された IAM ロールも指定する必要があります。

テンプレートには、一般的なパラメータ値が示されています。たとえば、AllocatedCapacity (DPU) はデフォルトで 5 になります。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a job using the public flights S3 table in a
public bucket
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:

# The name of the job to be created
CFNJobName:
  Type: String
  Default: cfn-job-S3-to-S3-2
# The name of the IAM role that the job assumes. It must have access to data, script,
temporary directory
CFNIAMRoleName:
  Type: String
  Default: AWSGlueServiceRoleGA
# The S3 path where the script for this job is located
CFNScriptLocation:
  Type: String
  Default: s3://aws-glue-scripts-123456789012-us-east-1/myid/sal-job-test2
#
#
# Resources section defines metadata for the Data Catalog
Resources:
```

```
# Create job to run script which accesses flightscsv table and write to S3 file as
parquet.
# The script already exists and is called by this job
CFNJobFlights:
  Type: AWS::Glue::Job
  Properties:
    Role: !Ref CFNIAMRoleName
    #DefaultArguments: JSON object
    # If script written in Scala, then set DefaultArguments={'--job-language';
'scala', '--class': 'your scala class'}
    #Connections: No connection needed for S3 to S3 job
    # ConnectionsList
    #MaxRetries: Double
    Description: Job created with CloudFormation
    #LogUri: String
    Command:
      Name: glueetl
      ScriptLocation: !Ref CFNScriptLocation
        # for access to directories use proper IAM role with permission to buckets
and folders that begin with "aws-glue-"
        # script uses temp directory from job definition if required (temp
directory not used S3 to S3)
        # script defines target for output as s3://aws-glue-target/sal
    AllocatedCapacity: 5
    ExecutionProperty:
      MaxConcurrentRuns: 1
    Name: !Ref CFNJobName
```

Amazon S3 に書き込む JDBC の AWS Glue ジョブ用のサンプル AWS CloudFormation テンプレート

Data Catalog の AWS Glue ジョブには、AWS Glue でスクリプトを実行するために必要なパラメータ値が含まれています。

このサンプルで作成するジョブでは、cfn-connection-mysql-flights-1 という接続で定義された MySQL JDBC データベースからフライトデータを読み取り、Amazon S3 の Parquet ファイルに書き込みます。このジョブで実行するスクリプトは既存する必要があります。環境に応じた ETL スクリプトを AWS Glue コンソールで生成できます。このジョブ実行時に、適切なアクセス許可が設定された IAM ロールも指定する必要があります。

テンプレートには、一般的なパラメータ値が示されています。たとえば、AllocatedCapacity (DPU) はデフォルトで 5 になります。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a job using a MySQL JDBC DB with the flights
# data to an S3 file
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:

# The name of the job to be created
  CFNJobName:
    Type: String
    Default: cfn-job-JDBC-to-S3-1
# The name of the IAM role that the job assumes. It must have access to data, script,
# temporary directory
  CFNIAMRoleName:
    Type: String
    Default: AWSGlueServiceRoleGA
# The S3 path where the script for this job is located
  CFNScriptLocation:
    Type: String
    Default: s3://aws-glue-scripts-123456789012-us-east-1/myid/sal-job-dec4a
# The name of the connection used for JDBC data source
  CFNConnectionName:
    Type: String
    Default: cfn-connection-mysql-flights-1
#
#
# Resources section defines metadata for the Data Catalog
Resources:
# Create job to run script which accesses JDBC flights table via a connection and write
# to S3 file as parquet.
# The script already exists and is called by this job
  CFNJobFlights:
    Type: AWS::Glue::Job
    Properties:
      Role: !Ref CFNIAMRoleName
      #DefaultArguments: JSON object
```

```
# For example, if required by script, set temporary directory as
DefaultArguments={'--TempDir'; 's3://aws-glue-temporary-xyz/sal'}
Connections:
  Connections:
    - !Ref CFNConnectionName
#MaxRetries: Double
Description: Job created with CloudFormation using existing script
#LogUri: String
Command:
  Name: glueetl
  ScriptLocation: !Ref CFNScriptLocation
    # for access to directories use proper IAM role with permission to buckets
and folders that begin with "aws-glue-"
    # if required, script defines temp directory as argument TempDir and used
in script like redshift_tmp_dir = args["TempDir"]
    # script defines target for output as s3://aws-glue-target/sal
AllocatedCapacity: 5
ExecutionProperty:
  MaxConcurrentRuns: 1
Name: !Ref CFNJobName
```

AWS Glue オンデマンドトリガー用のサンプル AWS CloudFormation テンプレート

Data Catalog の AWS Glue トリガーには、トリガーの発動によってジョブ実行を開始するために必要なパラメータ値が含まれています。オンデマンドトリガーは、このトリガーを有効にしたときに発生します。

このサンプルで作成するオンデマンドトリガーでは、`cfn-job-S3-to-S3-1` という 1 つのジョブを開始します。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating an on-demand trigger
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:
  # The existing job to be started by this trigger
  CFNJobName:
```

```
Type: String
Default: cfn-job-S3-to-S3-1
# The name of the trigger to be created
CFNTriggerName:
  Type: String
  Default: cfn-trigger-ondemand-flights-1
#
# Resources section defines metadata for the Data Catalog
# Sample CFN YAML to demonstrate creating an on-demand trigger for a job
Resources:
# Create trigger to run an existing job (CFNJobName) on an on-demand schedule.
CFNTriggerSample:
  Type: AWS::Glue::Trigger
  Properties:
    Name:
      Ref: CFNTriggerName
    Description: Trigger created with CloudFormation
    Type: ON_DEMAND
    Actions:
      - JobName: !Ref CFNJobName
      # Arguments: JSON object
    #Schedule:
    #Predicate:
```

AWS Glue のスケジュールされたトリガー用のサンプル AWS CloudFormation テンプレート

Data Catalog の AWS Glue トリガーには、トリガーの発動によってジョブ実行を開始するために必要なパラメータ値が含まれています。スケジュールされたトリガーは、このトリガーを有効にして cron タイマーがポップすると、発生します。

このサンプルで作成するスケジュールされたトリガーでは、cfn-job-S3-to-S3-1 という 1 つのジョブを開始します。このタイマーは、平日の 10 分ごとにジョブを実行する cron 式です。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a scheduled trigger
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
```

```
Parameters:
# The existing job to be started by this trigger
CFNJobName:
  Type: String
  Default: cfn-job-S3-to-S3-1
# The name of the trigger to be created
CFNTriggerName:
  Type: String
  Default: cfn-trigger-scheduled-flights-1
#
# Resources section defines metadata for the Data Catalog
# Sample CFN YAML to demonstrate creating a scheduled trigger for a job
#
Resources:
# Create trigger to run an existing job (CFNJobName) on a cron schedule.
TriggerSample1CFN:
  Type: AWS::Glue::Trigger
  Properties:
    Name:
      Ref: CFNTriggerName
    Description: Trigger created with CloudFormation
    Type: SCHEDULED
    Actions:
      - JobName: !Ref CFNJobName
        # Arguments: JSON object
        ## Run the trigger every 10 minutes on Monday to Friday
        Schedule: cron(0/10 * ? * MON-FRI *)
    #Predicate:
```

AWS Glue の条件付きトリガー用のサンプル AWS CloudFormation テンプレート

Data Catalog の AWS Glue トリガーには、トリガーの発動によってジョブ実行を開始するために必要なパラメータ値が含まれています。条件付きトリガーは、このトリガーを有効にして、その条件が満たされる (例: ジョブが正常に完了する) と、発生します。

このサンプルで作成する条件付きトリガーでは、cfn-job-S3-to-S3-1 という 1 つのジョブを開始します。このジョブは、cfn-job-S3-to-S3-2 というジョブが正常に完了すると、開始されます。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a conditional trigger for a job, which starts
  when another job completes
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:
  # The existing job to be started by this trigger
  CFNJobName:
    Type: String
    Default: cfn-job-S3-to-S3-1
  # The existing job that when it finishes causes trigger to fire
  CFNJobName2:
    Type: String
    Default: cfn-job-S3-to-S3-2
  # The name of the trigger to be created
  CFNTriggerName:
    Type: String
    Default: cfn-trigger-conditional-1
#
Resources:
# Create trigger to run an existing job (CFNJobName) when another job completes
  (CFNJobName2).
  CFNTriggerSample:
    Type: AWS::Glue::Trigger
    Properties:
      Name:
        Ref: CFNTriggerName
      Description: Trigger created with CloudFormation
      Type: CONDITIONAL
      Actions:
        - JobName: !Ref CFNJobName
          # Arguments: JSON object
      #Schedule: none
      Predicate:
        #Value for Logical is required if more than 1 job listed in Conditions
        Logical: AND
        Conditions:
          - LogicalOperator: EQUALS
            JobName: !Ref CFNJobName2
            State: SUCCEEDED
```

AWS Glue の開発エンドポイント用のサンプル AWS CloudFormation テンプレート

AWS Glue 機械学習変換は、データをクレンジングするためのカスタム変換です。現在、FindMatches という名前の変換が 1 つあります。FindMatches 変換を使用すると、レコードに共通の一意の識別子がなく、正確に一致するフィールドがない場合でも、データセット内の重複レコードまたは一致するレコードを識別できます。

このサンプルでは、機械学習変換を作成します。機械学習変換の作成に必要なパラメータの詳細については、「[AWS Lake Formation FindMatches によるレコードのマッチング](#)」を参照してください。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a machine learning transform
#
# Resources section defines metadata for the machine learning transform
Resources:
  MyMLTransform:
    Type: "AWS::Glue::MLTransform"
    Condition: "isGlueMLGARegion"
    Properties:
      Name: !Sub "MyTransform"
      Description: "The bestest transform ever"
      Role: !ImportValue MyMLTransformUserRole
      GlueVersion: "1.0"
      WorkerType: "Standard"
      NumberOfWorkers: 5
      Timeout: 120
      MaxRetries: 1
      InputRecordTables:
        GlueTables:
          - DatabaseName: !ImportValue MyMLTransformDatabase
            TableName: !ImportValue MyMLTransformTable
      TransformParameters:
        TransformType: "FIND_MATCHES"
      FindMatchesParameters:
        PrimaryKeyColumnName: "testcolumn"
        PrecisionRecallTradeoff: 0.5
        AccuracyCostTradeoff: 0.5
        EnforceProvidedLabels: True
    Tags:
```

```

    key1: "value1"
    key2: "value2"
  TransformEncryption:
    TaskRunSecurityConfigurationName: !ImportValue
MyMLTransformSecurityConfiguration
  MLUserDataEncryption:
    MLUserDataEncryptionMode: "SSE-KMS"
    KmsKeyId: !ImportValue MyMLTransformEncryptionKey

```

AWS Glue Data Quality ルールセット用のサンプル AWS CloudFormation テンプレート

AWS Glue Data Quality ルールセットには、データカタログ内のテーブルで評価できるルールが含まれています。ルールセットをターゲットテーブルに配置したら、データカタログにアクセスして、ルールセット内のルールに対してデータを実行する評価を実行できます。これらのルールは、行数の評価からデータの参照整合性の評価までさまざまです。

次のサンプルは、指定されたターゲットテーブルにさまざまなルールを含むルールセットを作成する CloudFormation テンプレートです。

```

AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a DataQualityRuleset
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:

  # The name of the ruleset to be created
  RulesetName:
    Type: String
    Default: "CFNRulesetName"
  RulesetDescription:
    Type: String
    Default: "CFN DataQualityRuleset"
# Rules that will be associated with this ruleset
Rules:
  Type: String
  Default: 'Rules = [
    RowCount > 100,
    IsUnique "id",
    IsComplete "nametype"

```

```
    ]'  
    # Name of database and table within Data Catalog which the ruleset will  
    # be applied too  
    DatabaseName:  
      Type: String  
      Default: "ExampleDatabaseName"  
    TableName:  
      Type: String  
      Default: "ExampleTableName"  
  
# Resources section defines metadata for the Data Catalog  
Resources:  
  # Creates a Data Quality ruleset under specified rules  
  DQRuleset:  
    Type: AWS::Glue::DataQualityRuleset  
    Properties:  
      Name: !Ref RulesetName  
      Description: !Ref RulesetDescription  
      # The String within rules must be formatted in DQDL, a language  
      # used specifically to make rules  
      Ruleset: !Ref Rules  
      # The targeted table must exist within Data Catalog alongside  
      # the correct database  
      TargetTable:  
        DatabaseName: !Ref DatabaseName  
        TableName: !Ref TableName
```

EventBridge スケジューラを使用する AWS Glue Data Quality ルールセット用の AWS CloudFormation テンプレート

AWS Glue Data Quality ルールセットには、データカタログ内のテーブルで評価できるルールが含まれています。ルールセットをターゲットテーブルに配置したら、データカタログにアクセスして、ルールセット内のルールに対してデータを実行する評価を実行できます。手動でデータカタログにアクセスしてルールセットを評価する代わりに、CloudFormation テンプレート内に EventBridge スケジューラを追加して、これらのルールセットの評価を一定間隔でスケジュールすることもできます。

次のサンプルは、前述のルールセットを 5 分ごとに評価する、Data Quality ルールセットと EventBridge スケジューラを作成する CloudFormation テンプレートです。

```
AWSTemplateFormatVersion: '2010-09-09'  
# Sample CFN YAML to demonstrate creating a DataQualityRuleset
```

```
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:

# The name of the ruleset to be created
RulesetName:
  Type: String
  Default: "CFNRulesetName"
# Rules that will be associated with this Ruleset
Rules:
  Type: String
  Default: 'Rules = [
    RowCount > 100,
    IsUnique "id",
    IsComplete "nametype"
  ]'
# The name of the Schedule to be created
ScheduleName:
  Type: String
  Default: "ScheduleDQRulsetEvaluation"
# This expression determines the rate at which the Schedule will evaluate
# your data using the above ruleset
ScheduleRate:
  Type: String
  Default: "rate(5 minutes)"
# The Request that being sent must match the details of the Data Quality Ruleset
ScheduleRequest:
  Type: String
  Default: '
    { "DataSource": { "GlueTable": { "DatabaseName": "ExampleDatabaseName",
      "TableName": "ExampleTableName" } },
      "Role": "role/AWSGlueServiceRoleDefault",
      "RulesetNames": [ ""CFNRulesetName"" ] }
    ,

# Resources section defines metadata for the Data Catalog
Resources:
# Creates a Data Quality ruleset under specified rules
DQRuleset:
  Type: AWS::Glue::DataQualityRuleset
  Properties:
    Name: !Ref RulesetName
    Description: "CFN DataQualityRuleset"
```

```
# The String within rules must be formatted in DQDL, a language
# used specifically to make rules
Ruleset: !Ref Rules
# The targeted table must exist within Data Catalog alongside
# the correct database
TargetTable:
  DatabaseName: "ExampleDatabaseName"
  TableName: "ExampleTableName"
# Create a Scheduler to schedule evaluation runs on the above ruleset
ScheduleDQEval:
  Type: AWS::Scheduler::Schedule
  Properties:
    Name: !Ref ScheduleName
    Description: "Schedule DataQualityRuleset Evaluations"
    FlexibleTimeWindow:
      Mode: "OFF"
    ScheduleExpression: !Ref ScheduleRate
    ScheduleExpressionTimezone: "America/New_York"
    State: "ENABLED"
    Target:
      # The ARN is the API that will be run, since we want to evaluate our ruleset
      # we want this specific ARN
      Arn: "arn:aws:scheduler::aws-sdk:glue:startDataQualityRulesetEvaluationRun"
      # Your RoleArn must have approval to schedule
      RoleArn: "arn:aws:iam::123456789012:role/AWSGlueServiceRoleDefault"
      # This is the Request that is being sent to the Arn
      Input: '
        { "DataSource": { "GlueTable": { "DatabaseName": "sampledb", "TableName":
"meteorite" } },
          "Role": "role/AWSGlueServiceRoleDefault",
          "RulesetNames": [ "TestCFN" ] }
      '
```

AWS Glue の開発エンドポイント用のサンプル AWS CloudFormation テンプレート

AWS Glue の開発エンドポイントは、AWS Glue スクリプトの開発およびテストに使用できる環境です。

このサンプルで作成する開発エンドポイントでは、正常な作成に最低限必要なネットワークパラメータ値を使用します。開発エンドポイントの設定に必要なパラメータの詳細については、「[AWS Glue のための開発用ネットワークの設定](#)」を参照してください。

開発エンドポイントを作成するには、既存の IAM ロール ARN (Amazon リソースネーム) を指定します。開発エンドポイントでノートブックサーバーを作成する場合は、有効な RSA パブリックキーを指定し、対応するプライベートキーを使用可能な状態に保持します。

Note

作成した開発エンドポイントに関連付けられているすべてのノートブックサーバーを管理します。したがって、開発エンドポイントを削除した場合、ノートブックサーバーを削除するには AWS CloudFormation コンソールで AWS CloudFormation スタックを削除する必要があります。

```
---
AWSTemplateFormatVersion: '2010-09-09'
# Sample CFN YAML to demonstrate creating a development endpoint
#
# Parameters section contains names that are substituted in the Resources section
# These parameters are the names the resources created in the Data Catalog
Parameters:

# The name of the crawler to be created
CFNEndpointName:
  Type: String
  Default: cfn-devendpoint-1
CFNIAMRoleArn:
  Type: String
  Default: arn:aws:iam::123456789012/role/AWSGlueServiceRoleGA
#
#
# Resources section defines metadata for the Data Catalog
Resources:
  CFNDevEndpoint:
    Type: AWS::Glue::DevEndpoint
    Properties:
      EndpointName: !Ref CFNEndpointName
      #ExtraJarsS3Path: String
```

```
#ExtraPythonLibsS3Path: String
NumberOfNodes: 5
PublicKey: ssh-rsa public.....key myuserid-key
RoleArn: !Ref CFNIAMRoleArn
SecurityGroupIds:
  - sg-64986c0b
SubnetId: subnet-c67cccac
```

AWS Glue プログラミングガイド

スクリプトには、ソースからデータを抽出して、変換し、ターゲットにロードするコードが含まれています。AWS Glue はジョブを開始するときにスクリプトを実行します。

AWS Glue ETL スクリプトは Python または Scala で記述されます。すべてのジョブタイプは Python で記述できますが、AWS Glue for Spark ジョブは Scala でも記述できます。AWS Glue Studio でジョブのソースコードロジックを自動で生成するときに、スクリプトが作成されます。このスクリプトを編集するか、または、独自のスクリプトを指定して ETL 作業を処理することができます。

独自のカスタムスクリプトの提供

スクリプトは、AWS Glue で抽出、変換、ロード (ETL) 作業を実行します。スクリプトは、自動でジョブのソースコードロジックを生成するときに作成されます。この生成されたスクリプトを編集することもできますし、独自のカスタムスクリプトを指定することもできます。

AWS Glue で独自のカスタムスクリプトを提供するには、以下の一般的な手順に従います。

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. [ETL Jobs] タブをクリックしてから、[ジョブを作成する] セクションを開きます。[script editor] オプションをクリックします。
3. [This job runs] (このジョブ実行) で、次のいずれかを選択します。
 - [Create a new script with boilerplate code]
 - [Upload and edit an existing script]
4. [ジョブの詳細] ページで、カスタムスクリプトを実行するのに必要な [IAM ロール] をクリックします。詳細については、「[AWS Glue の Identity and Access Management](#)」を参照してください。
5. スクリプトが参照する接続を選択します。これらのオブジェクトは、目的の JDBC データストアに接続するために必要です。

Elastic Network Interface は、Virtual Private Cloud (VPC) でインスタンスにアタッチできる、仮想ネットワークインターフェイスです。スクリプトで使用されているデータストアに接続するために必要な Elastic Network Interface を選択します。

6. ジョブタイプ固有のパラメータなどの追加設定を行います。ジョブタイプの設定の詳細については、「[AWS Glue Studio でビジュアル ETL ジョブを作成する](#)」セクションを参照してください。
7. [スクリプト] タブで、カスタムスクリプトを貼り付けるか、記述します。

カスタムスクリプトを記述するプロセスのガイドとして、このセクションの内容を使用してください。

AWS Glue におけるジョブ追加の詳細については、「[AWS Glue Studio でビジュアル ETL ジョブを作成する](#)」を参照してください。

詳細な手順については、コンソールの [Add jobAWS Glue] (ジョブ追加) チュートリアルを参照してください。

Spark スクリプトのプログラミング

AWS Glue では、テストや実行に加えて、抽出、変換、およびロード (ETL) スクリプトの書き込みや自動生成を簡単に行えます。このセクションでは、AWS Glue に導入された Apache Spark の拡張機能について説明し、ETL スクリプトを Python と Scala で記述し実行する方法の例を示します。

Important

AWS Glue の異なるバージョンでは、Apache Spark の異なるバージョンがサポートされます。カスタムスクリプトは、サポートされている Apache Spark バージョンと互換性がある必要があります。AWS Glue のバージョンの詳細に関しては、「[Glue version job property](#)」を参照してください。

トピック

- [チュートリアル: AWS Glue for Spark スクリプトの作成](#)
- [で AWS Glue ETL スクリプトをプログラムする PySpark](#)
- [Scala での AWS Glue ETL スクリプトのプログラミング](#)
- [AWS Glue for Spark ETL スクリプトのプログラミングの機能と最適化](#)

チュートリアル: AWS Glue for Spark スクリプトの作成

このチュートリアルでは、AWS Glue スクリプトを作成するプロセスを紹介します。スクリプトを、ジョブを使用してスケジュールに従って実行したり、インタラクティブセッションを使用してインタラクティブに実行したりできます。ジョブの詳細については、「[AWS Glue Studio でビジュアル ETL ジョブを作成する](#)」を参照してください。詳細については、「[the section called “AWS Glue インタラクティブセッションの概要”](#)」を参照してください。

AWS Glue Studio のビジュアルエディターには、AWS Glue ジョブを作成するためのグラフィカルでコード不要のインターフェイスが用意されています。AWS スクリプトをビジュアルジョブに貼り付けます。Glue スクリプトを使用すると、Apache Spark プログラムを操作するための拡張された一連のツールにアクセスできます。ネイティブの Spark API だけでなく、AWS Glue スクリプト内から抽出、変換、読み込み (ETL) AWS ワークフローを容易にする Glue ライブラリにもアクセスできます。

このチュートリアルでは、駐車チケットのデータセットを抽出、変換、ロードします。[この作業を行うスクリプトは、Glue Studio AWS のビジュアルエディタを紹介するビッグデータブログの「AWS Glue Studio で ETL を簡単にする」で生成したスクリプトの形式と機能は同じです](#)。AWS このスクリプトをジョブで実行することで、ビジュアルジョブと比較し、AWS Glue ETL スクリプトがどのように機能するかを確認できます。これにより、ビジュアルジョブでまだ使用可能ではない追加機能を使用する準備ができます。

このチュートリアルでは、Python 言語とライブラリを使用しています。Scala でも同様の機能が使用可能です。このチュートリアルを終えると、サンプル Scala スクリプトを生成して調べることができるようになり、Scala AWS Glue ETL スクリプトの記述プロセスを実行する方法を理解できるようになるはずです。

前提条件

このチュートリアルには、次のような前提条件があります。

- テンプレートの実行を指示する AWS Glue Studio のブログ記事と同じ前提条件です。AWS CloudFormation

このテンプレートは、AWS Glue データカタログを使用して、`s3://aws-bigdata-blog/artifacts/gluestudio/`で利用できる駐車券データセットを管理します。これによって参照される以下のリソースが作成されます。

- AWS Glue StudioRole - AWS Glue ジョブを実行する IAM ロール

- AWS Glue StudioAmazon S3Bucket - ブログ関連のファイルを保存する Amazon S3 バケットの名前
- AWS Glue StudioTicketsYYZDB – AWS Glue データカタログデータベース
- AWS Glue StudioTableTickets— ソースとして使用するデータカタログテーブル
- AWS Glue StudioTableTrials— ソースとして使用するデータカタログテーブル
- AWS Glue StudioParkingTicketCount — デステイネーションとして使用するデータカタログテーブル
- AWS Glue Studio のブログ投稿で生成されたスクリプト。ブログ記事が変更された場合、スクリプトは次のテキストからも入手できます。

サンプルスクリプトの生成

AWS Glue Studio のビジュアルエディターを強力なコード生成ツールとして使用して、作成するスクリプトの足場を作成できます。サンプルスクリプトは、このツールを使用して作成します。

これらのステップをスキップしたい場合は、提供されているスクリプトを使用してください。

このチュートリアルでのサンプルスクリプト

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job

args = getResolvedOptions(sys.argv, ["JOB_NAME"])
sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args["JOB_NAME"], args)

# Script generated for node S3 bucket
S3bucket_node1 = glueContext.create_dynamic_frame.from_catalog(
    database="yyz-tickets", table_name="tickets", transformation_ctx="S3bucket_node1"
)

# Script generated for node ApplyMapping
ApplyMapping_node2 = ApplyMapping.apply(
```

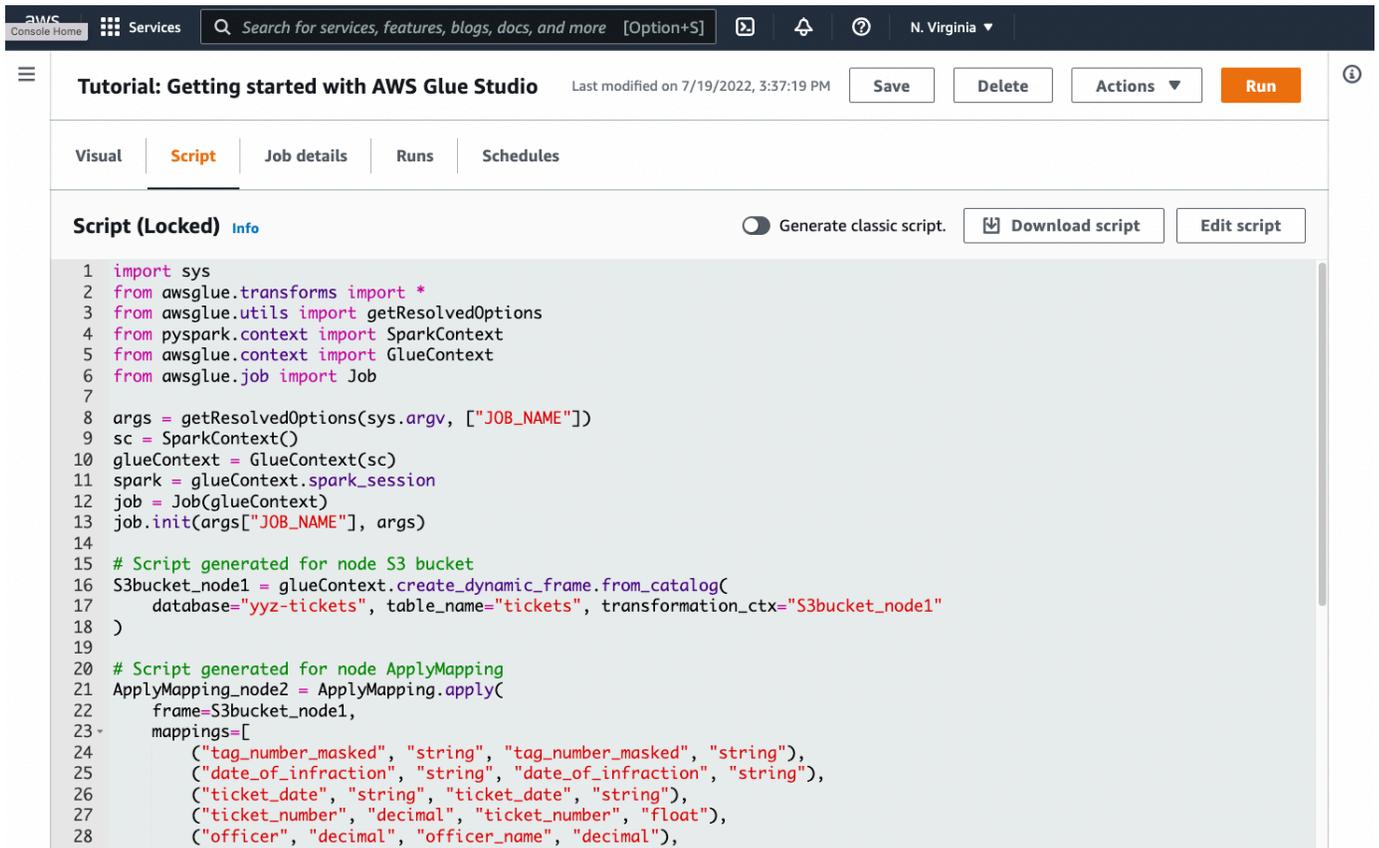
```
frame=S3bucket_node1,
mappings=[
    ("tag_number_masked", "string", "tag_number_masked", "string"),
    ("date_of_infraction", "string", "date_of_infraction", "string"),
    ("ticket_date", "string", "ticket_date", "string"),
    ("ticket_number", "decimal", "ticket_number", "float"),
    ("officer", "decimal", "officer_name", "decimal"),
    ("infraction_code", "decimal", "infraction_code", "decimal"),
    ("infraction_description", "string", "infraction_description", "string"),
    ("set_fine_amount", "decimal", "set_fine_amount", "float"),
    ("time_of_infraction", "decimal", "time_of_infraction", "decimal"),
],
transformation_ctx="ApplyMapping_node2",
)

# Script generated for node S3 bucket
S3bucket_node3 = glueContext.write_dynamic_frame.from_options(
    frame=ApplyMapping_node2,
    connection_type="s3",
    format="glueparquet",
    connection_options={"path": "s3://DOC-EXAMPLE-BUCKET", "partitionKeys": []},
    format_options={"compression": "gzip"},
    transformation_ctx="S3bucket_node3",
)

job.commit()
```

サンプルスクリプトを生成するには

1. AWS Glue スタジオのチュートリアルを完了してください。このチュートリアルを完了するには、[サンプルジョブから「AWS Glue Studio でのジョブの作成」](#)を参照してください。
2. 次のスクリーンショットに示すように、ジョブページの [Script] (スクリプト) タブに移動します。



The screenshot shows the AWS Glue Studio interface. At the top, there's a navigation bar with 'Services', a search bar, and a region dropdown set to 'N. Virginia'. Below that, a breadcrumb trail reads 'Tutorial: Getting started with AWS Glue Studio' with a timestamp 'Last modified on 7/19/2022, 3:37:19 PM'. There are buttons for 'Save', 'Delete', 'Actions', and 'Run'. The main content area has tabs for 'Visual', 'Script', 'Job details', 'Runs', and 'Schedules'. The 'Script' tab is active, showing a Python script. The script is titled 'Script (Locked)' and has a 'Generate classic script.' toggle. There are buttons for 'Download script' and 'Edit script'. The script content is as follows:

```
1 import sys
2 from awsglue.transforms import *
3 from awsglue.utils import getResolvedOptions
4 from pyspark.context import SparkContext
5 from awsglue.context import GlueContext
6 from awsglue.job import Job
7
8 args = getResolvedOptions(sys.argv, ["JOB_NAME"])
9 sc = SparkContext()
10 glueContext = GlueContext(sc)
11 spark = glueContext.spark_session
12 job = Job(glueContext)
13 job.init(args["JOB_NAME"], args)
14
15 # Script generated for node S3 bucket
16 S3bucket_node1 = glueContext.create_dynamic_frame_from_catalog(
17     database="yyz-tickets", table_name="tickets", transformation_ctx="S3bucket_node1"
18 )
19
20 # Script generated for node ApplyMapping
21 ApplyMapping_node2 = ApplyMapping.apply(
22     frame=S3bucket_node1,
23     mappings=[
24         ("tag_number_masked", "string", "tag_number_masked", "string"),
25         ("date_of_infraction", "string", "date_of_infraction", "string"),
26         ("ticket_date", "string", "ticket_date", "string"),
27         ("ticket_number", "decimal", "ticket_number", "float"),
28         ("officer", "decimal", "officer_name", "decimal"),
```

3. [Script] (スクリプト) タブのすべての内容をコピーします。[Job details] (Job の詳細) でスクリプト言語を設定することで、生成するコードを Python と Scala の間で自由に切り替えることができます。

Step 1. ジョブを作成してスクリプトを貼り付ける

このステップでは、AWS で Glue ジョブを作成します。AWS Management Console を使用して、これにより、AWS Glue がスクリプトを実行できるようにする構成が設定されます。また、スクリプトを保存および編集するための場所も作成されます。

ジョブを作成するには

1. AWS Management Console、AWS Glue のランディングページに移動します。
2. サイドナビゲーションペインで、[Jobs] (ジョブ) を選択します。
3. [Create job] (ジョブの作成) で [Spark script editor] (Spark スクリプトエディタ) を選択してから、[Create] (作成) を選択します。
4. オプション — スクリプトの全文を [Script] (スクリプト) ペインに貼り付けます。この代わりに、チュートリアルに従うこともできます。

Step 2. AWS Glue ライブラリのインポート

スクリプトの外部で定義されたコードや設定とやり取りするには、スクリプトの設定が必要です。この作業は AWS Glue Studio の舞台裏で行われています。

このステップでは、次のアクションを実行します。

- GlueContext オブジェクトをインポートして初期化します。スクリプト作成の観点から見ると、このインポートが最も重要なものとなります。これにより、ETL スクリプトの開始点となる、ソースデータセットとターゲットデータセットを定義するための標準的な方法が表示されます。GlueContext クラスの詳細については、「[GlueContext クラス](#)」を参照してください。
- SparkContext および SparkSession を初期化します。これらにより、AWS Glue ジョブ内で使用可能な Spark エンジンを設定できます。入門用の AWS Glue スクリプト内で直接使用する必要はありません。
- getResolvedOptions を呼び出し、スクリプト内でジョブ引数を使用するための準備を行います。ジョブパラメータの解決方法の詳細については、「[the section called “getResolvedOptions”](#)」を参照してください。
- Job を初期化します。Job オブジェクトは設定を設定し、AWS さまざまなオプションの Glue 機能の状態を追跡します。スクリプトは Job オブジェクトがなくても実行できますが、ベストプラクティスとして、後にこの機能が統合された場合に混乱が生じないようにオブジェクトを初期化しておきます。

ジョブのブックマークは、これらの機能の 1 つであり、このチュートリアルではオプションで設定可能です。ジョブのブックマークに関する説明は、「[the section called “オプション — ジョブのブックマークを有効にする”](#)」のセクションで確認できます。

この手順では、次のコードを作成します。このコードは、生成されるサンプルスクリプトの一部です。

```
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job

args = getResolvedOptions(sys.argv, ["JOB_NAME"])
sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
```

```
job = Job(glueContext)
job.init(args["JOB_NAME"], args)
```

AWS Glue ライブラリをインポートするには

- コードのこのセクションをコピーして、[Script] (スクリプト) エディタに貼り付けます。

Note

コードのコピーは技術的に推奨されないと思われるかもしれません。このチュートリアルでは、すべての AWS Glue ETL スクリプトでコア変数に一貫した名前を付けることを奨励するために、これをお勧めします。

ステップ 3. ソースからデータを抽出する

すべての ETL プロセスにおいて、変更するソースデータセットを最初に定義する必要があります。AWS Glue Studio のビジュアルエディターでは、ソースノードを作成してこの情報を入力します。

このステップでは、AWS Glue データカタログで設定されたソースからデータを抽出する `create_dynamic_frame.from_catalog` メソッドに `database` および `table_name` を提供します。

GlueContext オブジェクトは、前のステップで初期化が済んでいます。このオブジェクトを使用して、ソースの設定に使用される `create_dynamic_frame.from_catalog` などのメソッドを検索します。

この手順では、`create_dynamic_frame.from_catalog` を使用して次のコードを記述します。このコードは、生成されるサンプルスクリプトの一部です。

```
S3bucket_node1 = glueContext.create_dynamic_frame.from_catalog(
    database="yyz-tickets", table_name="tickets", transformation_ctx="S3bucket_node1"
)
```

ソースからデータを抽出するには

1. ドキュメントを調べて、AWS Glue GlueContext データカタログで定義されているソースからデータを抽出する方法を見つけてください。これらの方法については、「[the section called “GlueContext”](#)」を参照してください。[create_dynamic_frame.from_catalog](#) メソッドを選択します。glueContext でこのメソッドを呼び出します。

2. `create_dynamic_frame.from_catalog` に関するドキュメントで確認してください。このメソッドには、`database` および `table_name` パラメータが必要です。`create_dynamic_frame.from_catalog` に対し、必要なパラメータを指定します。

AWS Glue データカタログには、ソースデータの場所と形式に関する情報が保存されており、前提条件セクションで設定されています。スクリプトに情報を直接提供する必要はありません。

3. オプション — ジョブのブックマークをサポートするには、このメソッドに `transformation_ctx` パラメータを提供します。ジョブのブックマークに関する説明は、「[the section called “オプション — ジョブのブックマークを有効にする”](#)」のセクションで確認できます。

Note

データ抽出用の一般的なメソッド

[the section called “create_dynamic_frame_from_catalog”](#) AWS Glue データカタログのテーブルへの接続に使用されます。

ソースの構造と場所を記述する設定でジョブを直接提供する必要がある場合は、「[the section called “create_dynamic_frame_from_options”](#)」メソッドを参照してください。`create_dynamic_frame.from_catalog` を使用する場合に比べ、データを説明する詳細なパラメータの指定が必要になります。

必要なパラメータを特定するには、`format_options` および `connection_parameters` に関する補足文書を参照してください。スクリプトに対しソースデータの形式に関する情報を指定する方法については、「[the section called “データ形式に関するオプション”](#)」を参照してください。スクリプトに対しソースデータの場所に関する情報を指定する方法については、「[the section called “接続パラメータ”](#)」を参照してください。

ストリーミングソースから情報を読み取っている場合は、[the section called “create_data_frame_from_catalog”](#) または [the section called “create_data_frame_from_options”](#) メソッドを介して、ソースの情報をジョブに提供します。これらのメソッドは Apache Spark DataFrames を返すことに注意してください。

リファレンスドキュメントでは `create_dynamic_frame_from_catalog` を使用している部分のために、ここで生成されるコードは `create_dynamic_frame.from_catalog` を呼び出しています。これらのメソッドは、最終的には同じコードを呼び出すもので、よりクリーンなコードを作成できるようにインクルードされています。これは、[aws-glue-libs](#) にある Python のラッパーのソースを表示して確認できます。

Step 4. AWS Glue を使用してデータを変換する

ETL プロセスで抽出したソースデータについては、その変更方法を規定する必要があります。この情報を提供するには、AWS Glue Studio のビジュアルエディターでトランスフォームノードを作成します。

このステップでは、希望する現在のフィールドの名前とタイプのマップを `ApplyMapping` メソッドに提供し、`DynamicFrame` を変換します。

以下の変換を実行します。

- 4 つの `location` および `province` キーを削除します。
- `officer` の名前を `officer_name` に変更します。
- `ticket_number` および `set_fine_amount` の型を `float` に変更します。

`create_dynamic_frame.from_catalog` により、`DynamicFrame` オブジェクトが提供されます。A `DynamicFrame` は AWS Glue のデータセットを表します。AWS Glue `DynamicFrames` ートランスフォームは変化する操作です。

Note

`DynamicFrame` について

`DynamicFrame` とは抽象化の 1 つであり、データ内のエントリの名前とタイプに関する記述を、データセットに接続できるようにします。Apache Spark には、と呼ばれる同様の抽象化があります。DataFrameの説明については、「[Spark SQL DataFrames ガイド](#)」を参照してください。

`DynamicFrames` を使用すると、データセットのスキーマを動的に記述できます。価格列のあるデータセットを考えてみましょう。価格を文字列として格納するエントリもあれば、価格を `double` として格納するエントリもあります。AWS Glue はスキーマを計算し on-the-fly、行ごとに自己記述型のレコードを作成します。

一貫性のないフィールド (価格など) は、そのフレーム用のスキーマ内で、型 (`ChoiceType`) を使用しながら明示的に表されます。一貫性のないフィールドに対応するには、`DropFields` を使用して削除するか、`ResolveChoice` を使用して解決します。これらは、`DynamicFrame` で使用可能な変換です。その後、`writeDynamicFrame` を使用して、データをデータレイクに再度書き込むことができます。

DynamicFrame クラスにあるメソッドからは、多くの同様な変換を呼び出すことができます。これにより、スクリプトが読みやすくなります。DynamicFrame の詳細については、「[the section called “DynamicFrame”](#)」を参照してください。

この手順では、ApplyMapping を使用して次のコードを記述します。このコードは、生成されるサンプルスクリプトの一部です。

```
ApplyMapping_node2 = ApplyMapping.apply(  
    frame=S3bucket_node1,  
    mappings=[  
        ("tag_number_masked", "string", "tag_number_masked", "string"),  
        ("date_of_infraction", "string", "date_of_infraction", "string"),  
        ("ticket_date", "string", "ticket_date", "string"),  
        ("ticket_number", "decimal", "ticket_number", "float"),  
        ("officer", "decimal", "officer_name", "decimal"),  
        ("infraction_code", "decimal", "infraction_code", "decimal"),  
        ("infraction_description", "string", "infraction_description", "string"),  
        ("set_fine_amount", "decimal", "set_fine_amount", "float"),  
        ("time_of_infraction", "decimal", "time_of_infraction", "decimal"),  
    ],  
    transformation_ctx="ApplyMapping_node2",  
)
```

AWS Glue でデータを変換するには

1. フィールドを変更や削除するための変換を特定するには、ドキュメントを参照してください。詳細については、「[the section called “GlueTransform”](#)」を参照してください。ApplyMapping 変換を選択します。ApplyMapping の詳細については、「[the section called “ApplyMapping”](#)」を参照してください。ApplyMapping 変換オブジェクトから `apply` を呼び出します。

Note

ApplyMapping とは

ApplyMapping は DynamicFrame にアクセスし、それを変換します。これにより、フィールドの変換を表すタプルのリスト、つまり「マッピング」を取得します。最初の 2 つのタプル要素 (フィールドの名前と型) は、フレーム内のフィールドを識別するために使用されます。その次にある 2 つのパラメータも、フィールドの名前と型です。ApplyMapping ソースフィールドをターゲット名に変換し DynamicFrame、新しい名前を入力すると返されます。提供されていないフィールドは戻り値から削除されます。

`apply` を呼び出す代わりに、`DynamicFrame` にある `apply_mapping` メソッドを使用して同様の変換を呼び出すと、さらにスムーズで読みやすいコードを作成できます。詳細については、「[the section called “apply_mapping”](#)」を参照してください。

- 必要なパラメータを特定するには、`ApplyMapping` のドキュメントを参照してください。[the section called “ApplyMapping”](#) を参照してください。このメソッドでは、`frame` および `mappings` パラメータが必要です。`ApplyMapping` に対し、必要なパラメータを指定します。
- オプション — ジョブのブックマークをサポートするには、メソッドに `transformation_ctx` を提供します。ジョブのブックマークに関する説明は、「[the section called “オプション — ジョブのブックマークを有効にする”](#)」のセクションで確認できます。

Note

Apache Spark の機能

ジョブ内の ETL ワークフローを効率化するための変換機能が提供されています。また、より一般的な目的のために構築された、ジョブ内の Spark プログラムで使用できるライブラリにもアクセスできます。これらを使用するには、`DynamicFrame` と `DataFrame` の間で変換します。

[the section called “toDF”](#) を使用して `DataFrame` を作成できます。その後、`DataFrame` にあるメソッドを使用してデータセットを変換できます。これらのメソッドの詳細については、[を参照してください DataFrame](#)。その後、[を使用して逆方向に変換し](#)、AWS Glue 操作を使用してフレームをターゲットに読み込むことができます。[the section called “fromDF”](#)

Step 5. ターゲットにデータをロードする

通常、変換後のデータは、そのソースとは別の場所に保存します。この操作を実行するには、AWS Glue Studio のビジュアルエディターでターゲットノードを作成します。

このステップでは、データを Amazon S3 のターゲットバケットにロードするための `write_dynamic_frame.from_options` メソッド `connection_type`、`connection_options`、`format`、`format_options` を提供します。

ステップ 1 において、`GlueContext` オブジェクトは初期化済みです。AWS Glue では、ソースと同様にターゲットの設定に使用されるメソッドがここにあります。

この手順では、`write_dynamic_frame.from_options` を使用して次のコードを記述します。このコードは、生成されるサンプルスクリプトの一部です。

```
S3bucket_node3 = glueContext.write_dynamic_frame.from_options(  
    frame=ApplyMapping_node2,  
    connection_type="s3",  
    format="glueparquet",  
    connection_options={"path": "s3://DOC-EXAMPLE-BUCKET", "partitionKeys": []},  
    format_options={"compression": "gzip"},  
    transformation_ctx="S3bucket_node3",  
)
```

データをターゲットにロードするには

1. ターゲットの Amazon S3 バケットに対しデータをロードするメソッドは、ドキュメントで確認してください。これらの方法については、「[the section called “GlueContext”](#)」を参照してください。[the section called “write_dynamic_frame_from_options”](#) メソッドを選択します。glueContext でこのメソッドを呼び出します。

Note

データロード用の一般的なメソッド

write_dynamic_frame.from_options は、データのロードに使用される最も一般的なメソッドです。AWS Glue で利用できるすべてのターゲットをサポートします。AWS Glue 接続で定義された JDBC ターゲットに書き込む場合は、メソッドを使用してください。[the section called “write_dynamic_frame_from_jdbc_conf”](#) AWS Glue 接続には、データソースへの接続方法に関する情報が格納されます。これにより、connection_options でそれらの情報を入力する必要がなくなります。ただし、dbtable を提供するためには、引き続き connection_options を使用する必要があります。

write_dynamic_frame.from_catalog は、データロード用として、一般的なメソッドではありません。このメソッドは、基になるデータセットを更新せずに AWS Glue Data Catalog を更新し、基になるデータセットを変更する他のプロセスと組み合わせて使用されます。詳細については、「[the section called “スキーマを更新し、新規パーティションを追加する”](#)」を参照してください。

2. [the section called “write_dynamic_frame_from_options”](#) に関するドキュメントで確認してください。このメソッドでは、frame、connection_type、format、connection_options、format_options が必要です。glueContext でこのメソッドを呼び出します。

- a. 必要なパラメータを特定するには、`format_options` および `format` に関する補足文書を参照してください。データ形式の説明については、「[the section called “データ形式に関するオプション”](#)」を参照してください。
 - b. 必要なパラメータを特定するには、`connection_type` および `connection_options` に関する補足文書を参照してください。接続の説明については、「[the section called “接続パラメータ”](#)」を参照してください。
 - c. `write_dynamic_frame.from_options` に対し、必要なパラメータを指定します。このメソッドの設定は、`create_dynamic_frame.from_options` と似ています。
3. オプション — ジョブのブックマークをサポートするには、`transformation_ctx` を `write_dynamic_frame.from_options` に提供します。ジョブのブックマークに関する説明は、「[the section called “オプション — ジョブのブックマークを有効にする”](#)」のセクションで確認できます。

ステップ 6。Job オブジェクトのコミット

Job オブジェクトは、ステップ 1 で初期化してあります。スクリプトの終了時に、ライフサイクルを手動で閉じる必要があります。一部のオプション機能が正しく機能するために、これが必要となります。この作業は AWS Glue Studio の舞台裏で行われています。

このステップでは、Job オブジェクトの `commit` メソッドを呼び出します。

この手順では、次のコードを作成します。このコードは、生成されるサンプルスクリプトの一部です。

```
job.commit()
```

Job オブジェクトをコミットするには

1. まだ行っていない場合は、前のセクションで概説したオプションの手順を実行して `transformation_ctx` を含めます。
2. `commit` を呼び出します。

オプション — ジョブのブックマークを有効にする

これまでのすべてのステップでは、`transformation_ctx` パラメータを設定することを促しています。これは、ジョブのブックマークと呼ばれる機能に関連した設定です。

ジョブのブックマークを使用すると、以前の作業を簡単に追跡できるデータセットに対して定期的に実行されるジョブによって、時間とコストを節約できます。Job ブックマークは、以前に実行したデータセット全体の AWS Glue 変換の進行状況を追跡します。前回の実行が終了した場所を追跡することで、AWS Glue は処理をまだ処理していない行に限定できます。ブックマークの詳細については、「[the section called “ジョブのブックマークを使用した処理済みデータの追跡”](#)」を参照してください。

前の例で説明したように、ジョブのブックマークを有効にするには、始めに、提供されている関数に `transformation_ctx` ステートメントを追加します。ジョブのブックマークの状態は実行後も保持されます。`transformation_ctx` パラメータは、その状態にアクセスするために使用されるキーです。これらのステートメント単体では機能を持ちません。また、ジョブのための設定内で、その機能を有効にする必要があります。

この手順では、AWS Management Consoleを使用してジョブのブックマークを有効にします。

ジョブのブックマークを有効にするには

1. 対象となるジョブの、[Job details] (ジョブの詳細) セクションに移動します。
2. [Job bookmark] (ジョブのブックマーク) で [Enable] (有効化) 設定します。

ステップ 7。コードをジョブとして実行する

このステップでは、チュートリアルを正常に完了したことを確認するために、ジョブを実行します。これは、AWS Glue Studio のビジュアルエディターのようにボタンをクリックするだけで実行できます。

コードをジョブとして実行するには

1. タイトルバーで [Untitled job] (無題のジョブ) 選択し、ジョブ名を編集あるいは設定します。
2. [Job details] (ジョブの詳細) タブを開きます。ジョブに [IAM Role] (IAM ロール) を割り当てます。AWS CloudFormation テンプレートによって作成されたものは、AWS Glue Studio チュートリアルの前提条件で使用できます。このチュートリアルを完了していれば、通常は、AWS Glue StudioRole のように表示されています。
3. [Save] (保存) を選択して、スクリプトを保存します。
4. [Run] (実行) を選択して、ジョブを実行します。
5. [Runs] (実行) タブに移動して、ジョブが完了したことを確認します。

6. `write_dynamic_frame.from_options` のターゲットである、***DOC-EXAMPLE-BUCKET*** に移動します。想定どおりに出力されていることを確認します。

ジョブの設定と管理の詳細については、「[the section called “独自のカスタムスクリプトの提供”](#)」を参照してください。

詳細情報

Apache Spark のライブラリとメソッドは AWS Glue スクリプトで使用できます。Spark のドキュメントを参照することで、これらの付属ライブラリで何が行えるのか把握できます。詳細については、「[examples section of the Spark source repository](#)」(Spark ソースリポジトリのサンプルセクション)を参照してください。

AWS Glue 2.0+には、デフォルトでいくつかの一般的なPythonライブラリが含まれています。Scala または Python 環境の AWS Glue ジョブに独自の依存関係を読み込むメカニズムもあります。Python の依存関係については、「[the section called “Python ライブラリ”](#)」を参照してください。

Python で AWS Glue 機能を使用する方法のその他の例については、[を参照してくださいthe section called “Python サンプル”](#)。Scala と Python ジョブには同等の機能があるため、Python の例を基にして、同様の処理を Scala で行う方法についても考察することが可能です。

で AWS Glue ETL スクリプトをプログラムする PySpark

の Python コード例とユーティリティはAWS Glue、GitHub ウェブサイト[AWS Glueのサンプルリポジトリ](#)にあります。

AWS Glue での Python の使用

AWS Glue は、抽出、変換、ロード (ETL) ジョブをスクリプト化するための PySpark Python ダイアレクトの拡張をサポートしています。このセクションでは、ETL スクリプトと AWS Glue API で Python を使用する方法について説明します。

- [AWS Glue で Python を使用するためのセットアップ](#)
- [Python での AWS Glue API の呼び出し](#)
- [AWS Glue での Python ライブラリの使用](#)
- [AWS Glue Python コードサンプル](#)

AWS Glue PySpark 拡張機能

AWS Glue は、PySpark Python ダイアレクトに次の拡張機能を作成しました。

- [getResolvedOptions](#) を使用して、パラメータにアクセスする
- [PySpark 拡張子型](#)
- [DynamicFrame クラス](#)
- [DynamicFrameCollection クラス](#)
- [DynamicFrameWriter クラス](#)
- [DynamicFrameReader クラス](#)
- [GlueContext クラス](#)

AWS Glue PySpark 変換

AWS Glue は、PySpark ETL オペレーションで使用する次の変換クラスを作成しました。

- [GlueTransform 基本クラス](#)
- [ApplyMapping クラス](#)
- [DropFields クラス](#)
- [DropNullFields クラス](#)
- [ErrorsAsDynamicFrame クラス](#)
- [FillMissingValues クラス](#)
- [フィルタクラス](#)
- [FindIncrementalMatches クラス](#)
- [FindMatches クラス](#)
- [FlatMap クラス](#)
- [Join クラス](#)
- [マップクラス](#)
- [MapToCollection クラス](#)
- [mergeDynamicFrame](#)
- [クラスの関連付け](#)
- [RenameField クラス](#)

- [ResolveChoice クラス](#)
- [SelectFields クラス](#)
- [SelectFromCollection クラス](#)
- [スピゴットクラス](#)
- [SplitFields クラス](#)
- [SplitRows クラス](#)
- [Unbox クラス](#)
- [UnnestFrame クラス](#)

AWS Glue で Python を使用するためのセットアップ

Python を使用して Spark ジョブ用の ETL スクリプトを開発します。ETL ジョブでサポートされている Python のバージョンは、ジョブが使用する AWS Glue バージョンによって異なります。AWS Glue バージョンの詳細については、「[Glue version job property](#)」を参照してください。

AWS Glue で Python を使用するためのシステムを設定するには

Python をインストールするには、以下の手順を実行して AWS Glue API を呼び出せるようにします。

1. Python がインストールされていない場合は、[Python.org のダウンロードページ](#) からダウンロードしてインストールします。
2. [AWS CLI のドキュメント](#) で説明されているように AWS Command Line Interface (AWS CLI) をインストールします。

AWS CLI は、Python を使用するのに直接必要なわけではありません。ただし、インストールおよび設定を行うと、アカウントの認証情報を使用して AWS を設定し、その動作を確認するのに便利です。

3. [Boto 3 クイックスタート](#) で説明されているように、AWS SDK for Python (Boto 3) をインストールします。

Boto 3 リソース API は AWS Glue にはまだ使用できません。現時点では、Boto 3 クライアント API のみ使用することができます。

Boto 3 の詳細については、「[AWS SDK for Python \(Boto3\) Getting Started](#)」を参照してください。

Python のコード例や AWS Glue のユーティリティは、GitHub サイトの [AWS Glue サンプルリポジトリ](#) で公開されています。

Python での AWS Glue API の呼び出し

Boto 3 リソース API は AWS Glue にはまだ使用できないことに注意してください。現時点では、Boto 3 クライアント API のみ使用することができます。

Python の AWS Glue API 名

Java や他のプログラミング言語での AWS Glue API 名は、通常 CamelCased です。ただし、Python から呼び出されるとき、これらの一般名は、より「Python 的」にするために小文字に変更され、名前の一部がアンダースコア文字で区切られます。[AWS Glue API](#) リファレンスドキュメントでは、これらの Python 用の名前を一般的な CamelCased 形式の名前の後に括弧で囲んで一覧表示しています。

ただし、AWS Glue API 名自体は小文字に変換されますが、パラメータ名は大文字のままです。次のセクションで説明されますが、AWS Glue API の呼び出し時にパラメータを名前で渡す必要があるため、この点を覚えておくことが重要です。

AWS Glue の Python パラメータの受け渡しとアクセス

AWS Glue API の Python 呼び出しでは、明示的に名前でパラメータを渡すことが最善です。例:

```
job = glue.create_job(Name='sample', Role='Glue_DefaultRole',
                      Command={'Name': 'glueetl',
                               'ScriptLocation': 's3://my_script_bucket/scripts/
my_etl_script.py'})
```

[Job 構造](#) または [JobRun 構造](#) の ETL スクリプトに渡す引数として指定する名前と値のタプルのディクショナリを Python が作成することを理解しておくことは役立ちます。その後 Boto 3 が REST API 呼び出しを経由して JSON 形式でそれらを AWS Glue に渡します。つまり、スクリプトでそれらにアクセスするときは、引数の順序に依存することはできません。

たとえば、Python Lambda ハンドラ関数で JobRun を開始しようとしており、複数のパラメータを指定するとします。コードは以下のようになります。

```
from datetime import datetime, timedelta
```

```

client = boto3.client('glue')

def lambda_handler(event, context):
    last_hour_date_time = datetime.now() - timedelta(hours = 1)
    day_partition_value = last_hour_date_time.strftime("%Y-%m-%d")
    hour_partition_value = last_hour_date_time.strftime("%-H")

    response = client.start_job_run(
        JobName = 'my_test_job',
        Arguments = {
            '--day_partition_key': 'partition_0',
            '--hour_partition_key': 'partition_1',
            '--day_partition_value': day_partition_value,
            '--hour_partition_value': hour_partition_value } )

```

これらのパラメータに ETL スクリプトで確実にアクセスするには、AWS Glue の `getResolvedOptions` 関数を使用して名前指定し、その後作成されたディクショナリからアクセスします。

```

import sys
from awsglue.utils import getResolvedOptions

args = getResolvedOptions(sys.argv,
                          ['JOB_NAME',
                           'day_partition_key',
                           'hour_partition_key',
                           'day_partition_value',
                           'hour_partition_value'])

print "The day partition key is: ", args['day_partition_key']
print "and the day partition value is: ", args['day_partition_value']

```

ネストされた JSON 文字列により AWS Glue ETL ジョブに引数を渡す際は、そのパラメータ値を保持するために、ジョブの実行が開始される前にパラメータ文字列をエンコードし、ジョブスクリプトが参照される前にパラメータ文字列をデコードする必要があります。例えば、次の引数文字列を考えてみます。

```

glue_client.start_job_run(JobName = "gluejobname", Arguments={
    "--my_curly_braces_string": '{"a": {"b": {"c": [{"d": {"e": 42}}]}}'})
})

```

このパラメータを正しく渡すには、引数を Base64 の文字列としてエンコードする必要があります。

```
import base64
...
sample_string='{"a": {"b": {"c": [{"d": {"e": 42}}]}}}'
sample_string_bytes = sample_string.encode("ascii")

base64_bytes = base64.b64encode(sample_string_bytes)
base64_string = base64_bytes.decode("ascii")
...
glue_client.start_job_run(JobName = "gluejobname", Arguments={
"--my_curly_braces_string": base64_bytes})
...
sample_string_bytes = base64.b64decode(base64_bytes)
sample_string = sample_string_bytes.decode("ascii")
print(f"Decoded string: {sample_string}")
...
```

例: ジョブを作成し実行する

以下の例では、Python を使用して AWS Glue API を呼び出し、ETL ジョブを作成して実行する方法を示します。

ジョブを作成し実行するには

1. AWS Glue クライアントのインスタンスを作成します。

```
import boto3
glue = boto3.client(service_name='glue', region_name='us-east-1',
                    endpoint_url='https://glue.us-east-1.amazonaws.com')
```

2. ジョブを作成します。次のコードに示すように、ETL コマンドの名前として `glueetl` を使用する必要があります。

```
myJob = glue.create_job(Name='sample', Role='Glue_DefaultRole',
                        Command={'Name': 'glueetl',
                                'ScriptLocation': 's3://my_script_bucket/
scripts/my_etl_script.py'})
```

3. 前のステップで作成したジョブの新しい実行を開始します。

```
myNewJobRun = glue.start_job_run(JobName=myJob['Name'])
```

4. ジョブのステータスを取得します。

```
status = glue.get_job_run(JobName=myJob['Name'], RunId=myNewJobRun['JobRunId'])
```

5. ジョブ実行の現在の状態を出力します。

```
print(status['JobRun']['JobRunState'])
```

AWS Glue での Python ライブラリの使用

AWS Glue AWS Glue ETL で使用するための Python モジュールやライブラリを追加でインストールすることができます。

トピック

- [AWS Glue 2.0 の pip を使用した追加 Python モジュールのインストール](#)
- [PySpark ネイティブ機能を備えた Python ファイルを含める](#)
- [ビジュアル変換を使用するプログラミングスクリプト](#)
- [AWS Glue で提供済みの Python モジュール](#)
- [取り込みのためのライブラリの圧縮](#)
- [AWS Glue Studio ノートブックに Python ライブラリをロードする](#)
- [開発エンドポイントへの Python ライブラリのロード](#)
- [ジョブまたはでの Python ライブラリの使用 JobRun](#)

AWS Glue 2.0 の pip を使用した追加 Python モジュールのインストール

AWS Glue では、Python Package Installer (pip3) を使用して、AWS Glue ETL で使用するモジュールを追加でインストールします。--additional-python-modules パラメータでコマンド区切りの Python モジュールのリストを指定することで、新しいモジュールを追加したり、既存のモジュールのバージョンを変更したりできます。Amazon S3 にディストリビューションをアップロードすることでライブラリのカスタムディストリビューションをインストールできます。その後で、モジュールのリストに Amazon S3 オブジェクトへのパスを含めます。

--python-modules-installer-option パラメータを使用すると、pip3 に追加オプションを渡すことができます。例えば "--upgrade" を渡すことで、"--additional-python-modules" で指定されたパッケージをアップグレードできます。その他の例については、「[AWS Glue 2.0 を使用して Spark ETL ワークロードのホイールから Python モジュールを構築する](#)」を参照してください。

Python の依存関係がネイティブのコンパイル済みコードに推移的に依存している場合は、次の制限に反して実行できません。AWS Glue はジョブ環境でのネイティブコードのコンパイルをサポートしていません。ただし、AWS Glue ジョブは Amazon Linux 2 環境内で実行されます。ディストリビューション版 Wheel を通じて、ネイティブの依存関係をコンパイルされた形式で提供できる場合があります。

例えば、更新したり新しい scikit-learn モジュールを追加したりするには、"--additional-python-modules", "scikit-learn==0.21.3" のキー/値を使用します。

また、--additional-python-modules オプションの中で、Python ホイールモジュールへの Amazon S3 パスを指定できます。例:

```
--additional-python-modules s3://aws-glue-native-spark/tests/j4.2/ephem-3.7.7.1-cp37-cp37m-linux_x86_64.whl,s3://aws-glue-native-spark/tests/j4.2/fbprophet-0.6-py3-none-any.whl,scikit-learn==0.21.3
```

AWS Glue コンソールの ジョブパラメータ フィールド--additional-python-modulesで を指定するか、AWS SDK でジョブ引数を変更します。ジョブパラメータの設定については、「[the section called “ジョブのパラメータ”](#)」を参照してください。

PySpark ネイティブ機能を備えた Python ファイルを含める

AWS Glue は PySpark を使用して AWS Glue ETL ジョブに Python ファイルを含めます。可能な場合には、依存関係を管理するために、--additional-python-modules を使用することが必要になります。Python ファイルをインクルードするには、--extra-py-files ジョブパラメータを使用します。依存関係は Amazon S3 でホストされている必要があります。引数の値は、スペースを含まない Amazon S3 パスのカンマ区切りリストである必要があります。この機能は、Spark で使用する Python の依存関係管理のように動作します。Spark での Python 依存関係管理の詳細については、Apache Spark [ドキュメントの PySpark 「ネイティブ機能の使用」](#) ページを参照してください。--extra-py-files は、追加のコードがパッケージ化されていない場合や、依存関係を管理するための既存のツールチェーンを使用して Spark プログラムを移行する場合に役立ちます。依存関係ツールをメンテナンス可能にするには、送信する前に依存関係をバンドルする必要があります。

ビジュアル変換を使用するプログラミングスクリプト

AWS Glue Studio ビジュアルインターフェイスを使用して AWS Glue ジョブを作成すると、マネージドデータ変換ノードとカスタムビジュアル変換を使用してデータを変換できます。マネージドデータ変換ノードの詳細については、「[the section called “AWS Glue マネージドデータ変換ノードの編集”](#)」を参照してください。カスタムビジュアル変換の詳細については、「[the section called “カスタ](#)

[「ムビジュアル変換」](#)を参照してください。ビジュアル変換を使用するスクリプトは、ジョブの [言語] が Python を使用するように設定されている場合にのみ生成できます。

ビジュアル変換を使用して AWS Glue ジョブを生成する場合、AWS Glue Studio はジョブ設定の `--extra-py-files` パラメータを使用して、これらの変換をランタイム環境に含めます。ジョブパラメータについては、「[the section called “ジョブのパラメータ”](#)」を参照してください。生成されたスクリプトまたはランタイム環境を変更するとき、スクリプトを正常に実行するためにこのジョブ設定を保存する必要があります。

AWS Glue で提供済みの Python モジュール

これらの提供済みモジュールのバージョンを変更するには、`--additional-python-modules` ジョブパラメータにより新しいバージョンを指定します。

AWS Glue version 2.0

AWS Glue バージョン 2.0 には、すぐに使用できる次の Python モジュールが含まれています。

- `avro-python3==1.10.0`
- `awscli==1.27.60`
- `boto3==1.12.4`
- `botocore==1.15.4`
- `certifi==2019.11.28`
- `chardet==3.0.4`
- `click==8.1.3`
- `colorama==0.4.4`
- `cycler==0.10.0`
- `Cython==0.29.15`
- `docutils==0.15.2`
- `enum34==1.1.9`
- `fsspec==0.6.2`
- `idna==2.9`
- `importlib-metadata==6.0.0`
- `jmespath==0.9.4`

- joblib==0.14.1
- kiwisolver==1.1.0
- matplotlib==3.1.3
- mpmath==1.1.0
- nltk==3.5
- numpy==1.18.1
- pandas==1.0.1
- patsy==0.5.1
- pmdarima==1.5.3
- ptvsd==4.3.2
- pyarrow==0.16.0
- pyasn1==0.4.8
- pydevd==1.9.0
- pyhocon==0.3.54
- PyMySQL==0.9.3
- pyparsing==2.4.6
- python-dateutil==2.8.1
- pytz==2019.3
- PyYAML==5.3.1
- regex==2022.10.31
- requests==2.23.0
- rsa==4.7.2
- s3fs==0.4.0
- s3transfer==0.3.3
- scikit-learn==0.22.1
- scipy==1.4.1
- setuptools==45.2.0
- six==1.14.0
- Spark==1.0

- statsmodels==0.11.1
- subprocess32==3.5.4
- sympy==1.5.1
- tbats==1.0.9
- tqdm==4.64.1
- typing-extensions==4.4.0
- urllib3==1.25.8
- wheel==0.35.1
- zipp==3.12.0

AWS Glue バージョン 3.0

AWS Glue バージョン 3.0 には、すぐに使用できる次の Python モジュールが含まれています。

- aiobotocore==1.4.2
- aiohttp==3.8.3
- aioitertools==0.11.0
- aiosignal==1.3.1
- async-timeout==4.0.2
- asyncctest==0.13.0
- attrs==22.2.0
- avro-python3==1.10.2
- boto3==1.18.50
- botocore==1.21.50
- certifi==2021.5.30
- chardet==3.0.4
- charset-normalizer==2.1.1
- click==8.1.3
- cycler==0.10.0
- Cython==0.29.4

- docutils==0.17.1
- enum34==1.1.10
- frozenlist==1.3.3
- fsspec==2021.8.1
- idna==2.10
- importlib-metadata==6.0.0
- jmespath==0.10.0
- joblib==1.0.1
- kiwisolver==1.3.2
- matplotlib==3.4.3
- mpmath==1.2.1
- multidict==6.0.4
- nltk==3.6.3
- numpy==1.19.5
- packaging==23.0
- pandas==1.3.2
- patsy==0.5.1
- Pillow==9.4.0
- pip==23.0
- pmdarima==1.8.2
- ptvsd==4.3.2
- pyarrow==5.0.0
- pydevd==2.5.0
- pyhocon==0.3.58
- PyMySQL==1.0.2
- pyparsing==2.4.7
- python-dateutil==2.8.2
- pytz==2021.1
- PyYAML==5.4.1

- regex==2022.10.31
- requests==2.23.0
- s3fs==2021.8.1
- s3transfer==0.5.0
- scikit-learn==0.24.2
- scipy==1.7.1
- six==1.16.0
- Spark==1.0
- statsmodels==0.12.2
- subprocess32==3.5.4
- sympy==1.8
- tbats==1.1.0
- threadpoolctl==3.1.0
- tqdm==4.64.1
- typing_extensions==4.4.0
- urllib3==1.25.11
- wheel==0.37.0
- wrapt==1.14.1
- yarl==1.8.2
- zipp==3.12.0

AWS Glue バージョン 4.0

AWS Glue バージョン 4.0 には、すぐに使用できる次の Python モジュールが含まれています。

- aiobotocore==2.4.1
- aiohttp==3.8.3
- aioitertools==0.11.0
- aiosignal==1.3.1
- async-timeout==4.0.2
- asynctest==0.13.0

- attrs==22.2.0
- avro-python3==1.10.2
- boto3==1.24.70
- botocore==1.27.59
- certifi==2021.5.30
- chardet==3.0.4
- charset-normalizer==2.1.1
- click==8.1.3
- cycler==0.10.0
- Cython==0.29.32
- docutils==0.17.1
- enum34==1.1.10
- frozenlist==1.3.3
- fsspec==2021.8.1
- idna==2.10
- importlib-metadata==5.0.0
- jmespath==0.10.0
- joblib==1.0.1
- kaleido==0.2.1
- kiwisolver==1.4.4
- matplotlib==3.4.3
- mpmath==1.2.1
- multidict==6.0.4
- nltk==3.7
- numpy==1.23.5
- packaging==23.0
- pandas==1.5.1
- patsy==0.5.1
- Pillow==9.4.0

- pip==23.0.1
- plotly==5.16.0
- pmdarima==2.0.1
- ptvsd==4.3.2
- pyarrow==10.0.0
- pydevd==2.5.0
- pyhocon==0.3.58
- PyMySQL==1.0.2
- pyparsing==2.4.7
- python-dateutil==2.8.2
- pytz==2021.1
- PyYAML==6.0.1
- regex==2022.10.31
- requests==2.23.0
- s3fs==2022.11.0
- s3transfer==0.6.0
- scikit-learn==1.1.3
- scipy==1.9.3
- setuptools==49.1.3
- six==1.16.0
- statsmodels==0.13.5
- subprocess32==3.5.4
- sympy==1.8
- tbats==1.1.0
- threadpoolctl==3.1.0
- tqdm==4.64.1
- typing_extensions==4.4.0
- urllib3==1.25.11
- wheel==0.37.0

- wrapt==1.14.1
- yarl==1.8.2
- zipp==3.10.0

取り込みのためのライブラリの圧縮

ライブラリは、単一の .py ファイルに含まれていない限り、.zip アーカイブにパッケージ化される必要があります。パッケージディレクトリは、アーカイブのルートにあって、パッケージの `__init__.py` ファイルを含んでいる必要があります。そうすると、Python は通常の方法でパッケージをインポートできるようになります。

ライブラリが 1 つの .py ファイルにある単一の Python モジュールでのみ構成されている場合、.zip ファイルに入れる必要はありません。

AWS Glue Studio ノートブックに Python ライブラリをロードする

AWS Glue Studio ノートブックで Python ライブラリを指定するには、[「追加の Python モジュールのインストール」](#)を参照してください。

開発エンドポイントへの Python ライブラリのロード

異なる ETL スクリプトに異なるライブラリセットを使用している場合、各セットに別々の開発エンドポイントをセットアップするか、スクリプトを切り替えるたびに開発エンドポイントがロードするライブラリ .zip ファイルを上書きすることができます。

コンソールを使用して、作成時に開発エンドポイントに 1 つまたは複数のライブラリ .zip ファイルを指定できます。名前と IAM ロールを割り当てた後、[Script Libraries and job parameters (optional)] (スクリプトライブラリおよびジョブパラメータ (任意)) をクリックし、[Python library path] (Python ライブラリパス) ボックスに、ライブラリ .zip ファイルへの完全な Amazon S3 パスを入力します。例:

```
s3://bucket/prefix/site-packages.zip
```

必要に応じてファイルへの複数のフルパスを指定できますが、以下のように、スペースなしでカンマで区切ります。

```
s3://bucket/prefix/lib_A.zip,s3://bucket_B/prefix/lib_X.zip
```

これらの .zip ファイルを後で更新する場合は、コンソールを使用して開発エンドポイントにそのファイルを再インポートできます。該当する開発エンドポイントに移動し、横にあるチェックボックスをオンにして、[Action] (アクション) メニューから [Update ETL libraries] (ETL ライブラリの更新) を選択します。

同様の方法で、AWS Glue API を使用してライブラリファイルを指定できます。[CreateDevEndpoint アクション \(Python: create_dev_endpoint\)](#) を呼び出して開発エンドポイントを作成する場合、ExtraPythonLibsS3Path パラメータでライブラリへの 1 つ以上のフルパスを指定できます。以下のような呼び出しになります。

```
dep = glue.create_dev_endpoint(  
    EndpointName="testDevEndpoint",  
    RoleArn="arn:aws:iam::123456789012",  
    SecurityGroupIds="sg-7f5ad1ff",  
    SubnetId="subnet-c12fdb4",  
    PublicKey="ssh-rsa AAAAB3NzaC1yc2EAAAADAQABAAQCTp04H/y...",  
    NumberOfNodes=3,  
    ExtraPythonLibsS3Path="s3://bucket/prefix/lib_A.zip,s3://bucket_B/prefix/  
lib_X.zip")
```

開発エンドポイントを更新するときに、[DevEndpointCustomLibraries](#) オブジェクトを使用し [UpdateDevEndpoint \(update_dev_endpoint\)](#) の呼び出し時に UpdateEtlLibraries パラメータを True に設定して、ロードするライブラリも更新できます。

ジョブまたは の Python ライブラリの使用 JobRun

コンソールで新しいジョブを作成する際、[Script Libraries and job parameters (optional)] (スクリプトライブラリおよびジョブパラメータ (任意)) をクリックし、開発エンドポイント作成時と同じ方法で Amazon S3 ライブラリの完全なパスを入力することで、1 つ以上のライブラリ .zip ファイルを指定できます。

```
s3://bucket/prefix/lib_A.zip,s3://bucket_B/prefix/lib_X.zip
```

[CreateJob \(create_job\)](#) を呼び出している場合は、以下のようにデフォルトの --extra-py-files パラメータを使用してデフォルトのライブラリへの 1 つ以上のフルパスを指定できます。

```
job = glue.create_job(Name='sampleJob',  
    Role='Glue_DefaultRole',  
    Command={'Name': 'glueetl',
```

```
        'ScriptLocation': 's3://my_script_bucket/scripts/  
my_etl_script.py'},  
        DefaultArguments={'--extra-py-files': 's3://bucket/prefix/  
lib_A.zip,s3://bucket_B/prefix/lib_X.zip'})
```

その後、を開始するときに JobRun、デフォルトのライブラリ設定を別のライブラリ設定で上書きできます。

```
runId = glue.start_job_run(JobName='sampleJob',  
                           Arguments={'--extra-py-files': 's3://bucket/prefix/  
lib_B.zip'})
```

AWS Glue Python コードサンプル

- [コード例: データの結合と関係付け](#)
- [コード例: ResolveChoice、Lambda、および ApplyMapping を使用したデータ準備](#)

コード例: データの結合と関係付け

この例では、<http://everypolitician.org/> から、Amazon Simple Storage Service (Amazon S3) の sample-dataset バケットにデータセットをダウンロードして使用します。s3://awsglue-datasets/examples/us-legislators/all このデータセットには、米国議会議員や米国下院および上院議員の議席に関する JSON 形式のデータが含まれており、このチュートリアルの目的のため少し変更され、パブリック Amazon S3 バケットで利用可能になりました。

この例のソースコードは、GitHub ウェブサイトの [join_and_relationalize.py Glue サンプルリポジトリ](#) の AWS Glue ファイルにあります。

このデータを使用して、このチュートリアルでは以下のことを実行する方法を示します。

- AWS Glue クローラを使用して、パブリックな Amazon S3 バケットに保存されているオブジェクトを分類し、それらのスキーマを AWS Glue Data Catalog に保存します。
- クロールの結果のテーブルのメタデータとスキーマを調べます。
- Data Catalog のメタデータを使用して Python の抽出、転送、およびロード (ETL) スクリプトを記述し、以下の操作を行います。
 - 異なるソースファイル内のデータをまとめて単一のデータテーブルに結合します (つまり、データを非正規化します)。
 - 議員のタイプ別に、結合テーブルを別のテーブルにフィルタリングします。

- 生成されたデータを後で分析するために Apache Parquet ファイルに分割して書き出します。

AWS で実行中に Python または PySpark スクリプトをデバッグするための推奨方法は、[AWS Glue Studio のノートブック](#)を使用することです。

ステップ 1: Amazon S3 バケット内のデータをクローलする

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. [クローラーの設定](#) の手順に従って、`s3://awsglue-datasets/examples/us-legislators/all` データセットをクローラできる新しいクローラを、AWS Glue Data Catalog 内のデータベース `legislators` に作成します。サンプルデータは既に、このパブリックな Amazon S3 バケットに用意されています。
3. 新しいクローラを実行し、`legislators` データベースを確認します。

クローラは、次のメタデータテーブルを作成します。

- `persons_json`
- `memberships_json`
- `organizations_json`
- `events_json`
- `areas_json`
- `countries_r_json`

これは、議員とその履歴を含むテーブルの半正規化されたテーブルの集合です。

ステップ 2: 開発エンドポイントノートブックに共通スクリプトを追加する

次の共通スクリプトを開発エンドポイントノートブックに貼り付けて、必要な AWS Glue ライブラリをインポートし、単一の `GlueContext` を設定します。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
```

```
from awsglue.job import Job

glueContext = GlueContext(SparkContext.getOrCreate())
```

ステップ 3: Data Catalog 内のデータでスキーマを確認する

次に、簡単な手順で AWS Glue Data Catalog から DynamicFrame を作成し、そのデータのスキーマを調べます。例えば、persons_json テーブルのスキーマを表示するには、ノートブックに以下を追加します。

```
persons = glueContext.create_dynamic_frame.from_catalog(
    database="legislators",
    table_name="persons_json")
print "Count: ", persons.count()
persons.printSchema()
```

プリントコールの出力を以下に示します。

```
Count: 1961
root
|-- family_name: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- gender: string
|-- image: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- name: string
|   |   |-- lang: string
|-- sort_name: string
|-- images: array
|   |-- element: struct
```

```
|      |      |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
|-- contact_details: array
|      |-- element: struct
|      |      |-- type: string
|      |      |-- value: string
|-- death_date: string
```

テーブル内の各人は、米国議会のメンバーです。

memberships_json テーブルのスキーマを表示するには、次のように入力します。

```
memberships = glueContext.create_dynamic_frame.from_catalog(
    database="legislators",
    table_name="memberships_json")
print "Count: ", memberships.count()
memberships.printSchema()
```

出力は次のとおりです。

```
Count: 10439
root
|-- area_id: string
|-- on_behalf_of_id: string
|-- organization_id: string
|-- role: string
|-- person_id: string
|-- legislative_period_id: string
|-- start_date: string
|-- end_date: string
```

organizations は政党および上院と下院の2つの議会です。organizations_json テーブルのスキーマを表示するには、次のように入力します。

```
orgs = glueContext.create_dynamic_frame.from_catalog(
    database="legislators",
```

```
        table_name="organizations_json")
print "Count: ", orgs.count()
orgs.printSchema()
```

出力は次のとおりです。

```
Count:  13
root
 |-- classification: string
 |-- links: array
 |   |-- element: struct
 |     |-- note: string
 |     |-- url: string
 |-- image: string
 |-- identifiers: array
 |   |-- element: struct
 |     |-- scheme: string
 |     |-- identifier: string
 |-- other_names: array
 |   |-- element: struct
 |     |-- lang: string
 |     |-- note: string
 |     |-- name: string
 |-- id: string
 |-- name: string
 |-- seats: int
 |-- type: string
```

ステップ 4: データをフィルタリングする

次に、必要なフィールドのみを保持し、`id` の名前を `org_id` に変更します。データセットは、小さいため全体を表示することができます。

`toDF()` は `DynamicFrame` を Apache Spark に変換するので、Apache Spark SQL に既に存在する `DataFrame` 変換を適用できます。

```
orgs = orgs.drop_fields(['other_names',
                        'identifiers']).rename_field(
    'id', 'org_id').rename_field(
    'name', 'org_name')
```

```
orgs.toDF().show()
```

以下に出力を示します。

```
+-----+-----+-----+-----+-----+
+-----+-----+
|classification|          org_id|          org_name|          links|seats|
|      type|          image|                    |                    |
+-----+-----+-----+-----+-----+
+-----+-----+
|      party|      party/al|          AL|          null| null|
|      null|      null|                    |                    |
|      party|      party/democrat|      Democrat|[[website,http://...| null|
|      null|https://upload.wi...|                    |                    |
|      party|party/democrat-li...|      Democrat-Liberal|[[website,http://...| null|
|      null|      null|                    |                    |
|      legislature|d56acebe-8fdc-47b...|House of Represen...|          null| 435|
|      lower house|      null|                    |                    |
|      party|      party/independent|      Independent|          null| null|
|      null|      null|                    |                    |
|      party|party/new_progres...|      New Progressive|[[website,http://...| null|
|      null|https://upload.wi...|                    |                    |
|      party|party/popular_dem...|      Popular Democrat|[[website,http://...| null|
|      null|      null|                    |                    |
|      party|      party/republican|      Republican|[[website,http://...| null|
|      null|https://upload.wi...|                    |                    |
|      party|party/republican-...|      Republican-Conser...|[[website,http://...| null|
|      null|      null|                    |                    |
|      party|      party/democrat|      Democrat|[[website,http://...| null|
|      null|https://upload.wi...|                    |                    |
|      party|      party/independent|      Independent|          null| null|
|      null|      null|                    |                    |
|      party|      party/republican|      Republican|[[website,http://...| null|
|      null|https://upload.wi...|                    |                    |
|      legislature|8fa6c3d2-71dc-478...|          Senate|          null| 100|
|      upper house|      null|                    |                    |
+-----+-----+-----+-----+-----+
+-----+-----+
```

memberships に表示される organizations を表示するには、次のように入力します。

```
memberships.select_fields(['organization_id']).toDF().distinct().show()
```

以下に出力を示します。

```
+-----+
| organization_id|
+-----+
|d56acebe-8fdc-47b...|
|8fa6c3d2-71dc-478...|
+-----+
```

ステップ 5: すべてをまとめる

ここで AWS Glue を使用して、これらのリレーショナルテーブルを結合し、議員の memberships とそれに対応する organizations の 1 つの完全な履歴テーブルを作成します。

1. まず、persons および memberships を id および person_id と結合します。
2. 次に、結果を orgs と org_id および organization_id と結合します。
3. 次に、冗長なフィールド person_id および org_id を削除します。

これらの操作はすべて、1 行の (拡張された) コードで行うことができます。

```
l_history = Join.apply(orgs,
                      Join.apply(persons, memberships, 'id', 'person_id'),
                      'org_id', 'organization_id').drop_fields(['person_id',
                        'org_id'])
print "Count: ", l_history.count()
l_history.printSchema()
```

出力は次のとおりです。

```
Count: 10439
root
 |-- role: string
 |-- seats: int
 |-- org_name: string
 |-- links: array
```

```
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- type: string
|-- sort_name: string
|-- area_id: string
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- on_behalf_of_id: string
|-- other_names: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- name: string
|   |   |-- lang: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- name: string
|-- birth_date: string
|-- organization_id: string
|-- gender: string
|-- classification: string
|-- death_date: string
|-- legislative_period_id: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- image: string
|-- given_name: string
|-- family_name: string
|-- id: string
|-- start_date: string
|-- end_date: string
```

これで、分析に使用できる最終テーブルが作成されました。これは、分析のためのコンパクトで効率的な形式 (つまり Parquet) で記述することができ、AWS Glue、Amazon Athena、または Amazon Redshift Spectrum で SQL を実行できます。

次の呼び出しは、複数のファイルにわたってテーブルを書き込んで、後で解析するときに高速な並列読み込みをサポートします。

```
glueContext.write_dynamic_frame.from_options(frame = l_history,  
      connection_type = "s3",  
      connection_options = {"path": "s3://glue-sample-target/output-dir/  
legislator_history"},  
      format = "parquet")
```

すべての履歴データを単一のファイルにまとめるには、データフレームに変換し、再パーティション化して書き出す必要があります。

```
s_history = l_history.toDF().repartition(1)  
s_history.write.parquet('s3://glue-sample-target/output-dir/legislator_single')
```

または、上院と下院でそれを分けたい場合。

```
l_history.toDF().write.parquet('s3://glue-sample-target/output-dir/legislator_part',  
      partitionBy=['org_name'])
```

ステップ 6: リレーショナルデータベース向けにデータを変換する

AWS Glue では半構造化データでも Amazon Redshift のようなリレーショナルデータベースに簡単に書き込むことができます。これにより、フレーム内のオブジェクトがどれほど複雑であっても、DynamicFrames をフラット化する変換 `relationalize` が提供されます。

この例の `l_history` DynamicFrame を使用して、ルートテーブル (`hist_root`) の名前と一時的な作業パスを `relationalize` に渡します。これにより、DynamicFrameCollection が返されます。その後、そのコレクション内の DynamicFrames の名前を一覧表示できます。

```
dfc = l_history.relationalize("hist_root", "s3://glue-sample-target/temp-dir/")  
dfc.keys()
```

`keys` 呼び出しの出力は次のとおりです。

```
[u'hist_root', u'hist_root_contact_details', u'hist_root_links',  
u'hist_root_other_names', u'hist_root_images', u'hist_root_identifiers']
```

Relationalize は、履歴テーブルを 6 つの新しいテーブルに分割します。DynamicFrame の各オブジェクトのレコードを含むルートテーブル、および配列の補助テーブルです。リレーショナルデータベースでの配列の処理は、特に配列が大きくなる場合に、最適ではないことがあります。配列を別のテーブルに分けることで、クエリの実行速度が大幅に向上します。

次に、contact_details を調べて分離を確認します。

```
l_history.select_fields('contact_details').printSchema()
dfc.select('hist_root_contact_details').toDF().where("id = 10 or id =
75").orderBy(['id', 'index']).show()
```

show 呼び出しの出力は次のとおりです。

```
root
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
+-----+-----+-----+-----+
| id|index|contact_details.val.type|contact_details.val.value|
+-----+-----+-----+-----+
| 10|  0|          fax|          |
| 10|  1|          |      202-225-1314|
| 10|  2|        phone|          |
| 10|  3|          |      202-225-3772|
| 10|  4|       twitter|          |
| 10|  5|          |      MikeRossUpdates|
| 75|  0|          fax|          |
| 75|  1|          |      202-225-7856|
| 75|  2|        phone|          |
| 75|  3|          |      202-225-2711|
| 75|  4|       twitter|          |
| 75|  5|          |      SenCapito|
+-----+-----+-----+-----+
```

contact_details フィールドは、元の DynamicFrame の構造体の配列です。これらの配列の各要素は、index によってインデックス化された、補助テーブルの個別の行です。ここで id は、contact_details キーを使用する hist_root テーブルの外部キーです。

```
dfc.select('hist_root').toDF().where(
    "contact_details = 10 or contact_details = 75").select(
    ['id', 'given_name', 'family_name', 'contact_details']).show()
```

出力を次に示します。

```
+-----+-----+-----+-----+
|          id|given_name|family_name|contact_details|
+-----+-----+-----+-----+
|f4fc30ee-7b42-432...|      Mike|      Ross|          10|
|e3c60f34-7d1b-4c0...|  Shelley|    Capito|          75|
+-----+-----+-----+-----+
```

これらのコマンドでは、`toDF()` および `where` 式を使用して、表示する行をフィルタリングすることに注意してください。

したがって、`hist_root` テーブルを補助テーブルと結合すると、次のことが可能になります。

- 配列をサポートせずにデータベースにデータをロードします。
- SQL を使用して配列内の各項目にクエリを実行します。

AWS Glue 接続を使用して、Amazon Redshift の認証情報を安全に保存してアクセスします。独自の接続の作成方法については、「[データへの接続](#)」を参照してください。

`DynamicFrames` を 1 つずつ切り替えて、接続にデータを書き込みできるようになりました。

```
for df_name in dfc.keys():
    m_df = dfc.select(df_name)
    print "Writing to table: ", df_name
    glueContext.write_dynamic_frame.from_jdbc_conf(frame = m_df, connection settings here)
```

接続設定は、リレーショナルデータベースのタイプによって異なります。

- Amazon Redshift への書き込み手順については、「[the section called “Redshift 接続”](#)」を参照してください。
- その他のデータベースについては、「[AWS Glue for Spark での ETL の接続タイプとオプション](#)」を参照してください。

結論

全体として、AWS Glue は非常に柔軟です。通常は書くのに数日かかるところを、数行のコードで達成できます。ソースからターゲットへの ETL スクリプトの全体は、GitHub の [AWS Glue サンプル](#)内の Python ファイル `join_and_relationalize.py` にあります。

コード例: `ResolveChoice`、`Lambda`、および `ApplyMapping` を使用したデータ準備

この例で使用されているデータセットは、2 つの [Data.CMS.gov](#) データセット (Inpatient Prospective Payment System Provider Summary for the Top 100 Diagnosis-Related Groups - FY2011 および Inpatient Charge Data FY 2011) からダウンロードされた、メディケアプロバイダの支払いデータで構成されています。ダウンロードした後、データセットを修正してファイルの最後にいくつかの誤ったレコードを追加しました。この変更されたファイルは、パブリックな Amazon S3 バケット (`s3://awsglue-datasets/examples/medicare/Medicare_Hospital_Provider.csv`) 内に置かれています。

この例のソースコードは、`data_cleaning_and_lambda.py` [例AWS Glue GitHub リポジトリ](#)のファイルにあります。

AWS で実行中に Python または PySpark スクリプトをデバッグするための推奨方法は、[AWS Glue Studio のノートブック](#)を使用することです。

ステップ 1: Amazon S3 バケット内のデータをクローलする

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. [クローラーの設定](#) で説明されているプロセスに従って `s3://awsglue-datasets/examples/medicare/Medicare_Hospital_Provider.csv` ファイルをクロールできる新しいクローラを作成し、得られた結果のメタデータを AWS Glue Data Catalog の `payments` という名前のデータベースに配置します。
3. 新しいクローラを実行し、`payments` データベースを確認します。クローラは、最初のファイルを読み込んでファイルの形式と区切り記号を判断してから、データベースに `medicare` という名前のメタデータテーブルを作成したことがわかります。

新しい `medicare` テーブルのスキーマは次のようになります。

Column name	Data type
=====	
<code>drg definition</code>	<code>string</code>
<code>provider id</code>	<code>bigint</code>

```
provider name                string
provider street address     string
provider city                string
provider state               string
provider zip code            bigint
hospital referral region description string
total discharges             bigint
average covered charges      string
average total payments       string
average medicare payments    string
```

ステップ 2: 開発エンドポイントノートブックに共通スクリプトを追加する

次の共通スクリプトを開発エンドポイントノートブックに貼り付けて、必要な AWS Glue ライブラリをインポートし、単一の GlueContext を設定します。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job

glueContext = GlueContext(SparkContext.getOrCreate())
```

ステップ 3: 異なるスキーマ解析を比較する

次に、Apache Spark DataFrame によって認識されたスキーマが、AWS Glue クローラによって記録されたスキーマと同じかどうかを確認できます。以下のコードを実行します。

```
medicare = spark.read.format(
    "com.databricks.spark.csv").option(
    "header", "true").option(
    "inferSchema", "true").load(
    's3://awsglue-datasets/examples/medicare/Medicare_Hospital_Provider.csv')
medicare.printSchema()
```

printSchema 呼び出しの出力は次のとおりです。

```
root
|-- DRG Definition: string (nullable = true)
|-- Provider Id: string (nullable = true)
|-- Provider Name: string (nullable = true)
|-- Provider Street Address: string (nullable = true)
|-- Provider City: string (nullable = true)
|-- Provider State: string (nullable = true)
|-- Provider Zip Code: integer (nullable = true)
|-- Hospital Referral Region Description: string (nullable = true)
|-- Total Discharges : integer (nullable = true)
|-- Average Covered Charges : string (nullable = true)
|-- Average Total Payments : string (nullable = true)
|-- Average Medicare Payments: string (nullable = true)
```

次に、AWS Glue DynamicFrame によって生成されるスキーマを確認します。

```
medicare_dynamicframe = glueContext.create_dynamic_frame.from_catalog(
    database = "payments",
    table_name = "medicare")
medicare_dynamicframe.printSchema()
```

printSchema の出力は次のとおりです。

```
root
|-- drg definition: string
|-- provider id: choice
|   |-- long
|   |-- string
|-- provider name: string
|-- provider street address: string
|-- provider city: string
|-- provider state: string
|-- provider zip code: long
|-- hospital referral region description: string
|-- total discharges: long
|-- average covered charges: string
|-- average total payments: string
|-- average medicare payments: string
```

DynamicFrame は、provider id が long 型または string 型のいずれかであるスキーマを生成します。DataFrame スキーマは Provider Id を string 型としてリストし、Data Catalog は provider id を bigint 型としてリストします。

正しいものはどちらでしょうか。ファイルの末尾には、その列に string 値を持つ 2 つのレコード (160,000 レコードのうち) があります。これらは、問題を説明するために導入されたエラーのあるレコードです。

このような問題に対処するために、AWS Glue DynamicFrame では Choice 型の概念を導入しています。この場合、DynamicFrame は、その列に long 値と string 値の両方が存在することを示しています。AWS Glue クローラはデータの 2 MB のプレフィックスのみを考慮しているため、string 値を見落とししました。Apache Spark DataFrame はデータセット全体を考慮しましたが、最も一般的な型、つまり string 型を強制的に列に割り当てました。実際、慣れていない複雑な型やバリエーションがある場合にも、Spark は最も一般的なケースを使用することがあります。

provider id 列のクエリを実行するには、Choice 型をまず解決する必要があります。DynamicFrame で、cast:long オプションを指定して resolveChoice 変換メソッドを使用すると、これらの string 値を long 値に変換できます。

```
medicare_res = medicare_dynamicframe.resolveChoice(specs = [('provider
  id', 'cast:long']))
medicare_res.printSchema()
```

この場合、printSchema の出力は次のようになります。

```
root
 |-- drg definition: string
 |-- provider id: long
 |-- provider name: string
 |-- provider street address: string
 |-- provider city: string
 |-- provider state: string
 |-- provider zip code: long
 |-- hospital referral region description: string
 |-- total discharges: long
 |-- average covered charges: string
 |-- average total payments: string
 |-- average medicare payments: string
```

値がキャストできない string だった場合に、AWS Glue は null を挿入しました。

もう 1 つのオプションは、両方のタイプの値を保持する struct に Choice 型を変換することです。

次に、異常だった行を確認してみましょう。

```
medicare_res.toDF().where("'provider id' is NULL").show()
```

次のように表示されています。

```
+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+
|      drg definition|provider id|  provider name|provider street address|provider
city|provider state|provider zip code|hospital referral region description|total
discharges|average covered charges|average total payments|average medicare payments|
+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+
|948 - SIGNS & SYM...|      null|          INC|      1050 DIVISION ST|
MAUSTON|          WI|          53948|          WI - Madison|
      12|          $11961.41|          $4619.00|          $3775.33|
|948 - SIGNS & SYM...|      null| INC- ST JOSEPH|      5000 W CHAMBERS ST|
MILWAUKEE|          WI|          53210|          WI - Milwaukee|
      14|          $10514.28|          $5562.50|          $4522.78|
+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+
```

次のように、2 つの不正な形式のレコードを削除します。

```
medicare_dataframe = medicare_res.toDF()
medicare_dataframe = medicare_dataframe.where("'provider id' is NOT NULL")
```

ステップ 4: データのマッピングと Apache Spark Lambda 関数の使用

AWS Glue では、まだユーザー定義関数とも呼ばれる Lambda 関数が直接サポートされていません。しかし、いつでも DynamicFrame を Apache Spark DataFrame との間で変換して、DynamicFrames の特殊な機能に加えて Spark の機能を利用できます。

次に、支払い情報を数字に変換すると、Amazon Redshift や Amazon Athena のような分析エンジンが、より迅速に数値処理を実行できるようになります。

```

from pyspark.sql.functions import udf
from pyspark.sql.types import StringType

chop_f = udf(lambda x: x[1:], StringType())
medicare_dataframe = medicare_dataframe.withColumn(
    "ACC", chop_f(
        medicare_dataframe["average covered charges"]))
medicare_dataframe.withColumn(
    "ATP", chop_f(
        medicare_dataframe["average total payments"]))
medicare_dataframe.withColumn(
    "AMP", chop_f(
        medicare_dataframe["average medicare payments"]))
medicare_dataframe.select(['ACC', 'ATP', 'AMP']).show()

```

show 呼び出しの出力は次のとおりです。

```

+-----+-----+-----+
|  ACC|  ATP|  AMP|
+-----+-----+-----+
|32963.07|5777.24|4763.73|
|15131.85|5787.57|4976.71|
|37560.37|5434.95|4453.79|
|13998.28|5417.56|4129.16|
|31633.27|5658.33|4851.44|
|16920.79|6653.80|5374.14|
|11977.13|5834.74|4761.41|
|35841.09|8031.12|5858.50|
|28523.39|6113.38|5228.40|
|75233.38|5541.05|4386.94|
|67327.92|5461.57|4493.57|
|39607.28|5356.28|4408.20|
|22862.23|5374.65|4186.02|
|31110.85|5366.23|4376.23|
|25411.33|5282.93|4383.73|
| 9234.51|5676.55|4509.11|
|15895.85|5930.11|3972.85|
|19721.16|6192.54|5179.38|
|10710.88|4968.00|3898.88|
|51343.75|5996.00|4962.45|
+-----+-----+-----+
only showing top 20 rows

```

これらは、すべてデータ内ではまだ文字列です。強力な `apply_mapping` 変換メソッドを使用して、データをドロップ、名前変更、キャスト、およびネストし、他のデータプログラミング言語やシステムで容易にアクセスできるようにします。

```
from awsglue.dynamicframe import DynamicFrame
medicare_tmp_dyf = DynamicFrame.fromDF(medicare_dataframe, glueContext, "nested")
medicare_nest_dyf = medicare_tmp_dyf.apply_mapping([('drg definition', 'string', 'drg',
'string'),
          ('provider id', 'long', 'provider.id', 'long'),
          ('provider name', 'string', 'provider.name', 'string'),
          ('provider city', 'string', 'provider.city', 'string'),
          ('provider state', 'string', 'provider.state', 'string'),
          ('provider zip code', 'long', 'provider.zip', 'long'),
          ('hospital referral region description', 'string', 'rr', 'string'),
          ('ACC', 'string', 'charges.covered', 'double'),
          ('ATP', 'string', 'charges.total_pay', 'double'),
          ('AMP', 'string', 'charges.medicare_pay', 'double')])
medicare_nest_dyf.printSchema()
```

`printSchema` の出力は次のとおりです。

```
root
 |-- drg: string
 |-- provider: struct
 |   |-- id: long
 |   |-- name: string
 |   |-- city: string
 |   |-- state: string
 |   |-- zip: long
 |-- rr: string
 |-- charges: struct
 |   |-- covered: double
 |   |-- total_pay: double
 |   |-- medicare_pay: double
```

データを Spark DataFrame に戻すと、現在どのような状態かが分かります。

```
medicare_nest_dyf.toDF().show()
```

出力は次のとおりです。

```
+-----+-----+-----+-----+
```

	drg	provider	rr	charges
039 - EXTRACRANIA...	[10001,SOUTHEAST ...	AL - Dothan	[32963.07,5777.24...	
039 - EXTRACRANIA...	[10005,MARSHALL M...	AL - Birmingham	[15131.85,5787.57...	
039 - EXTRACRANIA...	[10006,ELIZA COFF...	AL - Birmingham	[37560.37,5434.95...	
039 - EXTRACRANIA...	[10011,ST VINCENT...	AL - Birmingham	[13998.28,5417.56...	
039 - EXTRACRANIA...	[10016,SHELBY BAP...	AL - Birmingham	[31633.27,5658.33...	
039 - EXTRACRANIA...	[10023,BAPTIST ME...	AL - Montgomery	[16920.79,6653.8,...	
039 - EXTRACRANIA...	[10029,EAST ALABA...	AL - Birmingham	[11977.13,5834.74...	
039 - EXTRACRANIA...	[10033,UNIVERSITY...	AL - Birmingham	[35841.09,8031.12...	
039 - EXTRACRANIA...	[10039,HUNTSVILLE...	AL - Huntsville	[28523.39,6113.38...	
039 - EXTRACRANIA...	[10040,GADSDEN RE...	AL - Birmingham	[75233.38,5541.05...	
039 - EXTRACRANIA...	[10046,RIVERVIEW ...	AL - Birmingham	[67327.92,5461.57...	
039 - EXTRACRANIA...	[10055,FLOWERS HO...	AL - Dothan	[39607.28,5356.28...	
039 - EXTRACRANIA...	[10056,ST VINCENT...	AL - Birmingham	[22862.23,5374.65...	
039 - EXTRACRANIA...	[10078,NORTHEAST ...	AL - Birmingham	[31110.85,5366.23...	
039 - EXTRACRANIA...	[10083,SOUTH BALD...	AL - Mobile	[25411.33,5282.93...	
039 - EXTRACRANIA...	[10085,DECATUR GE...	AL - Huntsville	[9234.51,5676.55,...	
039 - EXTRACRANIA...	[10090,PROVIDENCE...	AL - Mobile	[15895.85,5930.11...	
039 - EXTRACRANIA...	[10092,D C H REGI...	AL - Tuscaloosa	[19721.16,6192.54...	
039 - EXTRACRANIA...	[10100,THOMAS HOS...	AL - Mobile	[10710.88,4968.0,...	
039 - EXTRACRANIA...	[10103,BAPTIST ME...	AL - Birmingham	[51343.75,5996.0,...	

only showing top 20 rows

ステップ 5: Apache Parquet にデータを書き込む

AWS Glue は、リレーショナルデータベースが効果的に消費できる Apache Parquet のような形式でデータを書き込むことを容易にします。

```
glueContext.write_dynamic_frame.from_options(
    frame = medicare_nest_dyf,
    connection_type = "s3",
    connection_options = {"path": "s3://glue-sample-target/output-dir/
medicare_parquet"},
    format = "parquet")
```

AWS Glue PySpark 拡張機能リファレンス

AWS Glue は PySpark Python 方言に以下の拡張機能を作成しました。

- [getResolvedOptions](#) を使用して、パラメータにアクセスする

- [PySpark 拡張子型](#)
- [DynamicFrame クラス](#)
- [DynamicFrameCollection クラス](#)
- [DynamicFrameWriter クラス](#)
- [DynamicFrameReader クラス](#)
- [GlueContext クラス](#)

getResolvedOptions を使用して、パラメータにアクセスする

AWS Glue `getResolvedOptions(args, options)` ユーティリティ関数を使用すると、ジョブの実行時にスクリプトに渡される引数にアクセスできます。この関数を使用するには、まず AWS Glue `utils` モジュールと `sys` モジュールからインポートします。

```
import sys
from awsglue.utils import getResolvedOptions
```

getResolvedOptions(args, options)

- `args` - `sys.argv` に含まれる引数のリスト。
- `options` - 取得したい引数名の Python 配列。

Example JobRun に渡された引数を取得する

スクリプト (Lambda 関数など) の中で、JobRun を作成した場合を考えます。

```
response = client.start_job_run(
    JobName = 'my_test_job',
    Arguments = {
        '--day_partition_key': 'partition_0',
        '--hour_partition_key': 'partition_1',
        '--day_partition_value': day_partition_value,
        '--hour_partition_value': hour_partition_value } )
```

渡された引数を取得するには、次のように、`getResolvedOptions` 関数を使用できます。

```
import sys
```

```
from awsglue.utils import getResolvedOptions

args = getResolvedOptions(sys.argv,
                          ['JOB_NAME',
                           'day_partition_key',
                           'hour_partition_key',
                           'day_partition_value',
                           'hour_partition_value'])

print "The day-partition key is: ", args['day_partition_key']
print "and the day-partition value is: ", args['day_partition_value']
```

各引数は 2 つのハイフンで始まり、ハイフンなしでスクリプト内で参照されるように定義されていることに注意してください。引数では、ハイフンではなく、アンダースコアのみを使用します。引数の解決のために、この規則に従う必要があります。

PySpark 拡張子型

AWS Glue PySpark 拡張子で使用される型。

DataType

他の AWS Glue 型の基本クラス。

__init__(properties={})

- `properties` – データ型のプロパティ (オプション)。

typeName(cls)

AWS Glue 型クラスの種類 (つまり、"型" が末尾から削除されたクラス名) を返します。

- `cls` – AWS Glue から派生した `DataType` クラスインスタンス。

jsonValue()

データ型とクラスのプロパティが含まれる JSON オブジェクトを返します。

```
{
  "dataType": typeName,
```

```
"properties": properties
}
```

AtomicType およびシンプルデリバティブ

[DataType](#) クラスから継承して拡張し、すべての AWS Glue アトミックデータ型の基本クラスとして機能します。

fromJsonValue(cls, json_value)

JSON オブジェクトからの値を使用してクラスインスタンスを初期化します。

- `cls` – 初期化する AWS Glue 型のクラスインスタンス。
- `json_value` – キーと値のペアのロード元の JSON オブジェクト。

次の型は、[AtomicType](#) クラスのシンプルデリバティブです。

- `BinaryType` – バイナリデータ。
- `BooleanType` – ブール値。
- `ByteType` – バイト値。
- `DateType` – 日時値。
- `DoubleType` – 倍精度浮動小数点値。
- `IntegerType` – 整数値。
- `LongType` – 長整数値。
- `NullType` – null 値。
- `ShortType` – 短整数値。
- `StringType` – テキスト文字列。
- `TimestampType` – タイムスタンプ値 (通常は 1970 年 1 月 1 日からの秒数)。
- `UnknownType` – 未確認型の値。

DecimalType(AtomicType)

10 進数 (バイナリ 2 進数ではなく、10 進数で表記される数) を表わすため、[AtomicType](#) クラスから継承して拡張します。

__init__(precision=10, scale=2, properties={})

- `precision` – 10 進数の桁数 (オプション、デフォルトは 10)。
- `scale` – 小数点以下の桁数 (オプション、デフォルトは 2)。
- `properties` – 10 進数のプロパティ (オプション)。

EnumType(AtomicType)

有効なオプションの列挙を表すために、[AtomicType](#) クラスから継承して拡張します。

__init__(options)

- `options` – 列挙されているオプションのリスト。

コレクション型

- [ArrayType\(DataType\)](#)
- [ChoiceType\(DataType\)](#)
- [MapType\(DataType\)](#)
- [フィールド \(オブジェクト\)](#)
- [StructType\(DataType\)](#)
- [EntityType\(DataType\)](#)

ArrayType(DataType)

__init__(elementType=UnknownType(), properties={})

- `elementType` – 配列の要素の型 (オプション、デフォルトは UnknownType)。
- `properties` – 配列のプロパティ (オプション)。

ChoiceType(DataType)

__init__(choices=[], properties={})

- `choices` – 選択肢のリスト (オプション)。
- `properties` – これらのオプションのプロパティ (オプション)。

add(new_choice)

可能な選択肢のリストに、新しい選択肢を追加します。

- `new_choice` – 可能な選択肢のリストに追加する選択肢。

merge(new_choices)

新しい選択肢のリストを既存の選択肢のリストとマージします。

- `new_choices` – 既存の選択肢とマージする新しい選択肢のリスト。

MapType(DataType)

__init__(valueType=UnknownType, properties={})

- `valueType` – マップの値の型 (オプション、デフォルトは `UnknownType`)。
- `properties` – マップのプロパティ (オプション)。

フィールド (オブジェクト)

[DataType](#) から派生したオブジェクトからフィールドオブジェクトを作成します。

__init__(name, dataType, properties={})

- `name` – フィールドに割り当てる名前。
- `dataType` – フィールド作成元のオブジェクト。
- `properties` – フィールドのプロパティ (オプション)。

StructType(DataType)

データ構造を定義します (`struct`)。

__init__(fields=[], properties={})

- `fields` – フィールドのリスト (`Field` 型) 構造に含めます (オプション)。
- `properties` – 構造のプロパティ (オプション)。

add(field)

- field – 構造に追加するオブジェクトの Field 型。

hasField(field)

この構造に同じ名前のフィールドがある場合は True を、そうでない場合は False を返します。

- field – フィールド名、または名前が使用される Field 型のオブジェクト。

getField(field)

- field – 名前が使用される Field 型のフィールド名またはオブジェクト。構造に同じ名前のフィールドがある場合は、返されます。

EntityType(DataType)

`__init__(entity, base_type, properties)`

このクラスは、まだ実装されていません。

その他の型

- [DataSource \(オブジェクト\)](#)
- [DataSink \(オブジェクト\)](#)

DataSource (オブジェクト)

`__init__(j_source, sql_ctx, name)`

- j_source – データソース。
- sql_ctx – SQL コンテキスト。
- name – データソース名。

setFormat(format, **options)

- format – データソースを設定する形式。

- `options` – データソースに設定するオプションのコレクション。形式オプションの詳細については、「[the section called “データ形式に関するオプション”](#)」を参照してください。

`getFrame()`

データソースに `DynamicFrame` を返します。

`DataSink` (オブジェクト)

`__init__(j_sink, sql_ctx)`

- `j_sink` – 作成するシンク。
- `sql_ctx` – データシンクの SQL コンテキスト。

`setFormat(format, **options)`

- `format` – データシンクを設定する形式。
- `options` – データシンクに設定するオプションのコレクション。形式オプションの詳細については、「[the section called “データ形式に関するオプション”](#)」を参照してください。

`setAccumulableSize(size)`

- `size` – 設定する累積サイズ (バイト単位)。

`writeFrame(dynamic_frame, info="")`

- `dynamic_frame` – 書き込む `DynamicFrame`。
- `info` – `DynamicFrame` に関する情報 (オプション)。

`write(dynamic_frame_or_dfc, info="")`

`DynamicFrame` または `DynamicFrameCollection` を書き込みます。

- `dynamic_frame_or_dfc` – 書き込む `DynamicFrame` オブジェクトまたは `DynamicFrameCollection` オブジェクトのいずれか。

- `info` – 書き込む `DynamicFrame` または `DynamicFrames` に関する情報 (オプション)。

DynamicFrame クラス

Apache Spark の主要な抽象化の 1 つは SparkSQL `DataFrame` で、これは R と Pandas にある `DataFrame` 構造に似ています。`DataFrame` はテーブルと似ており、機能スタイル (マップ/リデュース/フィルター/その他) 操作と SQL 操作 (選択、プロジェクト、集計) をサポートしています。

`DataFrames` は、強力で広く使用されていますが、抽出、変換、ロード (ETL) 操作に関しては制限があります。最も重要なのは、データをロードする前にスキーマを指定する必要があることです。SparkSQL は、データに対してパスを 2 つ作ることによってこれを解決します。この 1 つ目はスキーマの推定を行い、2 つ目はデータをロードします。ただし、この推測は限定されており、実際の煩雑なデータには対応しません。例えば、同じフィールドが異なるレコードの異なるタイプである可能性があります。Apache Spark は、多くの場合、作業を中断して、元のフィールドテキストを使用して型を `string` として報告します。これは正しくない可能性があり、スキーマの不一致を解決する方法を細かくコントロールする必要があるかもしれません。また、大規模なデータセットの場合、ソースデータに対する追加パスが非常に高価になる可能性があります。

これらの制限に対応するために、AWS Glue では `DynamicFrame` を導入しています。`DynamicFrame` は、`DataFrame` と似ていますが、各レコードが自己記述できるため、最初はスキーマは必要ありません。代わりに、AWS Glue は必要に応じてオンザフライでスキーマを計算し、選択 (または共用) タイプを使用してスキーマの不一致を明示的にエンコードします。これらの不整合を解決して、固定スキーマを必要とするデータストアとデータセットを互換性のあるものにできます。

同様に、`DynamicRecord` は `DynamicFrame` 内の論理レコードを表します。これは、Spark `DataFrame` の行と似ていますが、自己記述型であり、固定スキーマに適合しないデータに使用できます。PySpark で AWS Glue を使用する場合、通常は独立した操作 `DynamicRecords` は行いません。むしろ、データセットをその `DynamicFrame` を使ってまとめて変換します。

スキーマの不一致を解決したら、`DynamicFrames` を `DataFrames` との間で変換することができます。

— construction —

- [__init__](#)
- [fromDF](#)
- [toDF](#)

`__init__`

`__init__(jdf, glue_ctx, name)`

- `jdf` - Java 仮想マシン (JVM) 内のデータフレームへの参照。
- `glue_ctx` - [GlueContext クラス](#) オブジェクト。
- `name` - オプションの名前文字列。デフォルトでは空。

`fromDF`

`fromDF(dataframe, glue_ctx, name)`

DataFrame フィールドを DynamicRecord に変換することにより、DataFrame を DynamicFrame に変換します。新しい DynamicFrame を返します。

DynamicRecord は DynamicFrame 内の論理レコードを表します。これは、自己記述型であり、固定スキーマに適合しないデータに使用できる点を除いて、Spark DataFrame の行に似ています。

この関数は、DataFrame で名前と重複する列がすでに解決されていることを前提としています。

- `dataframe` - 変換する Apache Spark SQL DataFrame (必須)。
- `glue_ctx` - この変換のコンテキストを指定する [GlueContext クラス](#) オブジェクト (必須)。
- `name` - 結果 DynamicFrame の名前 (AWS Glue 3.0 以降はオプション)。

`toDF`

`toDF(options)`

DynamicRecords を DataFrame フィールドに変換することにより、DynamicFrame を Apache Spark DataFrame に変換します。新しい DataFrame を返します。

DynamicRecord は DynamicFrame 内の論理レコードを表します。これは、自己記述型であり、固定スキーマに適合しないデータに使用できる点を除いて、Spark DataFrame の行に似ています。

- `options` - オプションのリスト。Project と Cast アクションタイプを選択した場合、ターゲットのタイプを指定します。次に例を示します。

```
>>>toDF([ResolveOption("a.b.c", "KeepAsStruct")])
>>>toDF([ResolveOption("a.b.c", "Project", DoubleType())])
```

— information —

- [count](#)
- [スキーマ](#)
- [printSchema](#)
- [show](#)
- [repartition](#)
- [coalesce](#)

count

count() - 基盤となる DataFrame の行数を返します。

スキーマ

schema() - この DynamicFrame のスキーマを返します。使用できない場合は、基盤となる DataFrame のスキーマを返します。

このスキーマを構成する DynamicFrame タイプの詳細については、「[the section called “型”](#)」を参照してください。

printSchema

printSchema() - 基盤となる DataFrame のスキーマを表示します。

show

show(num_rows) - 基盤となる DataFrame から、指定された行数を表示します。

repartition

repartition(numPartitions) - numPartitions 個のパーティションを含む新しい DynamicFrame を返します。

coalesce

coalesce(numPartitions) - numPartitions 個のパーティションを含む新しい DynamicFrame を返します。

— transforms —

- [apply_mapping](#)

- [drop_fields](#)
- [フィルター](#)
- [join](#)
- [マップ](#)
- [mergeDynamicFrame](#)
- [関係付け](#)
- [rename_field](#)
- [resolveChoice](#)
- [select_fields](#)
- [スピゴット](#)
- [split_fields](#)
- [split_rows](#)
- [アンボックス](#)
- [the section called “union”](#)
- [ネスト解除](#)
- [unnest_ddb_json](#)
- [書き込み](#)

apply_mapping

apply_mapping(mappings, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)

DynamicFrame に宣言型のマッピングを適用し、それらのマッピングが適用された新しい DynamicFrame を指定したフィールドに返します。未指定のフィールドは新しい DynamicFrame から除外されます。

- mappings — マッピングタプルのリスト (必須)。それぞれが (ソース列、ソースタイプ、ターゲット列、ターゲットタイプ) で構成されます。

ソース列で、名前にドット「.」が含まれている場合、バックティック「`」で囲む必要があります。例えば、this.old.name (文字列) を thisNewName にマッピングするには、次のタプルを使用します。

```
("`this.old.name`", "string", "thisNewName", "string")
```

- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - この変換のエラー報告に関連付ける文字列 (オプション)。
- `stageThreshold` - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- `totalThreshold` - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。

例: `apply_mapping` を使用してフィールドの名前を変更し、フィールドタイプを変更する

次のコード例は、`apply_mapping` メソッドを使用して、選択したフィールドの名前を変更し、フィールドタイプを変更する方法を示しています。

Note

この例で使用されているデータセットにアクセスするには、「[コード例: データの結合と関係付け](#)」を参照し、「[ステップ 1: Amazon S3 バケット内のデータをクローलする](#)」の手順に従います。

```
# Example: Use apply_mapping to reshape source data into
# the desired column names and types as a new DynamicFrame

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Create a DynamicFrame and view its schema
persons = glueContext.create_dynamic_frame.from_catalog(
    database="legislators", table_name="persons_json"
)
print("Schema for the persons DynamicFrame:")
persons.printSchema()
```

```
# Select and rename fields, change field type
print("Schema for the persons_mapped DynamicFrame, created with apply_mapping:")
persons_mapped = persons.apply_mapping(
    [
        ("family_name", "String", "last_name", "String"),
        ("name", "String", "first_name", "String"),
        ("birth_date", "String", "date_of_birth", "Date"),
    ]
)
persons_mapped.printSchema()
```

出力

```
Schema for the persons DynamicFrame:
root
|-- family_name: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- gender: string
|-- image: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- sort_name: string
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
```

```
|    |    |-- value: string
|-- death_date: string
```

Schema for the persons_mapped DynamicFrame, created with apply_mapping:

```
root
|-- last_name: string
|-- first_name: string
|-- date_of_birth: date
```

drop_fields

drop_fields(paths, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)

[FlatMap クラス](#) 変換を呼び出して、DynamicFrame からフィールドを削除します。指定されたフィールドが削除された新しい DynamicFrame を返します。

- paths - 文字列のリスト。それぞれに、ドロップするフィールドノードへのフルパスが含まれます。ドット表記を使用して、ネストされたフィールドを指定できます。例えば、フィールド first がツリー内のフィールド name の子である場合は、パスに "name.first" を指定します。

フィールドノードの名前にリテラル . が含まれている場合は、その名前をバックティック (`) で囲む必要があります。

- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - この変換のエラー報告に関連付ける文字列 (オプション)。
- stageThreshold - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- totalThreshold - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。

例: drop_fields を使用して **DynamicFrame** からフィールドを削除する

このコード例では、drop_fields メソッドを使用して、選択したトップレベルフィールドとネストされたフィールドを DynamicFrame から削除します。

データセットの例

この例では、コード内の EXAMPLE-FRIENDS-DATA テーブルで表される次のデータセットを使用します。

```

{"name": "Sally", "age": 23, "location": {"state": "WY", "county": "Fremont"},
 "friends": []}
{"name": "Varun", "age": 34, "location": {"state": "NE", "county": "Douglas"},
 "friends": [{"name": "Arjun", "age": 3}]}
{"name": "George", "age": 52, "location": {"state": "NY"}, "friends": [{"name":
 "Fred"}, {"name": "Amy", "age": 15}]}
{"name": "Haruki", "age": 21, "location": {"state": "AK", "county": "Denali"}}
{"name": "Sheila", "age": 63, "friends": [{"name": "Nancy", "age": 22}]}

```

コードの例

```

# Example: Use drop_fields to remove top-level and nested fields from a DynamicFrame.
# Replace MY-EXAMPLE-DATABASE with your Glue Data Catalog database name.
# Replace EXAMPLE-FRIENDS-DATA with your table name.

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Create a DynamicFrame from Glue Data Catalog
glue_source_database = "MY-EXAMPLE-DATABASE"
glue_source_table = "EXAMPLE-FRIENDS-DATA"

friends = glueContext.create_dynamic_frame.from_catalog(
    database=glue_source_database, table_name=glue_source_table
)
print("Schema for friends DynamicFrame before calling drop_fields:")
friends.printSchema()

# Remove location.county, remove friends.age, remove age
friends = friends.drop_fields(paths=["age", "location.county", "friends.age"])
print("Schema for friends DynamicFrame after removing age, county, and friend age:")
friends.printSchema()

```

出力

```

Schema for friends DynamicFrame before calling drop_fields:
root
 |-- name: string

```

```
|-- age: int
|-- location: struct
|   |-- state: string
|   |-- county: string
|-- friends: array
|   |-- element: struct
|   |   |-- name: string
|   |   |-- age: int
```

Schema for friends DynamicFrame after removing age, county, and friend age:

```
root
|-- name: string
|-- location: struct
|   |-- state: string
|-- friends: array
|   |-- element: struct
|   |   |-- name: string
```

フィルター

filter(f, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)

指定された述語関数 *f* を満たす入力 DynamicFrame 内のすべての DynamicRecords を含む、新しい DynamicFrame を返します。

- *f* - DynamicFrame に適用する述語関数。この関数は DynamicRecord を引数として取り、DynamicRecord がフィルター要件を満たす場合は True を返し、そうでない場合は False を返します (必須)。

DynamicRecord は DynamicFrame 内の論理レコードを表します。これは、自己記述型であり、固定スキーマに準拠しないデータに使用できる点を除いて、Spark DataFrame の行に似ています。

- *transformation_ctx* - 状態情報を識別するために使用される一意の文字列 (オプション)。
- *info* - この変換のエラー報告に関連付ける文字列 (オプション)。
- *stageThreshold* - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- *totalThreshold* - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。

例: フィルターを使用して、フィールドのフィルタリングされた選択内容を取得する

この例では、`filter` メソッドを使用して、別の `DynamicFrame` のフィールドのフィルタリングされた選択内容を含む新しい `DynamicFrame` を作成します。

`map` メソッドと同様に、`filter` は、元の `DynamicFrame` の各レコードに適用される引数として関数を取ります。この関数はレコードを入力として受け取り、ブール値を返します。戻り値が `true` の場合、レコードは結果として生じる `DynamicFrame` に含まれます。`false` の場合、レコードは除外されます。

Note

この例で使用されているデータセットにアクセスするには、[「コード例: ResolveChoice、Lambda、および ApplyMapping を使用したデータ準備」](#)を参照し、[「ステップ 1: Amazon S3 バケット内のデータをクローलする」](#)の手順に従います。

```
# Example: Use filter to create a new DynamicFrame
# with a filtered selection of records

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Create DynamicFrame from Glue Data Catalog
medicare = glueContext.create_dynamic_frame.from_options(
    "s3",
    {
        "paths": [
            "s3://awsglue-datasets/examples/medicare/Medicare_Hospital_Provider.csv"
        ]
    },
    "csv",
    {"withHeader": True},
)

# Create filtered DynamicFrame with custom lambda
# to filter records by Provider State and Provider City
sac_or_mon = medicare.filter(
```

```
f=lambda x: x["Provider State"] in ["CA", "AL"]
and x["Provider City"] in ["SACRAMENTO", "MONTGOMERY"]
)

# Compare record counts
print("Unfiltered record count: ", medicare.count())
print("Filtered record count:  ", sac_or_mon.count())
```

出力

```
Unfiltered record count: 163065
Filtered record count: 564
```

join

```
join(paths1, paths2, frame2, transformation_ctx="", info="",
stageThreshold=0, totalThreshold=0)
```

別の DynamicFrame と等価結合を実行し、結果の DynamicFrame を返します。

- paths1 - 結合するこのフレームのキーのリスト。
- paths2 - 結合する別のフレームのキーのリスト。
- frame2 - 結合する他の DynamicFrame。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - この変換のエラー報告に関連付ける文字列 (オプション)。
- stageThreshold - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- totalThreshold - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。

例: 結合を使用して DynamicFrames を結合する

この例では、join メソッドを使用して 3 つの DynamicFrames で結合を実行します。AWSGlue は、入力されたフィールドキーに基づいて結合を実行します。結果として生じる DynamicFrame には、指定されたキーが一致する元の 2 つのフレームからの行が含まれます。

join 変換では、すべてのフィールドがそのまま保持されることに注意してください。これは、一致するように指定したフィールドが冗長で同じキーが含まれていても、結果として生じる

DynamicFrame に表示されることを意味します。この例では、`drop_fields` を使用して、結合後にこれらの冗長なキーを削除します。

Note

この例で使用されているデータセットにアクセスするには、「[コード例: データの結合と関係付け](#)」を参照し、「[ステップ 1: Amazon S3 バケット内のデータをクロールする](#)」の手順に従います。

```
# Example: Use join to combine data from three DynamicFrames

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Load DynamicFrames from Glue Data Catalog
persons = glueContext.create_dynamic_frame.from_catalog(
    database="legislators", table_name="persons_json"
)
memberships = glueContext.create_dynamic_frame.from_catalog(
    database="legislators", table_name="memberships_json"
)
orgs = glueContext.create_dynamic_frame.from_catalog(
    database="legislators", table_name="organizations_json"
)
print("Schema for the persons DynamicFrame:")
persons.printSchema()
print("Schema for the memberships DynamicFrame:")
memberships.printSchema()
print("Schema for the orgs DynamicFrame:")
orgs.printSchema()

# Join persons and memberships by ID
persons_memberships = persons.join(
    paths1=["id"], paths2=["person_id"], frame2=memberships
)

# Rename and drop fields from orgs
```

```
# to prevent field name collisions with persons_memberships
orgs = (
    orgs.drop_fields(["other_names", "identifiers"])
    .rename_field("id", "org_id")
    .rename_field("name", "org_name")
)

# Create final join of all three DynamicFrames
legislators_combined = orgs.join(
    paths1=["org_id"], paths2=["organization_id"], frame2=persons_memberships
).drop_fields(["person_id", "org_id"])

# Inspect the schema for the joined data
print("Schema for the new legislators_combined DynamicFrame:")
legislators_combined.printSchema()
```

出力

```
Schema for the persons DynamicFrame:
root
|-- family_name: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- gender: string
|-- image: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- sort_name: string
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- given_name: string
|-- birth_date: string
```

```
|-- id: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- death_date: string
```

Schema for the memberships DynamicFrame:

```
root
|-- area_id: string
|-- on_behalf_of_id: string
|-- organization_id: string
|-- role: string
|-- person_id: string
|-- legislative_period_id: string
|-- start_date: string
|-- end_date: string
```

Schema for the orgs DynamicFrame:

```
root
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- id: string
|-- classification: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- image: string
|-- seats: int
|-- type: string
```

Schema for the new legislators_combined DynamicFrame:

```
root
|-- role: string
|-- seats: int
```

```
|-- org_name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- type: string
|-- sort_name: string
|-- area_id: string
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- on_behalf_of_id: string
|-- other_names: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- name: string
|   |   |-- lang: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- name: string
|-- birth_date: string
|-- organization_id: string
|-- gender: string
|-- classification: string
|-- legislative_period_id: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- image: string
|-- given_name: string
|-- start_date: string
|-- family_name: string
|-- id: string
|-- death_date: string
|-- end_date: string
```

マップ

map(f, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)

指定したマッピング関数を元の `DynamicFrame` のすべてのレコードに適用した結果の新しい `DynamicFrame` を返します。

- `f` - `DynamicFrame` 内のすべてのレコードに適用されるマッピング関数。この関数は、`DynamicRecord` を引数として取り、新しい `DynamicRecord` を返す必要があります (必須)。

`DynamicRecord` は `DynamicFrame` 内の論理レコードを表します。これは、自己記述型であり、固定スキーマに準拠しないデータに使用できる点を除いて、Apache Spark `DataFrame` の行に似ています。

- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` - エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` - エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

例: マッピングを使用して、`DynamicFrame` のすべてのレコードに関数を適用する

この例は、`map` メソッドを使用して、関数を `DynamicFrame` のすべてのレコードに適用する方法を示しています。具体的には、この例では、複数のアドレスフィールドを単一の `struct` タイプに結合するため、`MergeAddress` という名前の関数を各レコードに適用します。

Note

この例で使用されているデータセットにアクセスするには、[「コード例: ResolveChoice、Lambda、および ApplyMapping を使用したデータ準備」](#)を参照し、[「ステップ 1: Amazon S3 バケット内のデータをクロールする」](#)の手順に従います。

```
# Example: Use map to combine fields in all records
# of a DynamicFrame

from pyspark.context import SparkContext
```

```
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Create a DynamicFrame and view its schema
medicare = glueContext.create_dynamic_frame.from_options(
    "s3",
    {"paths": ["s3://awsglue-datasets/examples/medicare/
Medicare_Hospital_Provider.csv"]},
    "csv",
    {"withHeader": True})
print("Schema for medicare DynamicFrame:")
medicare.printSchema()

# Define a function to supply to the map transform
# that merges address fields into a single field
def MergeAddress(rec):
    rec["Address"] = {}
    rec["Address"]["Street"] = rec["Provider Street Address"]
    rec["Address"]["City"] = rec["Provider City"]
    rec["Address"]["State"] = rec["Provider State"]
    rec["Address"]["Zip.Code"] = rec["Provider Zip Code"]
    rec["Address"]["Array"] = [rec["Provider Street Address"], rec["Provider City"],
rec["Provider State"], rec["Provider Zip Code"]]
    del rec["Provider Street Address"]
    del rec["Provider City"]
    del rec["Provider State"]
    del rec["Provider Zip Code"]
    return rec

# Use map to apply MergeAddress to every record
mapped_medicare = medicare.map(f = MergeAddress)
print("Schema for mapped_medicare DynamicFrame:")
mapped_medicare.printSchema()
```

出力

```
Schema for medicare DynamicFrame:
root
|-- DRG Definition: string
```

```

|-- Provider Id: string
|-- Provider Name: string
|-- Provider Street Address: string
|-- Provider City: string
|-- Provider State: string
|-- Provider Zip Code: string
|-- Hospital Referral Region Description: string
|-- Total Discharges: string
|-- Average Covered Charges: string
|-- Average Total Payments: string
|-- Average Medicare Payments: string

```

Schema for mapped_medicare DynamicFrame:

```

root
|-- Average Total Payments: string
|-- Average Covered Charges: string
|-- DRG Definition: string
|-- Average Medicare Payments: string
|-- Hospital Referral Region Description: string
|-- Address: struct
|   |-- Zip.Code: string
|   |-- City: string
|   |-- Array: array
|   |   |-- element: string
|   |-- State: string
|   |-- Street: string
|-- Provider Id: string
|-- Total Discharges: string
|-- Provider Name: string

```

mergeDynamicFrame

```

mergeDynamicFrame(stage_dynamic_frame, primary_keys, transformation_ctx =
"", options = {}, info = "", stageThreshold = 0, totalThreshold = 0)

```

レコードを識別するために、この DynamicFrame を指定されたプライマリキーに基づくステージング DynamicFrame とマージします。重複レコード (同じプライマリキーを持つレコード) は重複除外されません。ステージングフレームに一致するレコードがない場合、すべてのレコード (重複を含む) がソースから保持されます。ステージングフレームに一致するレコードがある場合、ステージングフレームのレコードによって、AWS Glue のソースのレコードが上書きされます。

- stage_dynamic_frame – マージするステージング DynamicFrame。

- `primary_keys` – ソースおよびステージング動的フレームからのレコードを照合する、プライマリキーフィールドのリスト。
- `transformation_ctx` – 現在の変換に関するメタデータを取得するために使用される一意の文字列 (オプション)。
- `options` – この変換に関する追加情報を提供する、JSON の名前と値のペアを示す文字列。この引数は現在使用されていません。
- `info` – String。この変換でのエラーに関連付けられる任意の文字列。
- `stageThreshold` – Long。指定された変換で処理がエラーアウトする必要があるエラーの数。
- `totalThreshold` – Long。この変換までに発生したエラーのうち、処理でエラーを出力する必要があるエラーの合計数。

このメソッドは、この `DynamicFrame` をステージング `DynamicFrame` とマージして取得した新しい `DynamicFrame` を返します。

以下の場合、返される `DynamicFrame` にはレコード A が含まれます。

- A がソースフレームとステージングフレームの両方に存在する場合、ステージングフレームの A が返されます。
- A がソーステーブルに存在し、`A.primaryKeys` が `stagingDynamicFrame` に存在しない場合、A はステージングテーブルで更新されていません。

ソースフレームとステージングフレームが、同じスキーマを持つ必要はありません。

例: `mergeDynamicFrame` を使用して、プライマリキーに基づいて 2 つの `DynamicFrames` をマージする

次のコード例は、`mergeDynamicFrame` メソッドを使用して、プライマリキー `id` に基づいて `DynamicFrame` を「ステージング」`DynamicFrame` とマージする方法を示しています。

データセットの例

この例では、`split_rows_collection` と呼ばれる `DynamicFrameCollection` からの 2 つの `DynamicFrames` を使用しています。次のリストに示しているのは、`split_rows_collection` のキーです。

```
dict_keys(['high', 'low'])
```

コードの例

```
# Example: Use mergeDynamicFrame to merge DynamicFrames
# based on a set of specified primary keys

from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.transforms import SelectFromCollection

# Inspect the original DynamicFrames
frame_low = SelectFromCollection.apply(dfc=split_rows_collection, key="low")
print("Inspect the DynamicFrame that contains rows where ID < 10")
frame_low.toDF().show()

frame_high = SelectFromCollection.apply(dfc=split_rows_collection, key="high")
print("Inspect the DynamicFrame that contains rows where ID > 10")
frame_high.toDF().show()

# Merge the DynamicFrames based on the "id" primary key
merged_high_low = frame_high.mergeDynamicFrame(
    stage_dynamic_frame=frame_low, primary_keys=["id"]
)

# View the results where the ID is 1 or 20
print("Inspect the merged DynamicFrame that contains the combined rows")
merged_high_low.toDF().where("id = 1 or id= 20").orderBy("id").show()
```

出力

```
Inspect the DynamicFrame that contains rows where ID < 10
+---+-----+-----+-----+-----+
| id|index|contact_details.val.type|contact_details.val.value|
+---+-----+-----+-----+-----+
| 1| 0| fax| 202-225-3307|
| 1| 1| phone| 202-225-5731|
| 2| 0| fax| 202-225-3307|
| 2| 1| phone| 202-225-5731|
| 3| 0| fax| 202-225-3307|
| 3| 1| phone| 202-225-5731|
| 4| 0| fax| 202-225-3307|
| 4| 1| phone| 202-225-5731|
| 5| 0| fax| 202-225-3307|
| 5| 1| phone| 202-225-5731|
```

```

| 6| 0| fax| 202-225-3307|
| 6| 1| phone| 202-225-5731|
| 7| 0| fax| 202-225-3307|
| 7| 1| phone| 202-225-5731|
| 8| 0| fax| 202-225-3307|
| 8| 1| phone| 202-225-5731|
| 9| 0| fax| 202-225-3307|
| 9| 1| phone| 202-225-5731|
| 10| 0| fax| 202-225-6328|
| 10| 1| phone| 202-225-4576|
+---+-----+-----+-----+-----+

```

only showing top 20 rows

Inspect the DynamicFrame that contains rows where ID > 10

```

+---+-----+-----+-----+-----+
| id|index|contact_details.val.type|contact_details.val.value|
+---+-----+-----+-----+
| 11| 0| fax| 202-225-6328|
| 11| 1| phone| 202-225-4576|
| 11| 2| twitter| RepTrentFranks|
| 12| 0| fax| 202-225-6328|
| 12| 1| phone| 202-225-4576|
| 12| 2| twitter| RepTrentFranks|
| 13| 0| fax| 202-225-6328|
| 13| 1| phone| 202-225-4576|
| 13| 2| twitter| RepTrentFranks|
| 14| 0| fax| 202-225-6328|
| 14| 1| phone| 202-225-4576|
| 14| 2| twitter| RepTrentFranks|
| 15| 0| fax| 202-225-6328|
| 15| 1| phone| 202-225-4576|
| 15| 2| twitter| RepTrentFranks|
| 16| 0| fax| 202-225-6328|
| 16| 1| phone| 202-225-4576|
| 16| 2| twitter| RepTrentFranks|
| 17| 0| fax| 202-225-6328|
| 17| 1| phone| 202-225-4576|
+---+-----+-----+-----+

```

only showing top 20 rows

Inspect the merged DynamicFrame that contains the combined rows

```

+---+-----+-----+-----+-----+
| id|index|contact_details.val.type|contact_details.val.value|
+---+-----+-----+-----+

```

```

| 1| 0|          fax|          202-225-3307|
| 1| 1|         phone|          202-225-5731|
| 20| 0|          fax|          202-225-5604|
| 20| 1|         phone|          202-225-6536|
| 20| 2|        twitter|          USRepLong|
+---+-----+-----+-----+

```

関係付け

relationalize(root_table_name, staging_path, options, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)

DynamicFrame をリレーショナルデータベースに適合する形式に変換します。DynamicFrame の関係付けは、DynamoDB などの NoSQL 環境から MySQL などのリレーショナルデータベースにデータを移動する場合に特に便利です。

この変換は、ネストされた列をネスト解除し、配列の列をピボットすることでフレームのリストを生成します。フェーズのネスト解除時に生成された結合キーを使用して、ピボットされた配列の列をルートテーブルに結合できます。

- `root_table_name` - ルートテーブルの名前。
- `staging_path` - このメソッドがピボットテーブルのパーティションを CSV 形式で保存する保存先のパス (オプション)。ピボットされたテーブルはこのパスから読み取ります。
- `options` - オプションのパラメータのディクショナリ。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - この変換のエラー報告に関連付ける文字列 (オプション)。
- `stageThreshold` - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- `totalThreshold` - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。

例: `relationalize` を使用して、**DynamicFrame** のネストされたスキーマをフラット化する

このコード例では、`relationalize` メソッドを使用して、ネストされたスキーマをリレーショナルデータベースに適合する形式にフラット化します。

データセットの例

この例では、次のスキーマを持つ `legislators_combined` と呼ばれる `DynamicFrame` を使用しています。`legislators_combined` には、`relationalize` 変換によってフラット化される、`links`、`images`、`contact_details` などのネストされた複数のフィールドがあります。

```
root
|-- role: string
|-- seats: int
|-- org_name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- type: string
|-- sort_name: string
|-- area_id: string
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- on_behalf_of_id: string
|-- other_names: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- name: string
|   |   |-- lang: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- name: string
|-- birth_date: string
|-- organization_id: string
|-- gender: string
|-- classification: string
|-- legislative_period_id: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- image: string
|-- given_name: string
|-- start_date: string
|-- family_name: string
|-- id: string
```

```
|-- death_date: string
|-- end_date: string
```

コードの例

```
# Example: Use relationalize to flatten
# a nested schema into a format that fits
# into a relational database.
# Replace DOC-EXAMPLE-S3-BUCKET/tmpDir with your own location.

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Apply relationalize and inspect new tables
legislators_relationalized = legislators_combined.relationalize(
    "l_root", "s3://DOC-EXAMPLE-BUCKET/tmpDir"
)
legislators_relationalized.keys()

# Compare the schema of the contact_details
# nested field to the new relationalized table that
# represents it
legislators_combined.select_fields("contact_details").printSchema()
legislators_relationalized.select("l_root_contact_details").toDF().where(
    "id = 10 or id = 75"
).orderBy(["id", "index"]).show()
```

出力

次の出力では、`contact_details` と呼ばれるネストされたフィールドのスキーマを、`relationalize` 変換によって作成されたテーブルと比較できます。テーブルのレコードが `id` と呼ばれる外部キーと配列の位置を表す `index` 列を使用してメインテーブルにリンクされていることに注意してください。

```
dict_keys(['l_root', 'l_root_images', 'l_root_links', 'l_root_other_names',
'l_root_contact_details', 'l_root_identifiers'])

root
```

```
 |-- contact_details: array
 |   |-- element: struct
 |     |-- type: string
 |     |-- value: string
```

```
+-----+-----+-----+-----+-----+
| id|index|contact_details.val.type|contact_details.val.value|
+-----+-----+-----+-----+
| 10|  0|          fax|          202-225-4160|
| 10|  1|         phone|          202-225-3436|
| 75|  0|          fax|          202-225-6791|
| 75|  1|         phone|          202-225-2861|
| 75|  2|       twitter|          RepSamFarr|
+-----+-----+-----+-----+-----+
```

rename_field

rename_field(oldName, newName, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)

この DynamicFrame のフィールドの名前を変更し、フィールドの名前が変更された新しい DynamicFrame を返します。

- oldName - 名前を変更するノードへのフルパス。

古い名前にドットが含まれている場合、RenameField はバックティック (`) で囲まなければ機能しません。例えば、this.old.name を thisNewName に置き換えるには、rename_field を次のように呼び出します。

```
newDyF = oldDyF.rename_field("`this.old.name`", "thisNewName")
```

- newName - 完全パスとしての新しい名前。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - この変換のエラー報告に関連付ける文字列 (オプション)。
- stageThreshold - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- totalThreshold - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。

例: `rename_field` を使用して、**DynamicFrame** のフィールドの名前を変更する

このコード例では、`rename_field` メソッドを使用して `DynamicFrame` のフィールドの名前を変更します。この例では、メソッドの連鎖を使用して、複数のフィールドの名前を同時に変更していることに注意してください。

Note

この例で使用されているデータセットにアクセスするには、「[コード例: データの結合と関係付け](#)」を参照し、「[ステップ 1: Amazon S3 バケット内のデータをクロールする](#)」の手順に従います。

コードの例

```
# Example: Use rename_field to rename fields
# in a DynamicFrame

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Inspect the original orgs schema
orgs = glueContext.create_dynamic_frame.from_catalog(
    database="legislators", table_name="organizations_json"
)
print("Original orgs schema: ")
orgs.printSchema()

# Rename fields and view the new schema
orgs = orgs.rename_field("id", "org_id").rename_field("name", "org_name")
print("New orgs schema with renamed fields: ")
orgs.printSchema()
```

出力

```
Original orgs schema:
root
 |-- identifiers: array
```

```
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- id: string
|-- classification: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- image: string
|-- seats: int
|-- type: string
```

New orgs schema with renamed fields:

```
root
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- classification: string
|-- org_id: string
|-- org_name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- image: string
|-- seats: int
|-- type: string
```

resolveChoice

```
resolveChoice(specs = None, choice = "" , database = None , table_name = None , transformation_ctx="", info="", stageThreshold=0, totalThreshold=0, catalog_id = None)
```

この DynamicFrame 内で選択タイプを解決し、新しい DynamicFrame を返します。

- specs – それぞれがタプルの形式である、解決すべき特定の曖昧要素のリスト: (field_path, action)。

resolveChoice を使用するには 2 つの方法があります。最初の方法では、specs 引数により、特定のフィールドのシーケンスと解決方法を指定します。resolveChoice のもう 1 つのモードでは、choice 引数を使用して、すべての ChoiceTypes に対して単一の解決策を指定します。

specs の値は、(field_path, action) ペアで構成されたタプルとして指定されます。field_path 値は特定のあいまいな要素を識別し、action 値は対応する解決を識別します。以下のアクションを指定できます。

- cast:*type* – すべての値について、指定した型へのキャストを試みます。例: cast:int。
- make_cols – それぞれの異なるタイプを *columnName_type* という名前の列に変換します。データをフラット化することで潜在的なあいまいさを解消します。例えば、columnA が int または string の場合、解決策は、作成された DynamicFrame に columnA_int および columnA_string という名前の 2 つの列を生成することです。
- make_struct – struct を使用してデータを表現することで、潜在的なあいまいさを解決します。例えば、列のデータが int または string となる可能性がある場合、make_struct アクションを使用すると、作成された DynamicFrame に構造体の列が生成されます。各構造体には、int と string の両方が含まれています。
- project:*type* – 可能なデータ型の 1 つにすべてのデータを投影することで、潜在的なあいまいさを解消します。例えば、列のデータが int または string の場合、project:string アクションを使用すると、すべての int 値が文字列に変換されている、作成された DynamicFrame に列が生成されます。

field_path で配列を識別する場合は、あいまいさを避けるために配列名の後に空の角括弧を置きます。例えば、使用しているデータが次のように構造化されているとします。

```
"myList": [  
  { "price": 100.00 },  
  { "price": "$100.00" }  
]
```

]

field_path を "myList[].price" に設定し、action を "cast:double" に設定すると、文字列バージョンではなく、数値バージョンの料金を選択できます。

Note

specs パラメータおよび choice パラメータのうち 1 つのみを使用できます。specs パラメータが None ではない場合、choice パラメータは空の文字列である必要があります。逆に、choice が空の文字列ではない場合、specs パラメータは None である必要があります。

- choice – すべての ChoiceTypes について、単一の解決方法を指定します。このモードは、ランタイム前に ChoiceTypes の完全なリストが不明な場合に使用できます。この引数では、上記の specs 用のアクションに加えて、以下のアクションもサポートされています。
- match_catalog – 各 ChoiceType について、指定した Data Catalog テーブル内の対応する型へのキャストを試みます。
- database – match_catalog アクションで使用する Data Catalog データベース。
- table_name – match_catalog アクションで使用する Data Catalog テーブル。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - この変換のエラー報告に関連付ける文字列 (オプション)。
- stageThreshold - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- totalThreshold – この変換までに発生したエラーのうち、プロセスでエラーを出力する必要があるエラーの数 (オプション)。デフォルトはゼロで、プロセスでエラーを出力しないことを示します。
- catalog_id – 現在アクセスされている Data Catalog のカタログ ID (Data Catalog のアカウント ID)。None を設定した場合(デフォルト値)では、呼び出し元アカウントのカタログ ID が使用されます。

例: resolveChoice を使用して、複数の型を含む列を処理する

このコード例では、resolveChoice メソッドを使用して、複数の型の値を含む DynamicFrame 列の処理方法を指定します。この例は、型の異なる列を処理する一般的な 2 つの方法を示しています。

- 列を単一のデータ型にキャストします。
- すべての型を別々の列に保持します。

データセットの例

Note

この例で使用されているデータセットにアクセスするには、「[コード例: ResolveChoice、Lambda、および ApplyMapping を使用したデータ準備](#)」を参照し、「[ステップ 1: Amazon S3 バケット内のデータをクローलする](#)」の手順に従います。

この例では、次のスキーマを持つ medicare と呼ばれる DynamicFrame を使用しています。

```
root
|-- drg definition: string
|-- provider id: choice
|   |-- long
|   |-- string
|-- provider name: string
|-- provider street address: string
|-- provider city: string
|-- provider state: string
|-- provider zip code: long
|-- hospital referral region description: string
|-- total discharges: long
|-- average covered charges: string
|-- average total payments: string
|-- average medicare payments: string
```

コードの例

```
# Example: Use resolveChoice to handle
# a column that contains multiple types

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
```

```
glueContext = GlueContext(sc)

# Load the input data and inspect the "provider id" column
medicare = glueContext.create_dynamic_frame.from_catalog(
    database="payments", table_name="medicare_hospital_provider_csv"
)
print("Inspect the provider id column:")
medicare.toDF().select("provider id").show()

# Cast provider id to type long
medicare_resolved_long = medicare.resolveChoice(specs=[("provider id", "cast:long")])
print("Schema after casting provider id to type long:")
medicare_resolved_long.printSchema()
medicare_resolved_long.toDF().select("provider id").show()

# Create separate columns
# for each provider id type
medicare_resolved_cols = medicare.resolveChoice(choice="make_cols")
print("Schema after creating separate columns for each type:")
medicare_resolved_cols.printSchema()
medicare_resolved_cols.toDF().select("provider id_long", "provider id_string").show()
```

出力

```
Inspect the 'provider id' column:
+-----+
|provider id|
+-----+
| [10001,]|
| [10005,]|
| [10006,]|
| [10011,]|
| [10016,]|
| [10023,]|
| [10029,]|
| [10033,]|
| [10039,]|
| [10040,]|
| [10046,]|
| [10055,]|
| [10056,]|
| [10078,]|
| [10083,]|
```

```
| [10085,]|  
| [10090,]|  
| [10092,]|  
| [10100,]|  
| [10103,]|
```

```
+-----+
```

only showing top 20 rows

Schema after casting 'provider id' to type long:

root

```
|-- drg definition: string  
|-- provider id: long  
|-- provider name: string  
|-- provider street address: string  
|-- provider city: string  
|-- provider state: string  
|-- provider zip code: long  
|-- hospital referral region description: string  
|-- total discharges: long  
|-- average covered charges: string  
|-- average total payments: string  
|-- average medicare payments: string
```

```
+-----+
```

```
|provider id|
```

```
+-----+
```

```
| 10001|  
| 10005|  
| 10006|  
| 10011|  
| 10016|  
| 10023|  
| 10029|  
| 10033|  
| 10039|  
| 10040|  
| 10046|  
| 10055|  
| 10056|  
| 10078|  
| 10083|  
| 10085|  
| 10090|  
| 10092|
```

```
|      10100|
|      10103|
```

```
+-----+
```

only showing top 20 rows

Schema after creating separate columns for each type:

root

```
|-- drg definition: string
|-- provider id_string: string
|-- provider id_long: long
|-- provider name: string
|-- provider street address: string
|-- provider city: string
|-- provider state: string
|-- provider zip code: long
|-- hospital referral region description: string
|-- total discharges: long
|-- average covered charges: string
|-- average total payments: string
|-- average medicare payments: string
```

```
+-----+-----+
```

```
|provider id_long|provider id_string|
```

```
+-----+-----+
```

```
|      10001|      null|
|      10005|      null|
|      10006|      null|
|      10011|      null|
|      10016|      null|
|      10023|      null|
|      10029|      null|
|      10033|      null|
|      10039|      null|
|      10040|      null|
|      10046|      null|
|      10055|      null|
|      10056|      null|
|      10078|      null|
|      10083|      null|
|      10085|      null|
|      10090|      null|
|      10092|      null|
|      10100|      null|
|      10103|      null|
```

```
+-----+-----+
only showing top 20 rows
```

select_fields

```
select_fields(paths, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)
```

選択したフィールドを含む新しい `DynamicFrame` を返します。

- `paths` - 文字列のリスト。各文字列は、選択する最上位ノードへのパスです。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - この変換のエラー報告に関連付ける文字列 (オプション)。
- `stageThreshold` - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- `totalThreshold` - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。

例: `select_fields` を使用して、選択したフィールドで新しい `DynamicFrame` を作成する

次のコード例は、`select_fields` メソッドを使用して、既存の `DynamicFrame` から選択されたフィールドのリストを使用し、新しい `DynamicFrame` を作成する方法を示しています。

Note

この例で使用されているデータセットにアクセスするには、「[コード例: データの結合と関係付け](#)」を参照し、「[ステップ 1: Amazon S3 バケット内のデータをクロールする](#)」の手順に従います。

```
# Example: Use select_fields to select specific fields from a DynamicFrame

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
```

```
glueContext = GlueContext(sc)

# Create a DynamicFrame and view its schema
persons = glueContext.create_dynamic_frame.from_catalog(
    database="legislators", table_name="persons_json"
)
print("Schema for the persons DynamicFrame:")
persons.printSchema()

# Create a new DynamicFrame with chosen fields
names = persons.select_fields(paths=["family_name", "given_name"])
print("Schema for the names DynamicFrame, created with select_fields:")
names.printSchema()
names.toDF().show()
```

出力

```
Schema for the persons DynamicFrame:
```

```
root
|-- family_name: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- gender: string
|-- image: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- sort_name: string
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
```

```
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- death_date: string
```

Schema for the names DynamicFrame:

```
root
```

```
|-- family_name: string
|-- given_name: string
```

```
+-----+-----+
|family_name|given_name|
+-----+-----+
|   Collins|  Michael|
|  Huizenga|    Bill|
|   Clawson|  Curtis|
|   Solomon|  Gerald|
|   Rigell|   Edward|
|   Crapo|   Michael|
|   Hutto|    Earl|
|   Ertel|   Allen|
|   Minish|   Joseph|
|  Andrews|   Robert|
|   Walden|    Greg|
|   Kazen|   Abraham|
|   Turner|   Michael|
|   Kolbe|    James|
| Lowenthal|    Alan|
|   Capuano|  Michael|
|   Schrader|   Kurt|
|   Nadler|   Jerrold|
|   Graves|    Tom|
|  McMillan|    John|
+-----+-----+
only showing top 20 rows
```

`simplify_ddb_json`

`simplify_ddb_json(): DynamicFrame`

特に DynamoDB JSON 構造にある `DynamicFrame` でネスト化された列を単純化し、新しい単純化された `DynamicFrame` を返します。リストタイプに複数タイプまたは Map タイプがある場合、リ

ストの要素は単純化されません。これは通常の `unnest` 変換とは異なって動作する変換の特定タイプであり、データが既に DynamoDB JSON 構造に格納されている必要があることに注意してください。詳細については、「[DynamoDB JSON](#)」を参照してください。

例えば、DynamoDB JSON 構造体のあるエクスポートを読み取るスキーマは次のようになります。

```
root
|-- Item: struct
|   |-- parentMap: struct
|   |   |-- M: struct
|   |   |   |-- childMap: struct
|   |   |   |   |-- M: struct
|   |   |   |   |   |-- appName: struct
|   |   |   |   |   |   |-- S: string
|   |   |   |   |   |   |-- packageName: struct
|   |   |   |   |   |   |   |-- S: string
|   |   |   |   |   |   |   |-- updatedAt: struct
|   |   |   |   |   |   |   |   |-- N: string
|   |   |-- strings: struct
|   |   |   |-- SS: array
|   |   |   |   |-- element: string
|   |   |-- numbers: struct
|   |   |   |-- NS: array
|   |   |   |   |-- element: string
|   |   |-- binaries: struct
|   |   |   |-- BS: array
|   |   |   |   |-- element: string
|   |   |-- isDDBJson: struct
|   |   |   |-- BOOL: boolean
|   |   |-- nullValue: struct
|   |   |   |-- NULL: boolean
```

`simplify_ddb_json()` 変換は、以下のように変換します。

```
root
|-- parentMap: struct
|   |-- childMap: struct
|   |   |-- appName: string
|   |   |-- packageName: string
|   |   |-- updatedAt: string
|-- strings: array
|   |-- element: string
|-- numbers: array
```

```
|    |-- element: string
|-- binaries: array
|    |-- element: string
|-- isDDBJson: boolean
|-- nullValue: null
```

例: `simplify_ddb_json` を使用して DynamoDB JSON の単純化を呼び出す

このコード例では `simplify_ddb_json` メソッドを使用し、AWS Glue DynamoDB エクスポートコネクタを使用し、DynamoDB JSON 単純化を呼び出してパーティションの数を表示します。

コードの例

```
from pyspark.context import SparkContext
from awsglue.context import GlueContext

sc = SparkContext()
glueContext = GlueContext(sc)

dynamicFrame = glueContext.create_dynamic_frame.from_options(
    connection_type = "dynamodb",
    connection_options = {
        'dynamodb.export': 'ddb',
        'dynamodb.tableArn': '<table arn>',
        'dynamodb.s3.bucket': '<bucket name>',
        'dynamodb.s3.prefix': '<bucket prefix>',
        'dynamodb.s3.bucketOwner': '<account_id of bucket>'
    }
)
simplified = dynamicFrame.simplify_ddb_json()
print(simplified.getNumPartitions())
```

スピゴット

`spigot(path, options={})`

ジョブで実行された変換が確認しやすくなるように、サンプルレコードを指定した送信先に書き込みます。

- `path` – 書き込み先へのパス (必須)。

- `options` – オプションを指定するキーと値のペア (オプション)。"topk" オプションは、最初の k レコードを書き込むことを指定します。"prob" オプションは、指定されたレコードを選択する確率(10 進数)を指定します。これを使用して、書き込むレコードを選択できます。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。

例: `spigot` を使用して `DynamicFrame` のサンプルフィールドを Amazon S3 に書き込む

このコード例では、`spigot` メソッドを使用して、`select_fields` 変換を適用した後に Amazon S3 バケットにサンプルレコードを書き込みます。

データセットの例

Note

この例で使用されているデータセットにアクセスするには、[「コード例: データの結合と関係付け」](#)を参照し、[「ステップ 1: Amazon S3 バケット内のデータをクロールする」](#)の手順に従います。

この例では、次のスキーマを持つ `persons` と呼ばれる `DynamicFrame` を使用しています。

```
root
|-- family_name: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- gender: string
|-- image: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- sort_name: string
```

```
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- death_date: string
```

コードの例

```
# Example: Use spigot to write sample records
# to a destination during a transformation
# from pyspark.context import SparkContext.
# Replace DOC-EXAMPLE-BUCKET with your own location.

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Load table data into a DynamicFrame
persons = glueContext.create_dynamic_frame.from_catalog(
    database="legislators", table_name="persons_json"
)

# Perform the select_fields on the DynamicFrame
persons = persons.select_fields(paths=["family_name", "given_name", "birth_date"])

# Use spigot to write a sample of the transformed data
# (the first 10 records)
spigot_output = persons.spigot(
    path="s3://DOC-EXAMPLE-BUCKET", options={"topk": 10}
)
```

出力

spigot で Amazon S3 に書き込むデータの例を次に示します。サンプルコードで `options={"topk": 10}` を指定したため、サンプルデータには最初の 10 個のレコードが含まれません。

```
{"family_name":"Collins","given_name":"Michael","birth_date":"1944-10-15"}
{"family_name":"Huizenga","given_name":"Bill","birth_date":"1969-01-31"}
{"family_name":"Clawson","given_name":"Curtis","birth_date":"1959-09-28"}
{"family_name":"Solomon","given_name":"Gerald","birth_date":"1930-08-14"}
{"family_name":"Rigell","given_name":"Edward","birth_date":"1960-05-28"}
{"family_name":"Crapo","given_name":"Michael","birth_date":"1951-05-20"}
{"family_name":"Hutto","given_name":"Earl","birth_date":"1926-05-12"}
{"family_name":"Ertel","given_name":"Allen","birth_date":"1937-11-07"}
{"family_name":"Minish","given_name":"Joseph","birth_date":"1916-09-01"}
{"family_name":"Andrews","given_name":"Robert","birth_date":"1957-08-04"}
```

split_fields

split_fields(paths, name1, name2, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)

2 つの DynamicFrames を含む新しい DynamicFrameCollection を返します。1 つ目の DynamicFrame には分割されたすべてのノードが含まれ、2 つ目には残りのノードが含まれていません。

- paths - 文字列のリスト。各文字列は新しい DynamicFrame に分割するノードのフルパスです。
- name1 - 分割された DynamicFrame の名前文字列。
- name2 - 指定されたノードが分割された後に残る DynamicFrame の名前文字列。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - この変換のエラー報告に関連付ける文字列 (オプション)。
- stageThreshold - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- totalThreshold - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。

例: `split_fields` を使用して、選択したフィールドを別の **DynamicFrame** に分割する

このコード例では、`split_fields` メソッドを使用して、指定したフィールドのリストを別の **DynamicFrame** に分割します。

データセットの例

この例では、`legislators_relationalized` という名前のコレクションからの `l_root_contact_details` と呼ばれる **DynamicFrame** が使用されています。

`l_root_contact_details` には、次のスキーマとエントリが含まれています。

```
root
|-- id: long
|-- index: int
|-- contact_details.val.type: string
|-- contact_details.val.value: string
```

id	index	contact_details.val.type	contact_details.val.value
1	0	phone	202-225-5265
1	1	twitter	kathyhochul
2	0	phone	202-225-3252
2	1	twitter	repjackyrosen
3	0	fax	202-225-1314
3	1	phone	202-225-3772
...			

コードの例

```
# Example: Use split_fields to split selected
# fields into a separate DynamicFrame

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Load the input DynamicFrame and inspect its schema
frame_to_split = legislators_relationalized.select("l_root_contact_details")
```

```
print("Inspect the input DynamicFrame schema:")
frame_to_split.printSchema()

# Split id and index fields into a separate DynamicFrame
split_fields_collection = frame_to_split.split_fields(["id", "index"], "left", "right")

# Inspect the resulting DynamicFrames
print("Inspect the schemas of the DynamicFrames created with split_fields:")
split_fields_collection.select("left").printSchema()
split_fields_collection.select("right").printSchema()
```

出力

```
Inspect the input DynamicFrame's schema:
root
|-- id: long
|-- index: int
|-- contact_details.val.type: string
|-- contact_details.val.value: string

Inspect the schemas of the DynamicFrames created with split_fields:
root
|-- id: long
|-- index: int

root
|-- contact_details.val.type: string
|-- contact_details.val.value: string
```

split_rows

split_rows(comparison_dict, name1, name2, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)

DynamicFrame の 1 つ以上の行を、新しい DynamicFrame に分割します。

このメソッドは、2 つの DynamicFrames を含む新しい DynamicFrameCollection を返します。1 つ目の DynamicFrame には分割されたすべての行が含まれ、2 つ目には残りの行が含まれています。

- **comparison_dict** – 列へのパスを示すためのキーと、列の値の比較対象である値に対しコンパレータをマッピングする別のディクショナリを示すための値を持つ、ディクショナリ。例え

ば、{"age": {">": 10, "<": 20}} は、age 列の値が 10 より大きく 20 より小さいすべての行を分割します。

- name1 - 分割された DynamicFrame の名前文字列。
- name2 - 指定されたノードが分割された後に残る DynamicFrame の名前文字列。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - この変換のエラー報告に関連付ける文字列 (オプション)。
- stageThreshold - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- totalThreshold - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。

例: split_rows を使用して、**DynamicFrame** 内の行を分割する

このコード例では、split_rows メソッドを使用して、id フィールド値に基づいて DynamicFrame の行を分割します。

データセットの例

この例では、legislators_relationalized という名前のコレクションから選択された l_root_contact_details と呼ばれる DynamicFrame が使用されています。

l_root_contact_details には、次のスキーマとエントリが含まれています。

```
root
|-- id: long
|-- index: int
|-- contact_details.val.type: string
|-- contact_details.val.value: string
```

id	index	contact_details.val.type	contact_details.val.value
1	0	phone	202-225-5265
1	1	twitter	kathyhochul
2	0	phone	202-225-3252
2	1	twitter	repjackyrosen
3	0	fax	202-225-1314
3	1	phone	202-225-3772
3	2	twitter	MikeRossUpdates
4	0	fax	202-225-1314

```

| 4| 1| phone| 202-225-3772|
| 4| 2| twitter| MikeRossUpdates|
| 5| 0| fax| 202-225-1314|
| 5| 1| phone| 202-225-3772|
| 5| 2| twitter| MikeRossUpdates|
| 6| 0| fax| 202-225-1314|
| 6| 1| phone| 202-225-3772|
| 6| 2| twitter| MikeRossUpdates|
| 7| 0| fax| 202-225-1314|
| 7| 1| phone| 202-225-3772|
| 7| 2| twitter| MikeRossUpdates|
| 8| 0| fax| 202-225-1314|
+---+-----+-----+-----+

```

コードの例

```

# Example: Use split_rows to split up
# rows in a DynamicFrame based on value

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Retrieve the DynamicFrame to split
frame_to_split = legislators_relationalized.select("l_root_contact_details")

# Split up rows by ID
split_rows_collection = frame_to_split.split_rows({"id": {">": 10}}, "high", "low")

# Inspect the resulting DynamicFrames
print("Inspect the DynamicFrame that contains IDs < 10")
split_rows_collection.select("low").toDF().show()
print("Inspect the DynamicFrame that contains IDs > 10")
split_rows_collection.select("high").toDF().show()

```

出力

```

Inspect the DynamicFrame that contains IDs < 10
+---+-----+-----+-----+
| id|index|contact_details.val.type|contact_details.val.value|

```

```

+---+-----+-----+-----+
| 1|    0|           phone|           202-225-5265|
| 1|    1|           twitter|           kathyhochul|
| 2|    0|           phone|           202-225-3252|
| 2|    1|           twitter|           repjackyrosen|
| 3|    0|            fax|           202-225-1314|
| 3|    1|           phone|           202-225-3772|
| 3|    2|           twitter|           MikeRossUpdates|
| 4|    0|            fax|           202-225-1314|
| 4|    1|           phone|           202-225-3772|
| 4|    2|           twitter|           MikeRossUpdates|
| 5|    0|            fax|           202-225-1314|
| 5|    1|           phone|           202-225-3772|
| 5|    2|           twitter|           MikeRossUpdates|
| 6|    0|            fax|           202-225-1314|
| 6|    1|           phone|           202-225-3772|
| 6|    2|           twitter|           MikeRossUpdates|
| 7|    0|            fax|           202-225-1314|
| 7|    1|           phone|           202-225-3772|
| 7|    2|           twitter|           MikeRossUpdates|
| 8|    0|            fax|           202-225-1314|
+---+-----+-----+-----+

```

only showing top 20 rows

Inspect the DynamicFrame that contains IDs > 10

```

+---+-----+-----+-----+
| id|index|contact_details.val.type|contact_details.val.value|
+---+-----+-----+-----+
| 11|    0|           phone|           202-225-5476|
| 11|    1|           twitter|           RepDavidYoung|
| 12|    0|           phone|           202-225-4035|
| 12|    1|           twitter|           RepStephMurphy|
| 13|    0|            fax|           202-226-0774|
| 13|    1|           phone|           202-225-6335|
| 14|    0|            fax|           202-226-0774|
| 14|    1|           phone|           202-225-6335|
| 15|    0|            fax|           202-226-0774|
| 15|    1|           phone|           202-225-6335|
| 16|    0|            fax|           202-226-0774|
| 16|    1|           phone|           202-225-6335|
| 17|    0|            fax|           202-226-0774|
| 17|    1|           phone|           202-225-6335|
| 18|    0|            fax|           202-226-0774|
| 18|    1|           phone|           202-225-6335|

```

```

| 19|    0|                fax|                202-226-0774|
| 19|    1|                phone|               202-225-6335|
| 20|    0|                fax|                202-226-0774|
| 20|    1|                phone|               202-225-6335|
+---+-----+-----+-----+-----+
only showing top 20 rows

```

アンボックス

unbox(path, format, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0, **options)

DynamicFrame の文字列フィールドをアンボックス (再フォーマット) し、アンボックスされた DynamicRecords を含む新しい DynamicFrame を返します。

DynamicRecord は DynamicFrame 内の論理レコードを表します。これは、自己記述型であり、固定スキーマに準拠しないデータに使用できる点を除いて、Apache Spark DataFrame の行に似ています。

- path - アンボックスする文字列ノードへのフルパス。
- format - 形式の仕様 (オプション)。Amazon S3 や、複数の形式をサポートする AWS Glue 接続の場合に使用します。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - この変換のエラー報告に関連付ける文字列 (オプション)。
- stageThreshold - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- totalThreshold - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- options - 次の 1 つ以上。
 - separator - 区切り文字を含む文字列。
 - escaper - エスケープ文字を含む文字列。
 - skipFirst - 最初のインスタンスをスキップするかどうかを示すブール値。
 - withSchema - ノードのスキーマの JSON 表現を含む文字列。スキーマの JSON 表現の形式は、StructType.json() の出力によって定義されます。
 - withHeader - ヘッダーが含まれているかどうかを示すブール値。

例: unbox を使用して、文字列フィールドを構造体にアンボックスする

このコード例では、unbox メソッドを使用して、DynamicFrame の文字列フィールドを struct 型のフィールドにアンボックスまたは再フォーマットします。

データセットの例

この例では、次のスキーマとエントリを持つ mapped_with_string と呼ばれる DynamicFrame を使用しています。

AddressString という名前のフィールドに注目してください。この例では、このフィールドを構造体にアンボックスしています。

```

root
|-- Average Total Payments: string
|-- AddressString: string
|-- Average Covered Charges: string
|-- DRG Definition: string
|-- Average Medicare Payments: string
|-- Hospital Referral Region Description: string
|-- Address: struct
|   |-- Zip.Code: string
|   |-- City: string
|   |-- Array: array
|   |   |-- element: string
|   |-- State: string
|   |-- Street: string
|-- Provider Id: string
|-- Total Discharges: string
|-- Provider Name: string

+-----+-----+-----+-----+
+-----+-----+-----+-----+
+-----+-----+-----+-----+
|Average Total Payments|      AddressString|Average Covered Charges|      DRG
  Definition|Average Medicare Payments|Hospital Referral Region Description|
  Address|Provider Id|Total Discharges|      Provider Name|
+-----+-----+-----+-----+
+-----+-----+-----+-----+
+-----+-----+-----+-----+
|                $5777.24|{"Street": "1108 ...|                $32963.07|039 -
EXTRACRANIA...|                $4763.73|                AL - Dothan|[36301,
DOTHAN, [...]|                10001|                91|SOUTHEAST ALABAMA...|

```

```

|           $5787.57|{"Street": "2505 ...|           $15131.85|039 -
EXTRACRANIA...|           $4976.71|           AL - Birmingham|[35957,
BOAZ, [25...|           10005|           14|MARSHALL MEDICAL ...|
|           $5434.95|{"Street": "205 M...|           $37560.37|039 -
EXTRACRANIA...|           $4453.79|           AL - Birmingham|[35631,
FLORENCE,...|           10006|           24|ELIZA COFFEE MEMO...|
|           $5417.56|{"Street": "50 ME...|           $13998.28|039 -
EXTRACRANIA...|           $4129.16|           AL - Birmingham|[35235,
BIRMINGHA...|           10011|           25| ST VINCENT'S EAST|
...

```

コードの例

```

# Example: Use unbox to unbox a string field
# into a struct in a DynamicFrame

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

unboxed = mapped_with_string.unbox("AddressString", "json")
unboxed.printSchema()
unboxed.toDF().show()

```

出力

```

root
|-- Average Total Payments: string
|-- AddressString: struct
|   |-- Street: string
|   |-- City: string
|   |-- State: string
|   |-- Zip.Code: string
|   |-- Array: array
|   |   |-- element: string
|-- Average Covered Charges: string
|-- DRG Definition: string
|-- Average Medicare Payments: string
|-- Hospital Referral Region Description: string

```

```

|-- Address: struct
|   |-- Zip.Code: string
|   |-- City: string
|   |-- Array: array
|   |   |-- element: string
|   |-- State: string
|   |-- Street: string
|-- Provider Id: string
|-- Total Discharges: string
|-- Provider Name: string
    
```

```

+-----+-----+-----+
+-----+-----+-----+
+-----+-----+-----+
|Average Total Payments|      AddressString|Average Covered Charges|      DRG
|Definition|Average Medicare Payments|Hospital Referral Region Description|
|Address|Provider Id|Total Discharges|      Provider Name|
+-----+-----+-----+
+-----+-----+-----+
+-----+-----+-----+
|          $5777.24|[1108 ROSS CLARK ...|          $32963.07|039 -
EXTRACRANIA...|          $4763.73|          AL - Dothan|[36301,
DOTHAN, [...] 10001|          91|SOUTHEAST ALABAMA...|
|          $5787.57|[2505 U S HIGHWAY...|          $15131.85|039 -
EXTRACRANIA...|          $4976.71|          AL - Birmingham|[35957,
BOAZ, [25...| 10005|          14|MARSHALL MEDICAL ...|
|          $5434.95|[205 MARENGO STRE...|          $37560.37|039 -
EXTRACRANIA...|          $4453.79|          AL - Birmingham|[35631,
FLORENCE,...| 10006|          24|ELIZA COFFEE MEMO...|
|          $5417.56|[50 MEDICAL PARK ...|          $13998.28|039 -
EXTRACRANIA...|          $4129.16|          AL - Birmingham|[35235,
BIRMINGHA...| 10011|          25| ST VINCENT'S EAST|
|          $5658.33|[1000 FIRST STREE...|          $31633.27|039 -
EXTRACRANIA...|          $4851.44|          AL - Birmingham|[35007,
ALABASTER...| 10016|          18|SHELBY BAPTIST ME...|
|          $6653.80|[2105 EAST SOUTH ...|          $16920.79|039 -
EXTRACRANIA...|          $5374.14|          AL - Montgomery|[36116,
MONTGOMER...| 10023|          67|BAPTIST MEDICAL C...|
|          $5834.74|[2000 PEPPERELL P...|          $11977.13|039 -
EXTRACRANIA...|          $4761.41|          AL - Birmingham|[36801,
OPELIKA, ...| 10029|          51|EAST ALABAMA MEDI...|
|          $8031.12|[619 SOUTH 19TH S...|          $35841.09|039 -
EXTRACRANIA...|          $5858.50|          AL - Birmingham|[35233,
BIRMINGHA...| 10033|          32|UNIVERSITY OF ALA...|
    
```

```

|          $6113.38|[101 SIVLEY RD, H...|          $28523.39|039 -
EXTRACRANIA...|          $5228.40|          AL - Huntsville|[35801,
HUNTSVILL...|          10039|          135| HUNTSVILLE HOSPITAL|
|          $5541.05|[1007 GOODYEAR AV...|          $75233.38|039 -
EXTRACRANIA...|          $4386.94|          AL - Birmingham|[35903,
GADSDEN, ...|          10040|          34|GADSDEN REGIONAL ...|
|          $5461.57|[600 SOUTH THIRD ...|          $67327.92|039 -
EXTRACRANIA...|          $4493.57|          AL - Birmingham|[35901,
GADSDEN, ...|          10046|          14|RIVERVIEW REGIONA...|
|          $5356.28|[4370 WEST MAIN S...|          $39607.28|039 -
EXTRACRANIA...|          $4408.20|          AL - Dothan|[36305,
DOTHAN, [...|          10055|          45| FLOWERS HOSPITAL|
|          $5374.65|[810 ST VINCENT'S...|          $22862.23|039 -
EXTRACRANIA...|          $4186.02|          AL - Birmingham|[35205,
BIRMINGHA...|          10056|          43|ST VINCENT'S BIRM...|
|          $5366.23|[400 EAST 10TH ST...|          $31110.85|039 -
EXTRACRANIA...|          $4376.23|          AL - Birmingham|[36207,
ANNISTON,...|          10078|          21|NORTHEAST ALABAMA...|
|          $5282.93|[1613 NORTH MCKEN...|          $25411.33|039 -
EXTRACRANIA...|          $4383.73|          AL - Mobile|[36535,
FOLEY, [1...|          10083|          15|SOUTH BALDWIN REG...|
|          $5676.55|[1201 7TH STREET ...|          $9234.51|039 -
EXTRACRANIA...|          $4509.11|          AL - Huntsville|[35609,
DECATUR, ...|          10085|          27|DECATUR GENERAL H...|
|          $5930.11|[6801 AIRPORT BOU...|          $15895.85|039 -
EXTRACRANIA...|          $3972.85|          AL - Mobile|[36608,
MOBILE, [...|          10090|          27| PROVIDENCE HOSPITAL|
|          $6192.54|[809 UNIVERSITY B...|          $19721.16|039 -
EXTRACRANIA...|          $5179.38|          AL - Tuscaloosa|[35401,
TUSCALOOS...|          10092|          31|D C H REGIONAL ME...|
|          $4968.00|[750 MORPHY AVENU...|          $10710.88|039 -
EXTRACRANIA...|          $3898.88|          AL - Mobile|[36532,
FAIRHOPE,...|          10100|          18| THOMAS HOSPITAL|
|          $5996.00|[701 PRINCETON AV...|          $51343.75|039 -
EXTRACRANIA...|          $4962.45|          AL - Birmingham|[35211,
BIRMINGHA...|          10103|          33|BAPTIST MEDICAL C...|

```

```

+-----+-----+-----+
+-----+-----+-----+
+-----+-----+-----+

```

only showing top 20 rows

union

```
union(frame1, frame2, transformation_ctx = "", info = "", stageThreshold = 0, totalThreshold = 0)
```

2 つの DynamicFrames を結合します。両方の入力 DynamicFrame からのすべてのレコードを含む DynamicFrames を返します。この変換は、2 つの DataFrames rame を同等のデータで結合した結果とは異なる結果を返す場合があります。Spark DataFrame ユニオンの動作が必要な場合は、toDF の使用を検討してください。

- frame1 — 最初に結合する DynamicFrame。
- frame2 — 2 番目に結合する DynamicFrame。
- transformation_ctx - (オプション) 統計/ステータス情報を識別するために使用される一意の文字列。
- info - (オプション) 変換のエラーに関連付けられる文字列。
- stageThreshold — (オプション) 処理がエラーになるまでの変換中の最大エラー数。
- totalThreshold — (オプション) 処理がエラーになるまでの変換中の最大エラー数。

ネスト解除

```
unnest(transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)
```

DynamicFrame 内のネストされたオブジェクトをネスト解除して、最上位レベルのオブジェクトにし、新しくネスト解除された DynamicFrame を返します。

- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - この変換のエラー報告に関連付ける文字列 (オプション)。
- stageThreshold - この変換中に発生した、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。
- totalThreshold - この変換までに発生したエラーのうち、プロセスがエラーを出力するエラーの数 (オプション)。デフォルトはゼロで、エラーを出力しないことを示します。

例: unnest を使用して、ネストされたフィールドを最上位フィールドに変換する

このコード例では、unnest メソッドを使用して、DynamicFrame のネストされたすべてのフィールドを最上位フィールドにフラット化します。

データセットの例

この例では、次のスキーマを持つ `mapped_medicare` と呼ばれる `DynamicFrame` を使用しています。Address フィールドは、ネストされたデータを含む唯一のフィールドであることに注意してください。

```
root
|-- Average Total Payments: string
|-- Average Covered Charges: string
|-- DRG Definition: string
|-- Average Medicare Payments: string
|-- Hospital Referral Region Description: string
|-- Address: struct
|   |-- Zip.Code: string
|   |-- City: string
|   |-- Array: array
|   |   |-- element: string
|   |-- State: string
|   |-- Street: string
|-- Provider Id: string
|-- Total Discharges: string
|-- Provider Name: string
```

コードの例

```
# Example: Use unnest to unnest nested
# objects in a DynamicFrame

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Unnest all nested fields
unnested = mapped_medicare.unnest()
unnested.printSchema()
```

出力

```
root
|-- Average Total Payments: string
|-- Average Covered Charges: string
```

```

|-- DRG Definition: string
|-- Average Medicare Payments: string
|-- Hospital Referral Region Description: string
|-- Address.Zip.Code: string
|-- Address.City: string
|-- Address.Array: array
|   |-- element: string
|-- Address.State: string
|-- Address.Street: string
|-- Provider Id: string
|-- Total Discharges: string
|-- Provider Name: string

```

unnest_ddb_json

特に DynamoDB JSON 構造体にある DynamicFrame でネストされた列をネスト解除すると、新しくネスト解除された DynamicFrame が返されます。構造体型の配列のある列は、ネスト解除されません。これは、通常の unnest 変換とは異なる特殊なタイプのネスト解除変換で、データが DynamoDB JSON 構造体に格納されている必要があることに注意してください。詳細については、「[DynamoDB JSON](#)」を参照してください。

unnest_ddb_json(transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)

- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - この変換のエラー報告に関連付ける文字列 (オプション)。
- stageThreshold - この変換中にプロセスで発生するエラーの数 (オプション: デフォルトではゼロ、プロセスがエラーを出力しないことを示します)。
- totalThreshold - この変換までに発生したエラーのうち、プロセスでエラーとなるエラーの数 (オプション: デフォルトではゼロ、プロセスがエラーを出力しないことを示します)。

例えば、DynamoDB JSON 構造体のあるエクスポートを読み取るスキーマは次のようになります。

```

root
|-- Item: struct
|   |-- ColA: struct
|   |   |-- S: string
|   |-- ColB: struct

```

```
| | |-- S: string
| |-- ColC: struct
| | |-- N: string
| |-- ColD: struct
| | |-- L: array
| | | |-- element: null
```

`unnest_ddb_json()` 変換は、以下のように変換します。

```
root
|-- ColA: string
|-- ColB: string
|-- ColC: string
|-- ColD: array
| |-- element: null
```

次のコード例は、AWS Glue DynamoDB エクスポートコネクタを使用し、DynamoDB JSON `unnest` を呼び出し、パーティションの数を表示します。

```
import sys
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from awsglue.utils import getResolvedOptions

args = getResolvedOptions(sys.argv, ["JOB_NAME"])
glue_context= GlueContext(SparkContext.getOrCreate())
job = Job(glue_context)
job.init(args["JOB_NAME"], args)

dynamicFrame = glue_context.create_dynamic_frame.from_options(
    connection_type="dynamodb",
    connection_options={
        "dynamodb.export": "ddb",
        "dynamodb.tableArn": "<test_source>",
        "dynamodb.s3.bucket": "<bucket name>",
        "dynamodb.s3.prefix": "<bucket prefix>",
        "dynamodb.s3.bucketOwner": "<account_id>",
    }
)
unnested = dynamicFrame.unnest_ddb_json()
print(unnested.getNumPartitions())
```

```
job.commit()
```

書き込み

write(connection_type, connection_options, format, format_options, accumulator_size)

この DynamicFrame の [GlueContext クラス](#) から指定された接続タイプの [DataSink \(オブジェクト\)](#) を取得し、この DynamicFrame のコンテンツの書式設定および書き込みに使用します。指定されたとおりに書式設定され、書き込まれる新しい DynamicFrame を返します。

- connection_type - 使用する接続タイプ。有効な値には、s3、mysql、postgresql、redshift、sqlserver、および oracle があります。
- connection_options - 使用する接続オプション (オプション)。s3 の connection_type では、Amazon S3 パスが定義されています。

```
connection_options = {"path": "s3://aws-glue-target/temp"}
```

JDBC 接続の場合、いくつかのプロパティを定義する必要があります。データベース名は URL の一部である必要があることに注意してください。オプションで接続オプションに含めることができます。

Warning

スクリプトにパスワードを保存することはお勧めしません。AWS Secrets Manager または AWS Glue データカタログから取得する場合には、boto3 を使用することを検討してください。

```
connection_options = {"url": "jdbc-url/database", "user": "username",  
"password": passwordVariable, "dbtable": "table-name", "redshiftTmpDir": "s3-tempdir-path"}
```

- format - 形式の仕様 (オプション)。これは、Amazon Simple Storage Service (Amazon S3)、または複数の形式をサポートする AWS Glue 接続で使用されます。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- format_options - 指定された形式についてのオプション。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。

- `accumulator_size` - 使用するバイト単位の累積サイズ (オプション)。

— errors —

- [assertErrorThreshold](#)
- [errorsAsDynamicFrame](#)
- [errorsCount](#)
- [stageErrorsCount](#)

assertErrorThreshold

`assertErrorThreshold()` - この `DynamicFrame` を作成した変換エラーに対するアサーション。基盤になる `DataFrame` から `Exception` を返します。

errorsAsDynamicFrame

`errorsAsDynamicFrame()` - 内部にネストされたエラーレコードを持つ `DynamicFrame` を返します。

例: `errorsAsDynamicFrame` を使用してエラーレコードを表示する

次のコード例は、`errorsAsDynamicFrame` メソッドを使用して `DynamicFrame` のエラーレコードを表示する方法を示しています。

データセットの例

この例では、JSON として Amazon S3 にアップロードできる次のデータセットを使用します。2 つ目のレコードの形式に誤りがあることに注意してください。通常、SparkSQL を使用すると、不正な形式のデータによってファイルの解析が中断されます。ただし、`DynamicFrame` は、不正な形式の問題を認識し、不正な行を個別に処理できるエラーレコードに変換します。

```
{"id": 1, "name": "george", "surname": "washington", "height": 178}
{"id": 2, "name": "benjamin", "surname": "franklin",
{"id": 3, "name": "alexander", "surname": "hamilton", "height": 171}
{"id": 4, "name": "john", "surname": "jay", "height": 190}
```

コードの例

```
# Example: Use errorsAsDynamicFrame to view error records.
```

```
# Replace s3://DOC-EXAMPLE-S3-BUCKET/error_data.json with your location.

from pyspark.context import SparkContext
from awsglue.context import GlueContext

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Create errors DynamicFrame, view schema
errors = glueContext.create_dynamic_frame.from_options(
    "s3", {"paths": ["s3://DOC-EXAMPLE-S3-BUCKET/error_data.json"]}, "json"
)
print("Schema of errors DynamicFrame:")
errors.printSchema()

# Show that errors only contains valid entries from the dataset
print("errors contains only valid records from the input dataset (2 of 4 records)")
errors.toDF().show()

# View errors
print("Errors count:", str(errors.errorsCount()))
print("Errors:")
errors.errorsAsDynamicFrame().toDF().show()

# View error fields and error data
error_record = errors.errorsAsDynamicFrame().toDF().head()

error_fields = error_record["error"]
print("Error fields: ")
print(error_fields.asDict().keys())

print("\nError record data:")
for key in error_fields.asDict().keys():
    print("\n", key, ": ", str(error_fields[key]))
```

出力

```
Schema of errors DynamicFrame:
root
 |-- id: int
 |-- name: string
 |-- surname: string
```

```
|-- height: int
```

```
errors contains only valid records from the input dataset (2 of 4 records)
```

```
+---+-----+-----+-----+
| id|  name|  surname|height|
+---+-----+-----+-----+
|  1|george|washington|  178|
|  4|  john|      jay|   190|
+---+-----+-----+-----+
```

```
Errors count: 1
```

```
Errors:
```

```
+-----+
|                error|
+-----+
|[[  File "/tmp/20...|
+-----+
```

```
Error fields:
```

```
dict_keys(['callsite', 'msg', 'stackTrace', 'input', 'bytesread', 'source',
'dynamicRecord'])
```

```
Error record data:
```

```
callsite : Row(site=' File "/tmp/2060612586885849088", line 549, in <module>\n
sys.exit(main())\n File "/tmp/2060612586885849088", line 523, in main\n
response
= handler(content)\n File "/tmp/2060612586885849088", line 197, in execute_request
\n result = node.execute()\n File "/tmp/2060612586885849088", line 103, in
execute\n exec(code, global_dict)\n File "<stdin>", line 10, in <module>\n
File "/opt/amazon/lib/python3.6/site-packages/awsglue/dynamicframe.py", line 625, in
from_options\n format_options, transformation_ctx, push_down_predicate, **kwargs)\n
File "/opt/amazon/lib/python3.6/site-packages/awsglue/context.py", line 233, in
create_dynamic_frame_from_options\n source.setFormat(format, **format_options)\n',
info='')
```

```
msg : error in jackson reader
```

```
stackTrace : com.fasterxml.jackson.core.JsonParseException: Unexpected character
('{ ' (code 123)): was expecting either valid name character (for unquoted name) or
double-quote (for quoted) to start field name
at [Source: com.amazonaws.services.glue.readers.BufferedStream@73492578; line: 3,
column: 2]
at com.fasterxml.jackson.core.JsonParser._constructError(JsonParser.java:1581)
```

```
at
com.fasterxml.jackson.core.base.ParserMinimalBase._reportError(ParserMinimalBase.java:533)
at
com.fasterxml.jackson.core.base.ParserMinimalBase._reportUnexpectedChar(ParserMinimalBase.java:
at
com.fasterxml.jackson.core.json.UTF8StreamJsonParser._handleOddName(UTF8StreamJsonParser.java:
at
com.fasterxml.jackson.core.json.UTF8StreamJsonParser._parseName(UTF8StreamJsonParser.java:1650)
at
com.fasterxml.jackson.core.json.UTF8StreamJsonParser.nextToken(UTF8StreamJsonParser.java:740)
at com.amazonaws.services.glue.readers.JacksonReader$$anonfun$hasNextGoodToken
$1.apply(JacksonReader.scala:57)
at com.amazonaws.services.glue.readers.JacksonReader$$anonfun$hasNextGoodToken
$1.apply(JacksonReader.scala:57)
at scala.collection.Iterator$$anon$9.next(Iterator.scala:162)
at scala.collection.Iterator$$anon$16.hasNext(Iterator.scala:599)
at scala.collection.Iterator$$anon$16.hasNext(Iterator.scala:598)
at scala.collection.Iterator$class.foreach(Iterator.scala:891)
at scala.collection.AbstractIterator.foreach(Iterator.scala:1334)
at com.amazonaws.services.glue.readers.JacksonReader$$anonfun
$1.apply(JacksonReader.scala:120)
at com.amazonaws.services.glue.readers.JacksonReader$$anonfun
$1.apply(JacksonReader.scala:116)
at
com.amazonaws.services.glue.DynamicRecordBuilder.handleError(DynamicRecordBuilder.scala:209)
at
com.amazonaws.services.glue.DynamicRecordBuilder.handleErrorWithException(DynamicRecordBuilder
at
com.amazonaws.services.glue.readers.JacksonReader.nextFailSafe(JacksonReader.scala:116)
at com.amazonaws.services.glue.readers.JacksonReader.next(JacksonReader.scala:109)
at com.amazonaws.services.glue.readers.JSONReader.next(JSONReader.scala:247)
at
com.amazonaws.services.glue.hadoop.TapeHadoopRecordReaderSplittable.nextKeyValue(TapeHadoopRec
at org.apache.spark.rdd.NewHadoopRDD$$anon$1.hasNext(NewHadoopRDD.scala:230)
at org.apache.spark.InterruptibleIterator.hasNext(InterruptibleIterator.scala:37)
at scala.collection.Iterator$$anon$11.hasNext(Iterator.scala:409)
at scala.collection.Iterator$$anon$11.hasNext(Iterator.scala:409)
at scala.collection.Iterator$$anon$13.hasNext(Iterator.scala:462)
at scala.collection.Iterator$$anon$11.hasNext(Iterator.scala:409)
at scala.collection.Iterator$$anon$11.hasNext(Iterator.scala:409)
at scala.collection.Iterator$$anon$13.hasNext(Iterator.scala:462)
at scala.collection.Iterator$$anon$11.hasNext(Iterator.scala:409)
at scala.collection.Iterator$$anon$11.hasNext(Iterator.scala:409)
at org.apache.spark.sql.execution.SparkPlan$$anonfun$2.apply(SparkPlan.scala:255)
```

```
at org.apache.spark.sql.execution.SparkPlan$$anonfun$2.apply(SparkPlan.scala:247)
at org.apache.spark.rdd.RDD$$anonfun$mapPartitionsInternal$1$$anonfun$apply
$24.apply(RDD.scala:836)
at org.apache.spark.rdd.RDD$$anonfun$mapPartitionsInternal$1$$anonfun$apply
$24.apply(RDD.scala:836)
at org.apache.spark.rdd.MapPartitionsRDD.compute(MapPartitionsRDD.scala:52)
at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:324)
at org.apache.spark.rdd.RDD.iterator(RDD.scala:288)
at org.apache.spark.rdd.MapPartitionsRDD.compute(MapPartitionsRDD.scala:52)
at org.apache.spark.rdd.RDD.computeOrReadCheckpoint(RDD.scala:324)
at org.apache.spark.rdd.RDD.iterator(RDD.scala:288)
at org.apache.spark.scheduler.ResultTask.runTask(ResultTask.scala:90)
at org.apache.spark.scheduler.Task.run(Task.scala:121)
at org.apache.spark.executor.Executor$TaskRunner$$anonfun$10.apply(Executor.scala:408)
at org.apache.spark.util.Utils$.tryWithSafeFinally(Utils.scala:1360)
at org.apache.spark.executor.Executor$TaskRunner.run(Executor.scala:414)
at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1149)
at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:624)
at java.lang.Thread.run(Thread.java:750)
```

input :

bytesread : 252

source :

dynamicRecord : Row(id=2, name='benjamin', surname='franklin')

errorsCount

`errorsCount()` - DynamicFrame 内のエラーの総数を返します。

stageErrorsCount

`stageErrorsCount` - この DynamicFrame を生成するプロセスで発生したエラーの数を返します。

DynamicFrameCollection クラス

DynamicFrameCollection は [DynamicFrame クラス](#) オブジェクトのディクショナリで、そのキーは DynamicFrames の名前、値は DynamicFrame オブジェクトです。

`__init__`

`__init__(dynamic_frames, glue_ctx)`

- `dynamic_frames` – [DynamicFrame クラス](#) オブジェクトのディクショナリ。
- `glue_ctx` – [GlueContext クラス](#) オブジェクト。

キー

`keys()` – このコレクション内のキーのリストを返します。これは一般的に、対応する `DynamicFrame` 値の名前で構成されます。

[値]

`values(key)` – このコレクション内の `DynamicFrame` 値のリストを返します。

選択

`select(key)`

指定されたキー (一般に `DynamicFrame` の名前) に対応する `DynamicFrame` を返します。

- `key` – `DynamicFrameCollection` 内のキー。通常は `DynamicFrame` の名前を表します。

マッピング

`map(callable, transformation_ctx="")`

渡された関数を使用して、このコレクション内の `DynamicFrameCollection` に基づいた新しい `DynamicFrames` を作成して返します。

- `callable` – `DynamicFrame` と指定された変換コンテキストをパラメータとして取り、`DynamicFrame` を返す関数。
- `transformation_ctx` – 呼び出し可能なものによって使用される変換コンテキスト (省略可能)。

Flatmap

`flatmap(f, transformation_ctx="")`

渡された関数を使用して、このコレクション内の `DynamicFrameCollection` に基づいた新しい `DynamicFrames` を作成して返します。

- `f` - `DynamicFrame` をパラメータとして取り、`DynamicFrame` または `DynamicFrameCollection` を返す関数。
- `transformation_ctx` - 関数で使われる変換コンテキスト (省略可能)。

DynamicFrameWriter クラス

方法

- [__init__](#)
- [from_options](#)
- [from_catalog](#)
- [from_jdbc_conf](#)

`__init__`

`__init__(glue_context)`

- `glue_context` - 使用する [GlueContext クラス](#)。

`from_options`

`from_options(frame, connection_type, connection_options={}, format=None, format_options={}, transformation_ctx="")`

指定された接続と形式を使用して `DynamicFrame` を書き込みます。

- `frame` - 書き込む `DynamicFrame`。
- `connection_type` - 接続タイプ。有効な値には、`s3`、`mysql`、`postgresql`、`redshift`、`sqlserver`、および `oracle` があります。
- `connection_options` - 接続オプション (パスやデータベーステーブルなど) (オプション)。s3 の `connection_type` では、Amazon S3 パスが定義されています。

```
connection_options = {"path": "s3://aws-glue-target/temp"}
```

JDBC 接続の場合、いくつかのプロパティを定義する必要があります。データベース名は URL の一部である必要があることに注意してください。オプションで接続オプションに含めることができます。

⚠ Warning

スクリプトにパスワードを保存することはお勧めしません。AWS Secrets Manager または AWS Glue データカタログから取得する場合には、`boto3` を使用することを検討してください。

```
connection_options = {"url": "jdbc-url/database", "user": "username",  
  "password": passwordVariable, "dbtable": "table-name", "redshiftTmpDir": "s3-tempdir-path"}
```

`dbtable` プロパティは JDBC テーブルの名前です。データベース内でスキーマをサポートする JDBC データストアの場合、`schema.table-name` を指定します。スキーマを指定しない場合、デフォルトの「パブリック」スキーマが使用されます。

詳しくは、「[AWS Glue for Spark での ETL の接続タイプとオプション](#)」を参照してください。

- `format` - 形式の仕様 (オプション)。これは、Amazon Simple Storage Service (Amazon S3)、または複数の形式をサポートする AWS Glue 接続で使用されます。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `format_options` - 指定された形式についてのオプション。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。

`from_catalog`

```
from_catalog(frame, name_space, table_name, redshift_tmp_dir="",  
transformation_ctx="")
```

指定されたカタログデータベースとテーブル名を使用して `DynamicFrame` を書き込みます。

- `frame` - 書き込む `DynamicFrame`。
- `name_space` - 使用するデータベース。
- `table_name` - 使用する `table_name`。
- `redshift_tmp_dir` - 使用する Amazon Redshift の一時ディレクトリ (オプション)。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。

- `additional_options` – AWS Glue で使用する追加のオプション。

Lake Formation governed table に書き込むには、以下の追加オプションを使用することができます。

- `transactionId` — (文字列) governed table への書き込みを行うトランザクション ID。このトランザクションをコミットまたは中断することはできません。書き込みが失敗します。
- `callDeleteObjectsOnCancel` - (プール値、オプション) `true` (デフォルト) に設定すると、オブジェクトが Amazon S3 に書き込まれた後、AWS Glue で自動的に `DeleteObjectsOnCancel` API を呼び出します。詳細については、AWS Lake Formation Developer Guide の「[DeleteObjectsOnCancel](#)」を参照してください。

Example 例: Lake Formation の governed table への書き込み

```
txId = glueContext.start_transaction(read_only=False)
glueContext.write_dynamic_frame.from_catalog(
    frame=dyf,
    database = db,
    table_name = tbl,
    transformation_ctx = "datasource0",
    additional_options={"transactionId":txId})
...
glueContext.commit_transaction(txId)
```

`from_jdbc_conf`

```
from_jdbc_conf(frame, catalog_connection, connection_options={},
redshift_tmp_dir = "", transformation_ctx="")
```

指定された JDBC 接続情報を使用して `DynamicFrame` を書き込みます。

- `frame` - 書き込む `DynamicFrame`。
- `catalog_connection` - 使用するカタログ接続。
- `connection_options` - 接続オプション (パスやデータベーステーブルなど) (オプション)。
- `redshift_tmp_dir` - 使用する Amazon Redshift の一時ディレクトリ (オプション)。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。

write_dynamic_frame の例

この例では、POSIX パス引数を持つ S3 の connection_type を connection_options で使用して、ローカルストレージへの書き込みを許可し、出力をローカルに書き込みます。

```
glueContext.write_dynamic_frame.from_options(\
    frame = dyf_splitFields,\
    connection_options = {'path': '/home/glue/GlueLocalOutput/'},\
    connection_type = 's3',\
    format = 'json')
```

DynamicFrameReader クラス

— methods —

- [__init__](#)
- [from_rdd](#)
- [from_options](#)
- [from_catalog](#)

[__init__](#)

[__init__\(glue_context\)](#)

- glue_context - 使用する [GlueContext クラス](#)。

[from_rdd](#)

[from_rdd\(data, name, schema=None, sampleRatio=None\)](#)

Resilient Distributed Dataset (RDD) から DynamicFrame を読み取ります。

- data - 読み取り元のデータセット。
- name - 読み取り元の名前。
- schema - 読み取るスキーマ (オプション)。
- sampleRatio - サンプル比率 (オプション)。

from_options

```
from_options(connection_type, connection_options={}, format=None,
              format_options={}, transformation_ctx="")
```

指定された接続と形式を使用して DynamicFrame を読み込みます。

- `connection_type` - 接続タイプ。有効な値には、`s3`、`mysql`、`postgresql`、`redshift`、`sqlserver`、`oracle`、`dynamodb`、`snowflake` が含まれています。
- `connection_options` - 接続オプション (パスやデータベーステーブルなど) (オプション)。詳細については、「[AWS Glue for Spark](#)」の「[ETL の接続タイプとオプション](#)」を参照してください。s3 の `connection_type` の場合、Amazon S3 のパスは配列で定義されます。

```
connection_options = {"paths": [ "s3://mybucket/object_a", "s3://mybucket/object_b" ]}
```

JDBC 接続の場合、いくつかのプロパティを定義する必要があります。データベース名は URL の一部である必要があることに注意してください。オプションで接続オプションに含めることができます。

Warning

スクリプトにパスワードを保存することはお勧めしません。boto3 を使用して AWS Secrets Manager または AWS Glue データカタログからそれらを取得することを検討してください。

```
connection_options = {"url": "jdbc-url/database", "user": "username",
                      "password": passwordVariable, "dbtable": "table-name", "redshiftTmpDir": "s3-tempdir-path" }
```

並列読み込みを実行する JDBC 接続の場合、ハッシュフィールドオプションを設定できます。例:

```
connection_options = {"url": "jdbc-url/database", "user": "username",
                      "password": passwordVariable, "dbtable": "table-name", "redshiftTmpDir": "s3-tempdir-path" , "hashfield": "month" }
```

詳しくは、「[JDBC テーブルからの並列読み取り](#)」を参照してください。

- `format` – 形式の仕様 (オプション)。これは、Amazon Simple Storage Service (Amazon S3)、または複数の形式をサポートする AWS Glue 接続で使用されます。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `format_options` – 指定された形式についてのオプション。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。
- `push_down_predicate` – データセットのすべてのファイルをリストアップして読み取る必要がないフィルタパーティション。詳細については、「[Pre-Filtering Using Pushdown Predicates](#)」を参照してください。

`from_catalog`

```
from_catalog(database, table_name, redshift_tmp_dir="",  
transformation_ctx="", push_down_predicate="", additional_options={})
```

指定されたカタログの名前空間とテーブル名を使用して `DynamicFrame` を読み取ります。

- `database` - 読み取り元のデータベース。
- `table_name` - 読み取り元のテーブルの名前。
- `redshift_tmp_dir` – 使用する Amazon Redshift の一時ディレクトリ (Redshift からデータを読み取らない場合はオプション)。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。
- `push_down_predicate` – データセットのすべてのファイルをリストアップして読み取る必要がないフィルタパーティション。詳しくは、「[プッシュダウン述語を使用した事前フィルタ処理](#)」を参照してください。
- `additional_options` – AWS Glue で使用する追加のオプション。
 - 並列読み込みを実行する JDBC 接続を使用するには、`hashfield`、`hashexpression`、または `hashpartitions` オプションを設定できます。例:

```
additional_options = {"hashfield": "month"}
```

詳しくは、「[JDBC テーブルからの並列読み取り](#)」を参照してください。

- インデックス列に基づいたフィルタリングのためにカタログ式を渡すには、`catalogPartitionPredicate` オプションを使用できます。

`catalogPartitionPredicate` – カタログ式を渡して、インデックス列に基づいたフィルタリングができます。これにより、フィルタリングをサーバー側で処理できます。詳細については、「[AWS Glue パーティションインデックス](#)」を参照してください。`push_down_predicate` と `catalogPartitionPredicate` では、異なる構文が使用されることに注意してください。前者では Spark SQL の標準構文を使用し、後者では JSQL パーサーを使用します。

詳細については、「[AWS Glue での ETL 出力のパーティションの管理](#)」を参照してください。

GlueContext クラス

Apache Spark [SparkContext](#) オブジェクトをラップし、Apache Spark プラットフォームとやり取りするためのメカニズムを提供します。

`__init__`

`__init__(sparkContext)`

- `sparkContext` - 使用する Apache Spark のコンテキスト。

[作成中]

- [__init__](#)
- [getSource](#)
- [create_dynamic_frame_from_rdd](#)
- [create_dynamic_frame_from_catalog](#)
- [create_dynamic_frame_from_options](#)
- [create_sample_dynamic_frame_from_catalog](#)
- [create_sample_dynamic_frame_from_options](#)
- [add_ingestion_time_columns](#)
- [create_data_frame_from_catalog](#)
- [create_data_frame_from_options](#)
- [forEachBatch](#)

getSource

getSource(connection_type, transformation_ctx = "", **options)

外部ソースから DynamicFrames を読み取るために使用できる DataSource オブジェクトを作成します。

- connection_type - 使用する接続タイプ (Amazon Simple Storage Service (Amazon S3)、Amazon Redshift、JDBC など)。有効な値には、s3、mysql、postgresql、redshift、sqlserver、oracle および dynamodb があります。
- transformation_ctx - 使用する変換コンテキスト (オプション)。
- options - オプションの名前と値のペアのコレクション。詳しくは、[「AWS Glue for Spark での ETL の接続タイプとオプション」](#)を参照してください。

以下は、getSource の使用例です。

```
>>> data_source = context.getSource("file", paths=["/in/path"])
>>> data_source.setFormat("json")
>>> myFrame = data_source.getFrame()
```

create_dynamic_frame_from_rdd

create_dynamic_frame_from_rdd(data, name, schema=None, sample_ratio=None, transformation_ctx="")

Apache Spark Resilient Distributed Dataset (RDD) から作成された DynamicFrame を返します。

- data - 使用するデータソース。
- name - 使用するデータの名前。
- schema - 使用するスキーマ (オプション)。
- sample_ratio - 使用するサンプル比率 (オプション)。
- transformation_ctx - 使用する変換コンテキスト (オプション)。

```
create_dynamic_frame_from_catalog
```

```
create_dynamic_frame_from_catalog(database, table_name, redshift_tmp_dir,  
transformation_ctx = "", push_down_predicate= "", additional_options = {},  
catalog_id = None)
```

データカタログデータベースとテーブル名を使用して作成された DynamicFrame を返します。このメソッドを使用する場合は、指定された AWS Glue Data Catalog テーブルのテーブルプロパティ format_options を介して を指定し、 additional_options 引数を使用してその他のオプションを指定します。

- Database - 読み取り元のデータベース。
- table_name - 読み取り元のテーブルの名前。
- redshift_tmp_dir - 使用する Amazon Redshift の一時ディレクトリ (オプション)。
- transformation_ctx - 使用する変換コンテキスト (オプション)。
- push_down_predicate - データセットのすべてのファイルを一括アップロードして読み取る必要がないフィルタパーティション。サポートされているソースと制限については、[「AWS Glue ETL でのプッシュダウンによる読み取りの最適化」](#)を参照してください。詳細については、[「プッシュダウン述語を使用した事前フィルタ処理」](#)を参照してください。
- additional_options - オプションの名前と値のペアのコレクション。[AWS Glue for Spark での ETL の接続タイプとオプション](#) でリストされている使用可能なオプション (endpointUrl、streamName、bootstrap.servers、security.protocol、topicName、class および delimiter を除く)。別のオプションとして、catalogPartitionPredicate もサポートされています。

catalogPartitionPredicate - カタログ式を渡して、インデックス列に基づいたフィルタリングができます。これにより、フィルタリングをサーバー側で処理できます。詳細については、[「AWS Glue パーティションインデックス」](#)を参照してください。push_down_predicate と catalogPartitionPredicate では、異なる構文が使用されることに注意してください。前者では Spark SQL の標準構文を使用し、後者では JSQL パーサーを使用します。

- catalog_id - 現在アクセスされているデータカタログのカタログ ID (アカウント ID)。None の場合は、呼び出し元のアカウント ID のデフォルトが使用されます。

```
create_dynamic_frame_from_options
```

```
create_dynamic_frame_from_options(connection_type, connection_options={},  
format=None, format_options={}, transformation_ctx = "")
```

指定された接続と形式で作成された DynamicFrame を返します。

- `connection_type` – 接続タイプ (Amazon S3、Amazon Redshift、JDBC など)。有効な値には、`s3`、`mysql`、`postgresql`、`redshift`、`sqlserver`、`oracle` および `dynamodb` があります。
- `connection_options` – パスやデータベーステーブルなど接続オプション (オプション)。s3 の `connection_type` に関しては、Amazon S3 パスのリストが定義されています。

```
connection_options = {"paths": ["s3://aws-glue-target/temp"]}
```

JDBC 接続の場合、いくつかのプロパティを定義する必要があります。データベース名は URL の一部である必要があることに注意してください。オプションで接続オプションに含めることができます。

Warning

スクリプトにパスワードを保存することはお勧めしません。boto3 を使用して AWS Secrets Manager または AWS Glue データカタログからそれらを取得することを検討してください。

```
connection_options = {"url": "jdbc-url/database", "user": "username",  
"password": passwordVariable, "dbtable": "table-name", "redshiftTmpDir": "s3-tempdir-path"}
```

`dbtable` プロパティは JDBC テーブルの名前です。データベース内でスキーマをサポートする JDBC データストアの場合、`schema.table-name` を指定します。スキーマを指定しない場合、デフォルトの「パブリック」スキーマが使用されます。

詳細については、「[AWS Glue for Spark での ETL の接続タイプとオプション](#)」を参照してください。

- `format` – 形式の仕様。Amazon S3 や、複数のフォーマットをサポートする AWS Glue 接続の場合に使用します。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `format_options` – 指定された形式についてのオプション。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。
- `push_down_predicate` – データセットのすべてのファイルをリストアップして読み取る必要がないフィルタパーティション。サポートされているソースと制限については、「[AWS Glue ETL でのプッシュダウンによる読み取りの最適化](#)」を参照してください。詳細については、「[プッシュダウン述語を使用した事前フィルタ処理](#)」を参照してください。

`create_sample_dynamic_frame_from_catalog`

```
create_sample_dynamic_frame_from_catalog(database, table_name, num,
redshift_tmp_dir, transformation_ctx = "", push_down_predicate= "",
additional_options = {}, sample_options = {}, catalog_id = None)
```

データカタログデータベースとテーブル名を使用して作成されたサンプル `DynamicFrame` を返します。`DynamicFrame` にはデータソースからの最初の `num` レコードのみが含まれます。

- `database` - 読み取り元のデータベース。
- `table_name` - 読み取り元のテーブルの名前。
- `num` — 返されるサンプル動的フレーム内のレコードの最大数。
- `redshift_tmp_dir` - 使用する Amazon Redshift の一時ディレクトリ (オプション)。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。
- `push_down_predicate` – データセットのすべてのファイルをリストアップして読み取る必要がないフィルタパーティション。詳しくは、「[プッシュダウン述語を使用した事前フィルタ処理](#)」を参照してください。
- `additional_options` - オプションの名前と値のペアのコレクション。[AWS Glue for Spark での ETL の接続タイプとオプション](#) でリストされている使用可能なオプション (`endpointUrl`、`streamName`、`bootstrap.servers`、`security.protocol`、`topicName`、`className` および `delimiter` を除く)。
- `sample_options`— サンプルング動作を制御するパラメータ (オプション)。Simple Storage Service (Amazon S3) ソースで現在使用可能なパラメータ:

- `maxSamplePartitions` — サンプルングが読み取るパーティションの最大数。デフォルト値は 10 です
- `maxSampleFilesPerPartition` — サンプルングが 1 つのパーティションで読み取るファイルの最大数。デフォルト値は 10 です。

これらのパラメータは、ファイル一覧で消費される時間を短縮するのに役立ちます。例えば、データセットに 1000 個のパーティションがあり、各パーティションには 10 個のファイルがあるとします。10,000 個のファイルをすべてリスト表示する代わりに、`maxSamplePartitions = 10` および `maxSampleFilesPerPartition = 10` と設定した場合、サンプルングでは、最初の 10 個のパーティションの最初の 10 個のファイルのみがリスト表示されて読み込まれ、合計で $10 \times 10 = 100$ 個のファイルとなります。

- `catalog_id` – 現在アクセスされているデータカタログのカタログ ID (データカタログのアカウント ID)。デフォルトでは、None に設定されています。None のデフォルト値は、サービス内の呼び出し元アカウントのカタログ ID になります。

`create_sample_dynamic_frame_from_options`

```
create_sample_dynamic_frame_from_options(connection_type,  
connection_options={}, num, sample_options={}, format=None,  
format_options={}, transformation_ctx = "")
```

指定された接続と形式で作成されたサンプル `DynamicFrame` を返します。`DynamicFrame` にはデータソースからの最初の `num` レコードのみが含まれます。

- `connection_type` – 接続タイプ (Amazon S3、Amazon Redshift、JDBC など)。有効な値には、`s3`、`mysql`、`postgres`、`redshift`、`sqlserver`、`oracle` および `dynamodb` があります。
- `connection_options` – パスやデータベーステーブルなど接続オプション (オプション)。詳しくは、「[AWS Glue for Spark での ETL の接続タイプとオプション](#)」を参照してください。
- `num` — 返されるサンプル動的フレーム内のレコードの最大数。
- `sample_options` — サンプルング動作を制御するパラメータ (オプション)。Simple Storage Service (Amazon S3) ソースで現在使用可能なパラメータ:
 - `maxSamplePartitions` — サンプルングが読み取るパーティションの最大数。デフォルト値は 10 です
 - `maxSampleFilesPerPartition` — サンプルングが 1 つのパーティションで読み取るファイルの最大数。デフォルト値は 10 です。

これらのパラメータは、ファイル一覧で消費される時間を短縮するのに役立ちます。例えば、データセットに 1000 個のパーティションがあり、各パーティションには 10 個のファイルがあるとします。10,000 個のファイルをすべてリスト表示する代わりに、`maxSamplePartitions = 10` および `maxSampleFilesPerPartition = 10` と設定した場合、サンプリングでは、最初の 10 個のパーティションの最初の 10 個のファイルのみがリスト表示されて読み込まれ、合計で $10 \times 10 = 100$ 個のファイルとなります。

- `format` – 形式の仕様。Amazon S3 や、複数のフォーマットをサポートする AWS Glue 接続の場合に使用します。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `format_options` – 指定された形式についてのオプション。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。
- `push_down_predicate` – データセットのすべてのファイルをリストアップして読み取る必要がないフィルタパーティション。詳しくは、「[プッシュダウン述語を使用した事前フィルタ処理](#)」を参照してください。

`add_ingestion_time_columns`

`add_ingestion_time_columns(dataFrame, timeGranularity = "")`

入力 DataFrame への取り込み時間列

(`ingest_year`、`ingest_month`、`ingest_day`、`ingest_hour`、`ingest_minute`) を追加します。Amazon S3 のデータカタログテーブルをターゲットとして指定する場合、この関数は、AWS Glue により生成されたスクリプト内で自動的に生成されます。この関数は、出力テーブル上で、取り込み時間列があるパーティションを自動的に更新します。これにより、入力データに明示的な取り込み時間列を指定しなくても、取り込み時間において出力データの自動的なパーティション化が行えます。

- `dataFrame` – 取り込み時間列の追加先である `dataFrame`。
- `timeGranularity` – 時間列の詳細度。有効な値は「`day`」、「`hour`」、および「`minute`」です。例えば、関数に対し「`hour`」が渡された場合、元の `dataFrame` は「`ingest_year`」、「`ingest_month`」、「`ingest_day`」に加え「`ingest_hour`」の時間列を持つこととなります。

時間の詳細度列を追加した後、そのデータフレームを返します。

例：

```
dynamic_frame = DynamicFrame.fromDF(glueContext.add_ingestion_time_columns(dataFrame,
"hour"))
```

`create_data_frame_from_catalog`

```
create_data_frame_from_catalog(database, table_name, transformation_ctx =
"", additional_options = {})
```

データカタログテーブルからの情報を使用して作成された DataFrame を返します。

- `database` – 読み取り元のデータカタログデータベース。
- `table_name` – 読み取り元のデータカタログテーブルの名前。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。
- `additional_options` - オプションの名前と値のペアのコレクション。可能なオプションには、`startingPosition`、`maxFetchTimeInMs`、および `startingOffsets` など、ストリーミングソース用として [AWS Glue for Spark での ETL の接続タイプとオプション](#) にリストされているものが含まれます。
- `useSparkDataSource` – `true` に設定すると、AWS Glue はネイティブ Spark データソース API を使用してテーブルを読み取るように強制されます。Spark データソース API でサポートされている形式は、AVRO、バイナリ、CSV、JSON、ORC、Parquet、およびテキストです。データカタログテーブルでは、`classification` プロパティを使用して形式を指定します。Spark データソース API の詳細については、公式の [Apache Spark ドキュメント](#) を参照してください。

`create_data_frame_from_catalog` を `useSparkDataSource` と併用すると次のようなメリットがあります。

- DataFrame を直接返し、`create_dynamic_frame.from_catalog().toDF()` の代替手段になります。
- ネイティブ形式の AWS Lake Formation テーブルレベルのアクセス許可コントロールをサポートします。
- AWS Lake Formation テーブルレベルのアクセス許可制御なしでのデータレイク形式の読み取りをサポートします。詳細については、「[AWS Glue ETL ジョブでのデータレイクフレームワークの使用](#)」を参照してください。

を有効にするとuseSparkDataSource、additional_options必要に応じてに [Spark データソースオプション](#)を追加することもできます。AWS Glue はこれらのオプションを Spark リーダーに直接渡します。

- useCatalogSchema – true に設定すると、AWS Glue は結果の にデータカタログスキーマを適用しますDataFrame。true に設定しなければ、リーダーはデータからスキーマを推測します。useCatalogSchema を有効にする場合は、useSparkDataSource も true に設定する必要があります。

制約事項

useSparkDataSource オプションを使用する際には、次の制限事項を考慮してください。

- を使用するとuseSparkDataSource、AWS Glue はDataFrame元の Spark セッションとは異なる新しい Spark セッションを作成します。
- Spark DataFrame パーティションフィルタリングは、次の AWS Glue 機能では機能しません。
 - [ジョブのブックマーク](#)
 - [Amazon S3 ストレージクラスの除外](#)
 - [カタログパーティション述語](#)

これらの機能でパーティションフィルタリングを使用するには、AWS Glue プッシュダウン述語を使用できます。詳細については、「[プッシュダウン述語を使用した事前フィルタ処理](#)」を参照してください。分割されていない列のフィルタリングは影響を受けません。

次のスクリプト例は、excludeStorageClasses オプションを使用してパーティションフィルタリングを実行する誤った方法を示しています。

```
// Incorrect partition filtering using Spark filter with excludeStorageClasses
read_df = glueContext.create_data_frame.from_catalog(
    database=database_name,
    table_name=table_name,
    additional_options = {
        "useSparkDataSource": True,
        "excludeStorageClasses" : ["GLACIER", "DEEP_ARCHIVE"]
    }
)

// Suppose year and month are partition keys.
// Filtering on year and month won't work, the filtered_df will still
```

```
// contain data with other year/month values.
filtered_df = read_df.filter("year == '2017 and month == '04' and 'state == 'CA'")
```

次のスクリプト例は、`excludeStorageClasses` オプションを使用してパーティションフィルタリングを実行するためにプッシュダウン述語を使用する正しい方法を示しています。

```
// Correct partition filtering using the AWS Glue pushdown predicate
// with excludeStorageClasses
read_df = glueContext.create_data_frame.from_catalog(
    database=database_name,
    table_name=table_name,
    // Use AWS Glue pushdown predicate to perform partition filtering
    push_down_predicate = "(year=='2017' and month=='04')",
    additional_options = {
        "useSparkDataSource": True,
        "excludeStorageClasses" : ["GLACIER", "DEEP_ARCHIVE"]
    }
)

// Use Spark filter only on non-partitioned columns
filtered_df = read_df.filter("state == 'CA'")
```

例: Spark データソースリーダーを使用した CSV テーブルの作成

```
// Read a CSV table with '\t' as separator
read_df = glueContext.create_data_frame.from_catalog(
    database=<database_name>,
    table_name=<table_name>,
    additional_options = {"useSparkDataSource": True, "sep": '\t'}
)
```

`create_data_frame_from_options`

`create_data_frame_from_options(connection_type, connection_options={}, format=None, format_options={}, transformation_ctx = "")`

この API は廃止されました。代わりに、`getSource()` API を使用してください。指定された接続と形式で作成された `DataFrame` を返します。この関数は、AWS Glue ストリーミングソースのみで使用してください。

- `connection_type` – ストリーミング接続タイプ。有効な値は、`kinesis` および `kafka` です。

- `connection_options` – 接続オプション。Kinesis と Kafka では異なります。各ストリーミングデータソースのすべての接続オプションの一覧は、[AWS Glue for Spark での ETL の接続タイプとオプション](#) で確認いただけます。ストリーミング接続オプションについては、以下の違いに注意してください。
 - Kinesis ストリーミングのソースには `streamARN`、`startingPosition`、`inferSchema`、および `classification` が必要です。
 - Kafka ストリーミングのソースには `connectionName`、`topicName`、`startingOffsets`、`inferSchema`、および `classification` が必要です。
- `format` – 形式の仕様。Amazon S3 や、複数のフォーマットをサポートする AWS Glue 接続の場合に使用します。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `format_options` – 指定された形式についてのオプション。サポートされる形式オプションについては、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。

Amazon Kinesis ストリーミングソースの例:

```
kinesis_options =
  { "streamARN": "arn:aws:kinesis:us-east-2:777788889999:stream/fromOptionsStream",
    "startingPosition": "TRIM_HORIZON",
    "inferSchema": "true",
    "classification": "json"
  }
data_frame_datasource0 =
  glueContext.create_data_frame.from_options(connection_type="kinesis",
  connection_options=kinesis_options)
```

Kafka ストリーミングソースの例 :

```
kafka_options =
  { "connectionName": "ConfluentKafka",
    "topicName": "kafka-auth-topic",
    "startingOffsets": "earliest",
    "inferSchema": "true",
    "classification": "json"
  }
```

```
data_frame_datasource0 =  
    glueContext.create_data_frame.from_options(connection_type="kafka",  
    connection_options=kafka_options)
```

forEachBatch

forEachBatch(frame, batch_function, options)

ストリーミングソースから読み取られるすべてのマイクロバッチに渡される、batch_function を適用します。

- frame – 現在のマイクロバッチ DataFrame を含む。
- batch_function – すべてのマイクロバッチに適用される関数。
- options – マイクロバッチの処理方法に関する情報を保持している、キーと値のペアの集合。以下のような必須オプションがあります。
 - windowSize – 各バッチの処理にかかる時間。
 - checkpointLocation – ストリーミング ETL ジョブ用に、チェックポイントが格納される場所。
 - batchMaxRetries – 失敗した場合にこのバッチを再試行する最大回数。デフォルト値は 3 です。このオプションは、Glue バージョン 2.0 以降でのみ設定可能です。

例:

```
glueContext.forEachBatch(  
    frame = data_frame_datasource0,  
    batch_function = processBatch,  
    options = {  
        "windowSize": "100 seconds",  
        "checkpointLocation": "s3://kafka-auth-dataplane/confluent-test/output/  
checkpoint/"  
    }  
)  
  
def processBatch(data_frame, batchId):  
    if (data_frame.count() > 0):  
        datasource0 = DynamicFrame.fromDF(  
            glueContext.add_ingestion_time_columns(data_frame, "hour"),  
            glueContext, "from_data_frame"  
        )
```

```
additionalOptions_datasink1 = {"enableUpdateCatalog": True}
additionalOptions_datasink1["partitionKeys"] = ["ingest_yr", "ingest_mo",
"ingest_day"]
datasink1 = glueContext.write_dynamic_frame.from_catalog(
    frame = datasource0,
    database = "tempdb",
    table_name = "kafka-auth-table-output",
    transformation_ctx = "datasink1",
    additional_options = additionalOptions_datasink1
)
```

Simple Storage Service (Amazon S3) でのデータセットの操作

- [purge_table](#)
- [purge_s3_path](#)
- [transition_table](#)
- [transition_s3_path](#)

purge_table

```
purge_table(catalog_id=None, database="", table_name="", options={},
transformation_ctx="")
```

指定したカタログのデータベースとテーブルのファイルを Simple Storage Service (Amazon S3) から削除します。パーティション内のすべてのファイルが削除されると、そのパーティションもカタログから削除されます。

削除したオブジェクトを回復できるようにするには、Amazon S3 バケットで[オブジェクトのバージョニング](#)を有効にします。オブジェクトバージョニングが有効になっていないバケットからオブジェクトが削除された場合、そのオブジェクトは復元できません。バージョニングが有効にされているバケットで削除されたオブジェクトを復元する方法の詳細については、AWS Support ナレッジセンターで「[バージョニングが有効なバケットで削除された Simple Storage Service \(Amazon S3\) オブジェクトを取得するにはどうすればよいですか?](#)」を参照してください。

- catalog_id – 現在アクセスされているデータカタログのカタログ ID (データカタログのアカウント ID)。デフォルトでは、None に設定されています。None のデフォルト値は、サービス内の呼び出し元アカウントのカタログ ID になります。
- database – 使用するデータベース。

- `table_name` – 使用するテーブルの名前。
- `options` – 削除するファイルのフィルタリングと、マニフェストファイルの生成のためオプション。
 - `retentionPeriod` – ファイルを保持する期間を時間単位で指定します。保存期間より新しいファイルは保持されます。デフォルトでは 168 時間 (7 日) に設定されています。
 - `partitionPredicate` – この述語を満たすパーティションは削除されます。これらのパーティションの保存期間内のファイルは削除されません。"" を設定 – デフォルトでは空です。
 - `excludeStorageClasses` – `excludeStorageClasses` セット内のストレージクラスを持つファイルは削除されません。デフォルトは `Set()` – 空のセットです。
 - `manifestFilePath` – マニフェストファイルを生成するためのオプションのパス。正常にパーティジされたすべてのファイルが `Success.csv` に記録され、失敗したファイルは `Failed.csv` に記録されます。
- `transformation_ctx` – 使用する変換コンテキスト (オプション)。マニフェストファイルパスで使用されます。

Example

```
glueContext.purge_table("database", "table", {"partitionPredicate": "(month=='march')",  
  "retentionPeriod": 1, "excludeStorageClasses": ["STANDARD_IA"], "manifestFilePath":  
  "s3://bucketmanifest/"})
```

purge_s3_path

purge_s3_path(s3_path, options={}, transformation_ctx="")

指定された Amazon S3 パスからファイルを再帰的に削除します。

削除したオブジェクトを回復できるようにするには、Amazon S3 バケットで[オブジェクトのバージョニング](#)を有効にします。オブジェクトバージョニングが有効になっていないバケットからオブジェクトが削除された場合、そのオブジェクトは復元できません。バージョニングを使用してバケット内の削除されたオブジェクトを復元する方法の詳細については、AWS Support ナレッジセンターの「[削除された Amazon S3 オブジェクトを取得する方法](#)」を参照してください。

- `s3_path` – 削除するファイルを指す Simple Storage Service (Amazon S3) のパス (`s3://<bucket>/<prefix>/` 形式)
- `options` – 削除するファイルのフィルタリングと、マニフェストファイルの生成のためオプション。

- `retentionPeriod` – ファイルを保持する期間を時間単位で指定します。保存期間より新しいファイルは保持されません。デフォルトでは 168 時間 (7 日) に設定されています。
- `excludeStorageClasses` – `excludeStorageClasses` セット内のストレージクラスを持つファイルは削除されません。デフォルトは `Set()` – 空のセットです。
- `manifestFilePath` – マニフェストファイルを生成するためのオプションのパス。正常にパージされたすべてのファイルが `Success.csv` に記録され、失敗したファイルは `Failed.csv` に記録されます。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。マニフェストファイルパスで使用されます。

Example

```
glueContext.purge_s3_path("s3://bucket/path/", {"retentionPeriod": 1,
"excludeStorageClasses": ["STANDARD_IA"], "manifestFilePath": "s3://bucketmanifest/"})
```

transition_table

transition_table(database, table_name, transition_to, options={}, transformation_ctx="", catalog_id=None)

指定されたカタログのデータベースとテーブルのために、Simple Storage Service (Amazon S3) に格納されているファイルの、ストレージクラスを移行します。

任意の 2 つのストレージクラス間で移行できます。GLACIER と DEEP_ARCHIVE のストレージクラスでは、これらのクラスに移行できます。ただし、GLACIER と DEEP_ARCHIVE のストレージクラスからの移行には S3 RESTORE を使用します。

Amazon S3 からファイルまたはパーティションを読み取る AWS Glue ETL ジョブを実行している場合は、一部の Amazon S3 ストレージクラスタイプを除外できます。詳細については、「[Excluding Amazon S3 Storage Classes](#)」を参照してください。

- `database` – 使用するデータベース。
- `table_name` – 使用するテーブルの名前。
- `transition_to` – 移行する先の [Amazon S3 ストレージクラス](#)。
- `options` – 削除するファイルのフィルタリングと、マニフェストファイルの生成のためオプション。

- `retentionPeriod` – ファイルを保持する期間を時間単位で指定します。保存期間より新しいファイルは保持されます。デフォルトでは 168 時間 (7 日) に設定されています。
- `partitionPredicate` – この述語を満たすパーティションは移行されます。これらのパーティションの保存期間内のファイルは移行されません。"" を設定 – デフォルトでは空です。
- `excludeStorageClasses` – `excludeStorageClasses` セット内のストレージクラスを持つファイルは移行されません。デフォルトは `Set()` – 空のセットです。
- `manifestFilePath` – マニフェストファイルを生成するためのオプションのパス。正常に移行されたすべてのファイルが `Success.csv` に記録され、失敗したファイルは `Failed.csv` に記録されます。
- `accountId` – 移行変換を実行する Amazon Web Services アカウント ID。この変換には必須です。
- `roleArn` – 移行変換を実行する AWS ロール。この変換には必須です。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。マニフェストファイルパスで使用されます。
- `catalog_id` – 現在アクセスされているデータカタログのカタログ ID (データカタログのアカウント ID)。デフォルトでは、`None` に設定されています。`None` のデフォルト値は、サービス内の呼び出し元アカウントのカタログ ID になります。

Example

```
glueContext.transition_table("database", "table", "STANDARD_IA", {"retentionPeriod": 1,
"excludeStorageClasses": ["STANDARD_IA"], "manifestFilePath": "s3://bucketmanifest/",
"accountId": "12345678901", "roleArn": "arn:aws:iam::123456789012:user/example-username"})
```

transition_s3_path

```
transition_s3_path(s3_path, transition_to, options={},
transformation_ctx="")
```

指定された Simple Storage Service (Amazon S3) パス内のファイルのストレージクラスを再帰的に移行します。

任意の 2 つのストレージクラス間で移行できます。GLACIER と DEEP_ARCHIVE のストレージクラスでは、これらのクラスに移行できます。ただし、GLACIER と DEEP_ARCHIVE のストレージクラスからの移行には S3 RESTORE を使用します。

Amazon S3 からファイルまたはパーティションを読み取る AWS Glue ETL ジョブを実行している場合は、一部の Amazon S3 ストレージクラスタイプを除外できます。詳細については、「[Excluding Amazon S3 Storage Classes](#)」を参照してください。

- `s3_path` – 移行するファイルの Simple Storage Service (Amazon S3) のパス (`s3://<bucket>/<prefix>/` 形式)
- `transition_to` – 移行する先の [Amazon S3 ストレージクラス](#)。
- `options` – 削除するファイルのフィルタリングと、マニフェストファイルの生成のためオプション。
 - `retentionPeriod` – ファイルを保持する期間を時間単位で指定します。保存期間より新しいファイルは保持されます。デフォルトでは 168 時間 (7 日) に設定されています。
 - `partitionPredicate` – この述語を満たすパーティションは移行されます。これらのパーティションの保存期間内のファイルは移行されません。"" を設定 – デフォルトでは空です。
 - `excludeStorageClasses` – `excludeStorageClasses` セット内のストレージクラスを持つファイルは移行されません。デフォルトは `Set()` – 空のセットです。
 - `manifestFilePath` – マニフェストファイルを生成するためのオプションのパス。正常に移行されたすべてのファイルが `Success.csv` に記録され、失敗したファイルは `Failed.csv` に記録されます。
 - `accountId` – 移行変換を実行する Amazon Web Services アカウント ID。この変換には必須です。
 - `roleArn` – 移行変換を実行する AWS ロール。この変換には必須です。
- `transformation_ctx` – 使用する変換コンテキスト (オプション)。マニフェストファイルパスで使用されます。

Example

```
glueContext.transition_s3_path("s3://bucket/prefix/", "STANDARD_IA",
{"retentionPeriod": 1, "excludeStorageClasses": ["STANDARD_IA"],
"manifestFilePath": "s3://bucketmanifest/", "accountId": "12345678901", "roleArn":
"arn:aws:iam::123456789012:user/example-username"})
```

抽出

- [extract_jdbc_conf](#)

extract_jdbc_conf

extract_jdbc_conf(connection_name, catalog_id = None)

データカタログの AWS Glue 接続オブジェクトから、設定プロパティを使用するキーを持つ dict を返します。

- user - データベースのユーザー名です。
- password - データベースのパスワードです。
- vendor - ベンダーを指定します (mysql、postgresql、oracle、sqlserver など)。
- enforceSSL - 安全な接続が必要かどうかを示すブール文字列です。
- customJDBCCert - 提示された Amazon S3 パスからの特定のクライアント証明書を使用します。
- skipCustomJDBCCertValidation - customJDBCCert が CA によって検証される必要があるかどうかを示すブール文字列です。
- customJDBCCertString - ドライバーの種類に固有のカスタム証明書に関する追加情報です。
- url - (廃止) プロトコル、サーバー、ポートのみを含む JDBC URL です。
- fullUrl - 接続の作成時に入力された JDBC URL です (AWS Glue バージョン 3.0 以降で使用可能)。

JDBC 設定の取得例:

```
jdbc_conf = glueContext.extract_jdbc_conf(connection_name="your_glue_connection_name")
print(jdbc_conf)
>>> {'enforceSSL': 'false', 'skipCustomJDBCCertValidation': 'false', 'url':
'jdbc:mysql://myserver:3306', 'fullUrl': 'jdbc:mysql://myserver:3306/mydb',
'customJDBCCertString': '', 'user': 'admin', 'customJDBCCert': '', 'password': '1234',
'vendor': 'mysql'}
```

トランザクション

- [start_transaction](#)
- [commit_transaction](#)
- [cancel_transaction](#)

start_transaction

start_transaction(read_only)

新しいトランザクションの開始。Lake Formation [startTransaction](#) API を内部的に呼び出します。

- `read_only` — (Boolean) このトランザクションを読み取り専用にするか、または読み取りおよび書き込みを行うかを示します。読み取り専用のトランザクション ID を使用した書き込みは拒否されます。読み取り専用トランザクションはコミットする必要はありません。

トランザクション ID を返します。

commit_transaction

commit_transaction(transaction_id, wait_for_commit = True)

指定されたトランザクションをコミットしようとしています。commit_transaction でトランザクションのコミットが完了する前に戻ることがあります。Lake Formation [startTransaction](#) API を内部的に呼び出します。

- `transaction_id` — (文字列) コミットするトランザクション。
- `wait_for_commit` — (ブール値) commit_transaction がすぐに戻るかどうか指定します。デフォルト値は True です。false の場合、commit_transaction はトランザクションがコミットされるまでポーリングし待機します。最大で 6 回の再試行でエクスポネンシャルバックオフを使用すると、待機時間は 1 分に制限されます。

コミットが完了したかどうかを示すブール値を返します。

cancel_transaction

cancel_transaction(transaction_id)

指定されたトランザクションをキャンセルしようとしています。戻り値は、トランザクションが以前にコミットされた場合は TransactionCommittedException 例外です。Lake Formation [CancelTransaction](#) API を内部的に呼び出します。

- `transaction_id` — (文字列) キャンセルするトランザクション。

書き込み

- [getSink](#)
- [write_dynamic_frame_from_options](#)
- [write_from_options](#)
- [write_dynamic_frame_from_catalog](#)
- [write_data_frame_from_catalog](#)
- [write_dynamic_frame_from_jdbc_conf](#)
- [write_from_jdbc_conf](#)

getSink

getSink(connection_type, format = None, transformation_ctx = "", **options)

外部ソースに DynamicFrames を書き込むために使用できる DataSink オブジェクトを取得します。期待しているシンクを確実に取得するために、SparkSQL format を最初に確認します。

- connection_type – 使用する接続タイプ (Simple Storage Service (Amazon S3)、Amazon Redshift、JDBC など)。有効な値には、s3、mysql、postgresql、redshift、sqlserver、oracle、kinesis、kafka が含まれています。
- format – 使用する SparkSQL 形式 (オプション)。
- transformation_ctx – 使用する変換コンテキスト (オプション)。
- options — 接続オプションの指定に使用される名前と値のペアの集合。指定できる値は以下のとおりです。
 - user と password: 認可用
 - url: データストアのエンドポイント
 - dbtable: ターゲットテーブルの名前。
 - bulkSize: 挿入操作の並列度の度合い

指定できるオプションは、接続タイプによって異なります。追加の値と例については、「[AWS Glue for Spark での ETL の接続タイプとオプション](#)」を参照してください。

例 :

```
>>> data_sink = context.getSink("s3")
>>> data_sink.setFormat("json"),
>>> data_sink.writeFrame(myFrame)
```

`write_dynamic_frame_from_options`

```
write_dynamic_frame_from_options(frame, connection_type,  
connection_options={}, format=None, format_options={}, transformation_ctx =  
"")
```

指定された接続と形式を使用して `DynamicFrame` を書き込み、返します。

- `frame` - 書き込む `DynamicFrame`。
- `connection_type` - 接続タイプ (Amazon S3、Amazon Redshift、JDBC など)。有効な値には、`s3`、`mysql`、`postgresql`、`redshift`、`sqlserver`、`oracle`、`kinesis`、`kafka` が含まれています。
- `connection_options` - 接続オプション (パスやデータベーステーブルなど) (オプション)。`s3` の `connection_type` では、Amazon S3 パスが定義されています。

```
connection_options = {"path": "s3://aws-glue-target/temp"}
```

JDBC 接続の場合、いくつかのプロパティを定義する必要があります。データベース名は URL の一部である必要があることに注意してください。オプションで接続オプションに含めることができます。

Warning

スクリプトにパスワードを保存することはお勧めしません。boto3 を使用して AWS Secrets Manager または AWS Glue データカタログからそれらを取得することを検討してください。

```
connection_options = {"url": "jdbc-url/database", "user": "username",  
"password": passwordVariable, "dbtable": "table-name", "redshiftTmpDir": "s3-tempdir-  
path"}
```

dbtable プロパティは JDBC テーブルの名前です。データベース内でスキーマをサポートする JDBC データストアの場合、`schema.table-name` を指定します。スキーマを指定しない場合、デフォルトの「パブリック」スキーマが使用されます。

詳細については、「[AWS Glue for Spark での ETL の接続タイプとオプション](#)」を参照してください。

- `format` – 形式の仕様。Amazon S3 や、複数のフォーマットをサポートする AWS Glue 接続の場合に使用します。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `format_options` – 指定された形式についてのオプション。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。

`write_from_options`

```
write_from_options(frame_or_dfc, connection_type, connection_options={},  
format={}, format_options={}, transformation_ctx = "")
```

指定された接続および形式情報で作成された `DynamicFrame` または `DynamicFrameCollection` を書き込み、返します。

- `frame_or_dfc` - 書き込む `DynamicFrame` または `DynamicFrameCollection`。
- `connection_type` – 接続タイプ (Amazon S3、Amazon Redshift、JDBC など)。有効な値には、`s3`、`mysql`、`postgresql`、`redshift`、`sqlserver`、および `oracle` があります。
- `connection_options` - 接続オプション (パスやデータベーステーブルなど) (オプション)。s3 の `connection_type` では、Amazon S3 パスが定義されています。

```
connection_options = {"path": "s3://aws-glue-target/temp"}
```

JDBC 接続の場合、いくつかのプロパティを定義する必要があります。データベース名は URL の一部である必要があることに注意してください。オプションで接続オプションに含めることができます。

⚠ Warning

スクリプトにパスワードを保存することはお勧めしません。boto3 を使用して AWS Secrets Manager または AWS Glue データカタログからそれらを取得することを検討してください。

```
connection_options = {"url": "jdbc-url/database", "user": "username",  
"password": passwordVariable, "dbtable": "table-name", "redshiftTmpDir": "s3-tempdir-path"}
```

dbtable プロパティは JDBC テーブルの名前です。データベース内でスキーマをサポートする JDBC データストアの場合、schema.table-name を指定します。スキーマを指定しない場合、デフォルトの「パブリック」スキーマが使用されます。

詳細については、「[AWS Glue for Spark での ETL の接続タイプとオプション](#)」を参照してください。

- format – 形式の仕様。Amazon S3 や、複数のフォーマットをサポートする AWS Glue 接続の場合に使用します。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- format_options – 指定された形式についてのオプション。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- transformation_ctx - 使用する変換コンテキスト (オプション)。

write_dynamic_frame_from_catalog

```
write_dynamic_frame_from_catalog(frame, database, table_name,  
redshift_tmp_dir, transformation_ctx = "", additional_options = {},  
catalog_id = None)
```

データカタログデータベースとテーブルからの情報を使用して、記述した DynamicFrame を返します。

- frame - 書き込む DynamicFrame。
- Database – テーブルを含むデータカタログデータベース。
- table_name – ターゲットに関連付けられたデータカタログテーブルの名前。

- `redshift_tmp_dir` - 使用する Amazon Redshift の一時ディレクトリ (オプション)。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。
- `additional_options` - オプションの名前と値のペアのコレクション。
- `catalog_id` - 現在アクセスされているデータカタログのカタログ ID (アカウント ID)。None の場合は、呼び出し元のアカウント ID のデフォルトが使用されます。

`write_data_frame_from_catalog`

```
write_data_frame_from_catalog(frame, database, table_name,  
redshift_tmp_dir, transformation_ctx = "", additional_options = {},  
catalog_id = None)
```

データカタログデータベースとテーブルからの情報を使用して、記述した DataFrame を返します。このメソッドは、データレイク形式 (Hudi、Iceberg、および Delta Lake) への書き込みをサポートします。詳細については、「[AWS Glue ETL ジョブでのデータレイクフレームワークの使用](#)」を参照してください。

- `frame` - 書き込む DataFrame。
- `Database` - テーブルを含むデータカタログデータベース。
- `table_name` - ターゲットに関連付けられたデータカタログテーブルの名前。
- `redshift_tmp_dir` - 使用する Amazon Redshift の一時ディレクトリ (オプション)。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。
- `additional_options` - オプションの名前と値のペアのコレクション。
 - `useSparkDataSink` - true に設定すると、AWS Glue はネイティブ Spark Data Sink API を使用してテーブルに書き込むように強制されます。このオプションを有効にすると、`additional_options` 必要に応じて任意の [Spark データソースオプション](#) を追加できます。AWS Glue はこれらのオプションを Spark ライターに直接渡します。
- `catalog_id` - 現在アクセスされているデータカタログのカタログ ID (アカウント ID)。値を指定しない場合、呼び出し元のアカウント ID が使用されます。

制約事項

`useSparkDataSink` オプションを使用する際には、次の制限事項を考慮してください。

- `useSparkDataSink` オプションを使用する場合、[enableUpdateCatalog](#) オプションはサポートされません。

例: Spark データソースライターを使用して Hudi テーブルに書き込む

```
hudi_options = {
    'useSparkDataSink': True,
    'hoodie.table.name': <table_name>,
    'hoodie.datasource.write.storage.type': 'COPY_ON_WRITE',
    'hoodie.datasource.write.recordkey.field': 'product_id',
    'hoodie.datasource.write.table.name': <table_name>,
    'hoodie.datasource.write.operation': 'upsert',
    'hoodie.datasource.write.precombine.field': 'updated_at',
    'hoodie.datasource.write.hive_style_partitioning': 'true',
    'hoodie.upsert.shuffle.parallelism': 2,
    'hoodie.insert.shuffle.parallelism': 2,
    'hoodie.datasource.hive_sync.enable': 'true',
    'hoodie.datasource.hive_sync.database': <database_name>,
    'hoodie.datasource.hive_sync.table': <table_name>,
    'hoodie.datasource.hive_sync.use_jdbc': 'false',
    'hoodie.datasource.hive_sync.mode': 'hms'}

glueContext.write_data_frame.from_catalog(
    frame = <df_product_inserts>,
    database = <database_name>,
    table_name = <table_name>,
    additional_options = hudi_options
)
```

write_dynamic_frame_from_jdbc_conf

```
write_dynamic_frame_from_jdbc_conf(frame, catalog_connection,  
connection_options={}, redshift_tmp_dir = "", transformation_ctx = "",  
catalog_id = None)
```

指定された JDBC 接続情報を使用して DynamicFrame を書き込み、返します。

- frame - 書き込む DynamicFrame。
- catalog_connection - 使用するカタログ接続。
- connection_options - 接続オプション (パスやデータベーステーブルなど) (オプション)。詳しくは、「[AWS Glue for Spark での ETL の接続タイプとオプション](#)」を参照してください。
- redshift_tmp_dir - 使用する Amazon Redshift の一時ディレクトリ (オプション)。
- transformation_ctx - 使用する変換コンテキスト (オプション)。

- `catalog_id` – 現在アクセスされているデータカタログのカタログ ID (アカウント ID)。None の場合は、呼び出し元のアカウント ID のデフォルトが使用されます。

`write_from_jdbc_conf`

```
write_from_jdbc_conf(frame_or_dfc, catalog_connection,  
connection_options={}, redshift_tmp_dir = "", transformation_ctx = "",  
catalog_id = None)
```

指定された JDBC 接続情報を使用して `DynamicFrame` または `DynamicFrameCollection` を書き込み、返します。

- `frame_or_dfc` - 書き込む `DynamicFrame` または `DynamicFrameCollection`。
- `catalog_connection` - 使用するカタログ接続。
- `connection_options` - 接続オプション (パスやデータベーステーブルなど) (オプション)。詳しくは、「[AWS Glue for Spark での ETL の接続タイプとオプション](#)」を参照してください。
- `redshift_tmp_dir` - 使用する Amazon Redshift の一時ディレクトリ (オプション)。
- `transformation_ctx` - 使用する変換コンテキスト (オプション)。
- `catalog_id` – 現在アクセスされているデータカタログのカタログ ID (アカウント ID)。None の場合は、呼び出し元のアカウント ID のデフォルトが使用されます。

AWS Glue PySpark 変換リファレンス

AWS Glue には、PySpark ETL オペレーションで使用できる以下の組み込み変換が用意されています。データは、Apache Spark SQL への拡張である と呼ばれるデータ構造で `DynamicFrame` 変換から変換に渡されます `DataFrame`。 `DynamicFrame` にはデータが含まれており、データを処理するためにそのスキーマを参照します。

これらの変換のほとんどは、`DynamicFrame` クラスのメソッドとしても存在します。詳細については、「[DynamicFrame 変換](#)」を参照してください。

- [GlueTransform 基本クラス](#)
- [ApplyMapping クラス](#)
- [DropFields クラス](#)
- [DropNullFields クラス](#)
- [ErrorsAsDynamicFrame クラス](#)

- [EvaluateDataQuality クラス](#)
- [FillMissingValues クラス](#)
- [フィルタクラス](#)
- [FindIncrementalMatches クラス](#)
- [FindMatches クラス](#)
- [FlatMap クラス](#)
- [Join クラス](#)
- [マップクラス](#)
- [MapToCollection クラス](#)
- [mergeDynamicFrame](#)
- [クラスの関連付け](#)
- [RenameField クラス](#)
- [ResolveChoice クラス](#)
- [SelectFields クラス](#)
- [SelectFromCollection クラス](#)
- [Simplify_ddb_json class](#)
- [スピゴットクラス](#)
- [SplitFields クラス](#)
- [SplitRows クラス](#)
- [Unbox クラス](#)
- [UnnestFrame クラス](#)

GlueTransform 基本クラス

すべての `awsglue.transforms` クラスが継承する基本クラス。

クラスはすべて `__call__` メソッドを定義します。次のセクションにリストされている GlueTransform クラスのメソッドを上書きするか、デフォルトでクラス名を使用して呼び出されます。

方法

- [apply\(cls, *args, **kwargs\)](#)

- [name\(cls\)](#)
- [describeArgs\(cls\)](#)
- [describeReturn\(cls\)](#)
- [describeTransform\(cls\)](#)
- [describeErrors\(cls\)](#)
- [describe\(cls\)](#)

`apply(cls, *args, **kwargs)`

変換クラスを呼び出して変換を適用し、結果を返します。

- `cls` - `self` クラスオブジェクト。

`name(cls)`

派生変換クラスの名前を返します。

- `cls` - `self` クラスオブジェクト。

`describeArgs(cls)`

- `cls` - `self` クラスオブジェクト。

名前付き引数にそれぞれ対応する辞書のリストを次の形式で返します。

```
[
  {
    "name": "(name of argument)",
    "type": "(type of argument)",
    "description": "(description of argument)",
    "optional": "(Boolean, True if the argument is optional)",
    "defaultValue": "(Default value string, or None)(String; the default value, or None)"
  },
  ...
]
```

実装されていない派生変換で呼び出されたときに `NotImplementedError` 例外が発生します。

describeReturn(cls)

- cls - self クラスオブジェクト。

戻り型に関する情報を含む辞書を次の形式で返します。

```
{
  "type": "(return type)",
  "description": "(description of output)"
}
```

実装されていない派生変換で呼び出されたときに `NotImplementedError` 例外が発生します。

describeTransform(cls)

変換について説明する文字列を返します。

- cls - self クラスオブジェクト。

実装されていない派生変換で呼び出されたときに `NotImplementedError` 例外が発生します。

describeErrors(cls)

- cls - self クラスオブジェクト。

この変換によってスローされる可能性のある例外をそれぞれ説明する辞書のリストを、次の形式で返します。

```
[
  {
    "type": "(type of error)",
    "description": "(description of error)"
  },
  ...
]
```

describe(cls)

- cls - self クラスオブジェクト。

次の形式のオブジェクトを返します。

```
{
  "transform" : {
    "name" : cls.name( ),
    "args" : cls.describeArgs( ),
    "returns" : cls.describeReturn( ),
    "raises" : cls.describeErrors( ),
    "location" : "internal"
  }
}
```

ApplyMapping クラス

DynamicFrame でマッピングを適用します。

例

[DynamicFrame.apply_mapping\(\)](#) メソッドを使用して、DynamicFrame でマッピングを適用することをお勧めします。コード例については、「[例: apply_mapping を使用してフィールドの名前を変更し、フィールドタイプを変更する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(frame, mappings, transformation_ctx = "", info = "", stageThreshold = 0, totalThreshold = 0)`

指定された DynamicFrame に宣言型のマッピングを適用します。

- `frame` – マッピングを適用するための DynamicFrame (必須)。

- `mappings` — マッピングタプルのリスト (必須)。それぞれが (ソース列、ソースタイプ、ターゲット列、ターゲットタイプ) で構成されます。

ソース列で、名前にドット「`.`」が含まれている場合、バックティック「```」で囲む必要があります。例えば、`this.old.name` (文字列) を `thisNewName` にマッピングするには、次のタプルを使用します。

```
("`this.old.name`", "string", "thisNewName", "string")
```

- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` - エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` - エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

「マッピング」タプルで指定された `DynamicFrame` のフィールドのみを返します。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

DropFields クラス

DynamicFrame 内のフィールドを削除します。

例

[DynamicFrame.drop_fields\(\)](#) メソッドを使用して、DynamicFrame からフィールドを削除することをお勧めします。コード例については、「[例: drop_fields を使用して DynamicFrame からフィールドを削除する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, paths, transformation_ctx = "", info = "", stageThreshold = 0, totalThreshold = 0)
```

DynamicFrame 内のノードを削除します。

- frame – ノードを削除する DynamicFrame (必須)。
- paths – 削除するノードへの完全パスのリスト (必須)。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - 変換のエラーに関連付けられた文字列 (オプション)。
- stageThreshold – エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- totalThreshold – エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

指定したフィールドを除く新しい `DynamicFrame` を返します。

```
apply(cls, *args, **kwargs)
```

継承元は `GlueTransform` [apply](#)。

```
name(cls)
```

継承元は `GlueTransform` [name](#)。

```
describeArgs(cls)
```

継承元は `GlueTransform` [describeArgs](#)。

```
describeReturn(cls)
```

継承元は `GlueTransform` [describeReturn](#)。

```
describeTransform(cls)
```

継承元は `GlueTransform` [describeTransform](#)。

```
describeErrors(cls)
```

継承元は `GlueTransform` [describeErrors](#)。

```
describe(cls)
```

継承元は `GlueTransform` [説明](#)。

DropNullFields クラス

`DynamicFrame` でタイプが `NullType` のすべての null フィールドを削除します。これらは、`DynamicFrame` データセット内のすべてのレコードに、値がないか null 値があるフィールドです。

例

この例では、`DropNullFields` を使用して、タイプ `NullType` のフィールドが削除された新しい `DynamicFrame` を作成します。`DropNullFields` を示すため、null 型の `empty_column` という名前の新しい列を、ロード済みの `persons` データセットに追加します。

Note

この例で使用されているデータセットにアクセスするには、「[コード例: データの結合と関係付け](#)」を参照し、「[ステップ 1: Amazon S3 バケット内のデータをクローलする](#)」の手順に従います。

```
# Example: Use DropNullFields to create a new DynamicFrame without NullType fields

from pyspark.context import SparkContext
from awsglue.context import GlueContext
from pyspark.sql.functions import lit
from pyspark.sql.types import NullType
from awsglue.dynamicframe import DynamicFrame
from awsglue.transforms import DropNullFields

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Create DynamicFrame
persons = glueContext.create_dynamic_frame.from_catalog(
    database="legislators", table_name="persons_json"
)
print("Schema for the persons DynamicFrame:")
persons.printSchema()

# Add new column "empty_column" with NullType
persons_with_nulls = persons.toDF().withColumn("empty_column",
    lit(None).cast(NullType()))
persons_with_nulls_dyf = DynamicFrame.fromDF(persons_with_nulls, glueContext,
    "persons_with_nulls")
print("Schema for the persons_with_nulls_dyf DynamicFrame:")
persons_with_nulls_dyf.printSchema()

# Remove the NullType field
persons_no_nulls = DropNullFields.apply(persons_with_nulls_dyf)
print("Schema for the persons_no_nulls DynamicFrame:")
persons_no_nulls.printSchema()
```

出力

Schema for the persons DynamicFrame:

```
root
|-- family_name: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- gender: string
|-- image: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- sort_name: string
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- death_date: string
```

Schema for the persons_with_nulls_dyf DynamicFrame:

```
root
|-- family_name: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- gender: string
|-- image: string
```

```
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- sort_name: string
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- death_date: string
|-- empty_column: null
```

```
null_fields ['empty_column']
```

```
Schema for the persons_no_nulls DynamicFrame:
```

```
root
```

```
|-- family_name: string
|-- name: string
|-- links: array
|   |-- element: struct
|   |   |-- note: string
|   |   |-- url: string
|-- gender: string
|-- image: string
|-- identifiers: array
|   |-- element: struct
|   |   |-- scheme: string
|   |   |-- identifier: string
|-- other_names: array
|   |-- element: struct
|   |   |-- lang: string
|   |   |-- note: string
|   |   |-- name: string
|-- sort_name: string
```

```
|-- images: array
|   |-- element: struct
|   |   |-- url: string
|-- given_name: string
|-- birth_date: string
|-- id: string
|-- contact_details: array
|   |-- element: struct
|   |   |-- type: string
|   |   |-- value: string
|-- death_date: string
```

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(frame, transformation_ctx = "", info = "", stageThreshold = 0, totalThreshold = 0)`

DynamicFrame でタイプが NullType のすべての null フィールドを削除します。これらは、DynamicFrame データセット内のすべてのレコードに、値がないか null 値があるフィールドです。

- `frame` – null フィールドを削除する DynamicFrame (必須)。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` – エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` – エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

null フィールドのない新しい `DynamicFrame` を返します。

`apply(cls, *args, **kwargs)`

- `cls` – `cls`

`name(cls)`

- `cls` – `cls`

`describeArgs(cls)`

- `cls` – `cls`

`describeReturn(cls)`

- `cls` – `cls`

`describeTransform(cls)`

- `cls` – `cls`

`describeErrors(cls)`

- `cls` – `cls`

`describe(cls)`

- `cls` – `cls`

`ErrorsAsDynamicFrame` クラス

ソース `DynamicFrame` の作成中に発生したエラーのネストされたレコードを含む `DynamicFrame` を返します。

例

[DynamicFrame.errorsAsDynamicFrame\(\)](#) メソッドを使用して、エラーレコードを取得および表示することをお勧めします。コード例については、「[例: errorsAsDynamicFrame を使用してエラーレコードを表示する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(frame)`

ソース DynamicFrame に関連する、ネストされたエラーレコードを含む DynamicFrame を返します。

- `frame` – ソース DynamicFrame (必須)。

`apply(cls, *args, **kwargs)`

- `cls` – cls

`name(cls)`

- `cls` – cls

`describeArgs(cls)`

- `cls` – cls

describeReturn(cls)

- cls – cls

describeTransform(cls)

- cls – cls

describeErrors(cls)

- cls – cls

describe(cls)

- cls – cls

EvaluateDataQuality クラス

DynamicFrame に対してデータ品質ルールセットを評価し、評価結果を含む新しい DynamicFrame を返します。

例

次のコード例は、DynamicFrame のデータ品質を評価し、データ品質結果を表示する方法を示しています。

```
from awsglue.transforms import *
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsgluedq.transforms import EvaluateDataQuality

#Create Glue context
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Define DynamicFrame
legislatorsAreas = glueContext.create_dynamic_frame.from_catalog(
    database="legislators", table_name="areas_json")

# Create data quality ruleset
```

```
ruleset = """"Rules = [ColumnExists "id", IsComplete "id"]""""

# Evaluate data quality
dqResults = EvaluateDataQuality.apply(
  frame=legislatorsAreas,
  ruleset=ruleset,
  publishing_options={
    "dataQualityEvaluationContext": "legislatorsAreas",
    "enableDataQualityCloudWatchMetrics": True,
    "enableDataQualityResultsPublishing": True,
    "resultsS3Prefix": "DOC-EXAMPLE-BUCKET1",
  },
)

# Inspect data quality results
dqResults.printSchema()
dqResults.toDF().show()
```

出力

```
root
 |-- Rule: string
 |-- Outcome: string
 |-- FailureReason: string
 |-- EvaluatedMetrics: map
 |   |-- keyType: string
 |   |-- valueType: double

+-----+-----+-----+-----+
|Rule          |Outcome|FailureReason|EvaluatedMetrics          |
+-----+-----+-----+-----+
|ColumnExists "id"  |Passed |null         |{}                          |
|IsComplete "id"   |Passed |null         |{Column.first_name.Completeness -> 1.0}|
+-----+-----+-----+-----+
```

方法

- [call](#)
- [apply](#)
- [name](#)

- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, ruleset, publishing_options = {})
```

- `frame` – データ品質を評価したい `DynamicFrame`。
- `ruleset` – 文字列形式のデータ品質定義言語 (DQDL) ルールセット。DQDL の詳細については、[データ品質定義言語 \(DQDL\) リファレンス](#) のガイドを参照してください。
- `publishing_options` – 評価結果とメトリクスを発行する次のオプションを指定するディクショナリ。
 - `dataQualityEvaluationContext` – AWS Glue が Amazon CloudWatch メトリクスとデータ品質結果を発行する名前空間を指定する文字列。集計されたメトリクスは CloudWatch に表示され、完全な結果は AWS Glue Studio インターフェイスに表示されます。
 - 必須: いいえ
 - デフォルト値: `default_context`
 - `enableDataQualityCloudWatchMetrics` – データ品質評価の結果を CloudWatch に発行するかどうかを指定します。 `dataQualityEvaluationContext` オプションを使用してメトリクスの名前空間を指定します。
 - 必須: いいえ
 - デフォルト値: `False`
 - `enableDataQualityResultsPublishing` – データ品質結果を AWS Glue Studio インターフェイスの [Data Quality] (データ品質) タブに表示するかどうかを指定します。
 - 必須: いいえ
 - デフォルト値: `True`
 - `resultsS3Prefix` – AWS Glue がデータ品質評価結果を書き込める Amazon S3 ロケーションを指定します。
 - 必須: いいえ
 - デフォルト値: `""` (空の文字列)

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

FillMissingValues クラス

FillMissingValues クラスでは、指定された `DynamicFrame` の中で Null 値や空の文字列を特定し、欠落した値を推定するために、線形回帰やランダムフォレストなどの機械学習手法を使用します。ETL ジョブは、入力データセットの値を使用して機械学習モデルをトレーニングします。その後、欠落した値が予測されます。

Tip

増分データセットを使用する場合、各増分セットが機械学習モデルのトレーニングデータとして使用されます。したがって、結果から正確性が損なわれることがあります。

インポートの対象:

```
from awsglueml.transforms import FillMissingValues
```

方法

- [適用](#)

```
apply(frame, missing_values_column, output_column = "", transformation_ctx = "", info = "",
stageThreshold = 0, totalThreshold = 0)
```

指定された列の動的フレームの欠落値を埋め、推定値を含む新しい列を持つ新しいフレームを返します。欠落した値がない行の場合、指定した列の値が新しい列に複製されます。

- `frame` - 欠落値を埋める `DynamicFrame`。必須。
- `missing_values_column` - 欠落値を含む列 (`null` 値や空の文字列)。必須。
- `output_column` - 値が欠落しているすべての行に推定値が埋められた新しい列の名前。デフォルトは、末尾に `"_filled"` が付いた `missing_values_column` の名前です。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` - エラーが発生する前に変換で発生する可能性のあるエラーの最大数 (オプション、デフォルト値は 0)。
- `totalThreshold` - 処理がエラーを出す前に全体的に発生する可能性のあるエラーの最大数 (オプション、デフォルト値は 0)。

欠落値を持つ行の推定値と、他の行の現在値を含む、追加的な列を持つ新しい `DynamicFrame` を返します。

フィルタクラス

指定された述語関数を満たす、入力 `DynamicFrame` からのレコードが含まれた新しい `DynamicFrame` を構築します。

例

[DynamicFrame.filter\(\)](#) メソッドを使用して、`DynamicFrame` のレコードをフィルタリングすることをお勧めします。コード例については、「[例: フィルターを使用して、フィールドのフィルタリングされた選択内容を取得する](#)」を参照してください。

方法

- [__call__](#)

- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, f, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0))
```

指定された述語関数を満たす入力 DynamicFrame からのレコードを選択することにより構築された新しい DynamicFrame を返します。

- `frame` – 指定されたフィルタ関数を適用する先のソース DynamicFrame (必須)。
- `f` – DynamicFrame のそれぞれの DynamicRecord に適用する述語関数。この関数は、引数として DynamicRecord を取り、DynamicRecord がフィルタリングの要件を満たす場合は True を返し、そうでない場合は False を返す必要があります (必須)。

DynamicRecord は DynamicFrame 内の論理レコードを表します。これは、自己記述型であり、固定スキーマに準拠しないデータに使用できる点を除いて、Spark DataFrame の行に似ています。

- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` – エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` – エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

```
apply(cls, *args, **kwargs)
```

継承元は GlueTransform [apply](#)。

```
name(cls)
```

継承元は GlueTransform [name](#)。

describeArgs(cls)

継承元は GlueTransform [describeArgs](#)。

describeReturn(cls)

継承元は GlueTransform [describeReturn](#)。

describeTransform(cls)

継承元は GlueTransform [describeTransform](#)。

describeErrors(cls)

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

FindIncrementalMatches クラス

既存と増分の DynamicFrame 内で一致するレコードを識別します。さらに、一致するレコードの各グループに割り当てられる一意の識別子を使用して新しい DynamicFrame を作成します。

インポートの対象:

```
from awsglueml.transforms import FindIncrementalMatches
```

方法

- [適用](#)

```
apply(existingFrame, incrementalFrame, transformId, transformation_ctx = "", info = "",  
stageThreshold = 0, totalThreshold = 0, enforcedMatches = None, computeMatchConfidenceScores  
= 0)
```

入力の DynamicFrame 内で一致するレコードを特定し、そのレコードの各グループに割り当てられている一意の識別子を含む新しい DynamicFrame を作成します。

- existingFrame – FindIncrementalMatches 変換を適用する、既存の、および事前一致済みの DynamicFrame 。必須。

- `incrementalFrame - existingFrame` に一致させるために `FindIncrementalMatches` 変換を適用する増分 `DynamicFrame`。必須。
- `transformId` - `DynamicFrames` 内のレコードに適用するために `FindIncrementalMatches` 変換に関連付けられた一意の ID。必須。
- `transformation_ctx` - 統計/ステータス情報を識別するために使用される一意の文字列。オプション。
- `info` - 変換のエラーに関連付けられる文字列。オプション。
- `stageThreshold` - エラーで終了する前に、変換で許容されるエラーの最大発生数。オプション。デフォルト値は 0 です。
- `totalThreshold` - エラーで終了する前に、処理で全体的に許容されるエラーの最大発生数。オプション。デフォルト値は 0 です。
- `enforcedMatches` - 一致を強制するために使用される `DynamicFrame`。オプション。デフォルトは [なし] です。
- `computeMatchConfidenceScores` - 一致するレコードの各グループの信頼スコアをコンピューティングするかどうかを示すブール値。オプション。デフォルトは `False` です。

一致するレコードの各グループに割り当てられている一意の識別子を含む、新しい `DynamicFrame` を返します。

FindMatches クラス

入力の `DynamicFrame` 内で一致するレコードを特定し、そのレコードの各グループに割り当てられている一意の識別子を含む新しい `DynamicFrame` を作成します。

インポートの対象:

```
from awsglueml.transforms import FindMatches
```

方法

- [適用](#)

```
apply(frame, transformId, transformation_ctx = "", info = "", stageThreshold = 0, totalThreshold = 0, enforcedMatches = None, computeMatchConfidenceScores = 0)
```

入力の `DynamicFrame` 内で一致するレコードを特定し、そのレコードの各グループに割り当てられている一意の識別子を含む新しい `DynamicFrame` を作成します。

- `frame` – FindMatches 変換を適用する DynamicFrame。必須。
- `transformId` – DynamicFrame 内のレコードに適用するために FindMatches 変換に関連付けられた一意の ID。必須。
- `transformation_ctx` – 統計/ステータス情報を識別するために使用される一意の文字列。オプション。
- `info` – 変換のエラーに関連付けられる文字列。オプション。
- `stageThreshold` – エラーで終了する前に、変換で許容されるエラーの最大発生数。オプション。デフォルト値は 0 です。
- `totalThreshold` – エラーで終了する前に、処理で全体的に許容されるエラーの最大発生数。オプション。デフォルト値は 0 です。
- `enforcedMatches` – 一致を強制するために使用される DynamicFrame。オプション。デフォルトは [なし] です。
- `computeMatchConfidenceScores` – 一致するレコードの各グループの信頼スコアをコンピューティングするかどうかを示すブール値。オプション。デフォルトは False です。

一致するレコードの各グループに割り当てられている一意の識別子を含む、新しい DynamicFrame を返します。

FlatMap クラス

コレクションの各 DynamicFrame に変換が適用されます。結果が単一の DynamicFrame にフラット化されるのではなく、コレクションとして保存されます。

FlatMap の例

次のスニペットの例は、FlatMap に適用したときに、ResolveChoice 変換を動的フレームのコレクションで使用方法を示しています。入力に使用されるデータは、プレースホルダ Amazon S3 アドレス `s3://bucket/path-for-data/sample.json` にある JSON の中であって、次のデータが含まれています。

JSON データの例

```
[{
  "firstname": "Arnav",
  "lastname": "Desai",
  "address": {
    "street": "6 Anyroad Avenue",
    "city": "London",
```

```
        "state": "England",
        "country": "UK"
    },
    "phone": 17235550101,
    "affiliations": [
        "General Anonymous Example Products",
        "Example Independent Research",
        "Government Department of Examples"
    ]
},
{
    "firstname": "Mary",
    "lastname": "Major",
    "address": {
        "street": "7821 Spot Place",
        "city": "Centerville",
        "state": "OK",
        "country": "US"
    },
    "phone": 19185550023,
    "affiliations": [
        "Example Dot Com",
        "Example Independent Research",
        "Example.io"
    ]
},
{
    "firstname": "Paulo",
    "lastname": "Santos",
    "address": {
        "street": "123 Maple Street",
        "city": "London",
        "state": "Ontario",
        "country": "CA"
    },
    "phone": 12175550181,
    "affiliations": [
        "General Anonymous Example Products",
        "Example Dot Com"
    ]
}
]]
```

Example ResolveChoice を DynamicFrameCollection に適用し、出力を表示します。

```
#Read DynamicFrame
datasource = glueContext.create_dynamic_frame_from_options("s3", connection_options =
  {"paths":["s3://bucket/path/to/file/mysamplejson.json"]}, format="json")
datasource.printSchema()
datasource.show()

## Split to create a DynamicFrameCollection
split_frame=datasource.split_fields(["firstname","lastname","address"],"personal_info","business_info")
split_frame.keys()
print("---")

## Use FlatMap to run ResolveChoice
kwargs = {"choice": "cast:string"}
flat = FlatMap.apply(split_frame, ResolveChoice, frame_name="frame",
  transformation_ctx='tcx', **kwargs)
flat.keys()

##Select one of the DynamicFrames
personal_info = flat.select("personal_info")
personal_info.printSchema()
personal_info.show()
print("---")

business_info = flat.select("business_info")
business_info.printSchema()
business_info.show()
```

Important

FlatMap.apply を呼び出すとき、frame_name パラメータは "frame" である必要があります。現在、他の値は受け付けられません。

出力例

```
root
|-- firstname: string
|-- lastname: string
|-- address: struct
|   |-- street: string
```

```
|   |-- city: string
|   |-- state: string
|   |-- country: string
|-- phone: long
|-- affiliations: array
|   |-- element: string
---
{
  "firstname": "Mary",
  "lastname": "Major",
  "address": {
    "street": "7821 Spot Place",
    "city": "Centerville",
    "state": "OK",
    "country": "US"
  },
  "phone": 19185550023,
  "affiliations": [
    "Example Dot Com",
    "Example Independent Research",
    "Example.io"
  ]
}

{
  "firstname": "Paulo",
  "lastname": "Santos",
  "address": {
    "street": "123 Maple Street",
    "city": "London",
    "state": "Ontario",
    "country": "CA"
  },
  "phone": 12175550181,
  "affiliations": [
    "General Anonymous Example Products",
    "Example Dot Com"
  ]
}
---
root
|-- firstname: string
|-- lastname: string
|-- address: struct
```

```
|   |-- street: string
|   |-- city: string
|   |-- state: string
|   |-- country: string

{
  "firstname": "Mary",
  "lastname": "Major",
  "address": {
    "street": "7821 Spot Place",
    "city": "Centerville",
    "state": "OK",
    "country": "US"
  }
}

{
  "firstname": "Paulo",
  "lastname": "Santos",
  "address": {
    "street": "123 Maple Street",
    "city": "London",
    "state": "Ontario",
    "country": "CA"
  }
}

---
root
|-- phone: long
|-- affiliations: array
|   |-- element: string

{
  "phone": 19185550023,
  "affiliations": [
    "Example Dot Com",
    "Example Independent Research",
    "Example.io"
  ]
}

{
  "phone": 12175550181,
  "affiliations": [
```

```
        "General Anonymous Example Products",  
        "Example Dot Com"  
    ]  
}
```

方法

- [__call__](#)
- [適用](#)
- [名前](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(dfc, BaseTransform, frame_name, transformation_ctx = "", **base_kwargs)`

変換をコレクション内の各 DynamicFrame に適用し、結果をフラット化します。

- `dfc` – フラットマップする DynamicFrameCollection (必須)。
- `BaseTransform` – コレクションの各メンバーに適用する GlueTransform から派生した変換 (必須)。
- `frame_name` – コレクションの要素を渡す引数名 (必須)。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `base_kwargs` – ベース変換に渡す引数 (必須)。

ソース DynamicFrameCollection の各 DynamicFrame に変換を適用して作成された新しい DynamicFrameCollection を返します。

`apply(cls, *args, **kwargs)`

継承元は GlueTransform [apply](#)。

`name(cls)`

継承元は GlueTransform [name](#)。

describeArgs(cls)

継承元は GlueTransform [describeArgs](#)。

describeReturn(cls)

継承元は GlueTransform [describeReturn](#)。

describeTransform(cls)

継承元は GlueTransform [describeTransform](#)。

describeErrors(cls)

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

Join クラス

2 つの DynamicFrames の等値結合を実行します。

例

DynamicFrames を結合するには、[DynamicFrame.join\(\)](#) を使用することをお勧めします。コード例については、「[例: 結合を使用して DynamicFrames を結合する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame1, frame2, keys1, keys2, transformation_ctx = "")
```

2 つの DynamicFrames の等値結合を実行します。

- frame1 - 結合する最初の DynamicFrame (必須)。
- frame2 - 結合する 2 番目の DynamicFrame (必須)。
- keys1 - 最初のフレームで結合するキー (必須)。
- keys2 - 2 番目のフレームで結合するキー (必須)。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。

2 つの DynamicFrames を結合することで作成された新しい DynamicFrame を返します。

```
apply(cls, *args, **kwargs)
```

継承元は GlueTransform [apply](#)。

```
name(cls)
```

継承元は GlueTransform [name](#)。

```
describeArgs(cls)
```

継承元は GlueTransform [describeArgs](#)。

```
describeReturn(cls)
```

継承元は GlueTransform [describeReturn](#)。

```
describeTransform(cls)
```

継承元は GlueTransform [describeTransform](#)。

```
describeErrors(cls)
```

継承元は GlueTransform [describeErrors](#)。

```
describe(cls)
```

継承元は GlueTransform [説明](#)。

マップクラス

入力 `DynamicFrame` ですべてのレコードに関数を適用して、新しい `DynamicFrame` をビルドします。

例

[DynamicFrame.map\(\)](#) メソッドを使用して、`DynamicFrame` のすべてのレコードに関数を適用することをお勧めします。コード例については、「[例: マッピングを使用して、DynamicFrame のすべてのレコードに関数を適用する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, f, transformation_ctx="", info="", stageThreshold=0, totalThreshold=0)
```

指定された関数を元の `DynamicFrame` で、すべての `DynamicRecords` に適用した結果の新しい `DynamicFrame` を返します。

- `frame` - マッピング関数を適用するための元の `DynamicFrame` (必須)。
- `f` - `DynamicRecords` 内のすべての `DynamicFrame` に適用する関数。この関数は、引数として `DynamicRecord` を取り、マッピングによって生成された新しい `DynamicRecord` を返す必要があります (必須)。

`DynamicRecord` は `DynamicFrame` 内の論理レコードを表します。これは、自己記述型であり、固定スキーマに準拠しないデータに使用できる点を除いて、Apache Spark `DataFrame` の行に似ています。

- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。

- `stageThreshold` – エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` – エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

指定された関数を元の `DynamicFrame` で、すべての `DynamicRecords` に適用した結果の新しい `DynamicFrame` を返します。

```
apply(cls, *args, **kwargs)
```

継承元は `GlueTransform` [apply](#)。

```
name(cls)
```

継承元は `GlueTransform` [name](#)。

```
describeArgs(cls)
```

継承元は `GlueTransform` [describeArgs](#)。

```
describeReturn(cls)
```

継承元は `GlueTransform` [describeReturn](#)。

```
describeTransform(cls)
```

継承元は `GlueTransform` [describeTransform](#)。

```
describeErrors(cls)
```

継承元は `GlueTransform` [describeErrors](#)。

```
describe(cls)
```

継承元は `GlueTransform` [説明](#)。

MapToCollection クラス

指定された `DynamicFrameCollection` の各 `DynamicFrame` に変換が適用されます。

方法

- [__call__](#)

- [適用](#)
- [名前](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(dfc, BaseTransform, frame_name, transformation_ctx = "", **base_kwargs)
```

指定された DynamicFrameCollection の各 DynamicFrame に変換関数が適用されます。

- `dfc` - 変換関数を適用する DynamicFrameCollection (必須)。
- `callable` - コレクションの各メンバーに適用するコール可能な変換関数 (必須)。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。

ソース DynamicFrameCollection の各 DynamicFrame に変換を適用して作成された新しい DynamicFrameCollection を返します。

```
apply(cls, *args, **kwargs)
```

GlueTransform [apply](#) から継承されました。

```
name(cls)
```

継承元は GlueTransform [name](#)。

```
describeArgs(cls)
```

継承元は GlueTransform [describeArgs](#)。

```
describeReturn(cls)
```

継承元は GlueTransform [describeReturn](#)。

```
describeTransform(cls)
```

継承元は GlueTransform [describeTransform](#)。

`describeErrors(cls)`

継承元は GlueTransform [describeErrors](#)。

`describe(cls)`

継承元は GlueTransform [説明](#)。

クラスの関連付け

DynamicFrame のネストされたスキーマをフラット化し、フラット化されたフレームから配列の列をピボットアウトします。

例

[DynamicFrame.relationalize\(\)](#) メソッドを使用して、DynamicFrame を関係付けすることをお勧めします。コード例については、「[例: relationalize を使用して、DynamicFrame のネストされたスキーマをフラット化する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, staging_path=None, name='roottable', options=None, transformation_ctx = "", info = "", stageThreshold = 0, totalThreshold = 0)
```

DynamicFrame を関係付け、ネストされた列をネスト解除し、配列の列をピボットすることによってフレームのリストを生成します。unnest のフェーズで生成された結合キーを使用して、ピボットされた配列の列をルートテーブルに結合できます。

- frame – 関連付ける DynamicFrame (必須)。

- `staging_path` – このメソッドを使用して、ピボットされたテーブルのパーティションを CSV 形式で保存する先を示すパス (オプション)。ピボットされたテーブルはこのパスから読み取ります。
- `name` – ルートテーブルの名前 (オプション)。
- `options` - オプションのパラメータのディクショナリ。現在使用されていません。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` – エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` – エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

`RenameField` クラス

`DynamicFrame` 内のノードの名前を変更します。

例

[DynamicFrame.rename_field\(\)](#) メソッドを使用して、DynamicFrame のフィールドの名前を変更することをお勧めします。コード例については、「[例: rename_field を使用して、DynamicFrame のフィールドの名前を変更する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, old_name, new_name, transformation_ctx = "", info = "", stageThreshold = 0, totalThreshold = 0)
```

DynamicFrame 内のノードの名前を変更します。

- frame - ノードの名前を変更する DynamicFrame (必須)。
- old_name - 名前を変更するノードへのフルパス (必須)。

古い名前にドットが含まれている場合、RenameField はバックティック (``) で囲まなければ機能しません。たとえば、this.old.name を thisNewName に置き換えるには、RenameField を次のように呼び出します。

```
newDyF = RenameField(oldDyF, ``this.old.name``, "thisNewName")
```

- new_name - フルパスを含む新しい名前 (必須)。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - 変換のエラーに関連付けられた文字列 (オプション)。
- stageThreshold - エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

- `totalThreshold` – エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

ResolveChoice クラス

DynamicFrame 内で Choice 型を解決します。

例

[DynamicFrame.resolveChoice\(\)](#) メソッドを使用して、DynamicFrame の複数の型を含むフィールドを処理することをお勧めします。コード例については、「[例: resolveChoice を使用して、複数の型を含む列を処理する](#)」を参照してください。

方法

- [__call__](#)

- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, specs = None, choice = "", transformation_ctx = "", info = "", stageThreshold = 0, totalThreshold = 0)
```

DynamicFrame 内のあいまいな型を解決するための情報を提供します。作成された DynamicFrame を返します。

- `frame` – Choice 型を解決する DynamicFrame (必須)。
- `specs` – 解決する特定のあいまいさのリスト、それぞれがタプルの形式: (`path`, `action`)。 `path` 値は特定のあいまいな要素を識別し、`action` 値は対応する解決を識別します。

`spec` パラメータおよび `choice` パラメータのうち 1 つのみを使用できます。 `spec` パラメータが `None` ではない場合、`choice` パラメータは空の文字列である必要があります。逆に、`choice` が空の文字列ではない場合、`spec` パラメータは `None` である必要があります。どちらのパラメータも指定されていない場合、AWS Glue はスキーマを解析し、それを使用してあいまいさを解決します。

`specs` タプルの `action` 部分で、次の解決策のうちの 1 つを指定できます。

- `cast` – キャスト先の型を指定できます (例: `cast:int`)。
- `make_cols` – データをフラット化することで潜在的なあいまいさを解消します。たとえば、`columnA` が `int` または `string` の場合、解決策は、DynamicFrame に `columnA_int` および `columnA_string` という名前の 2 つの列を生成することです。
- `make_struct` – 構造体を使用してデータを表現することで、潜在的なあいまいさを解決します。たとえば、列のデータが `int` または `string` の場合、`make_struct` アクションを使用すると、作成された DynamicFrame に `int` および `string` の両方を含む構造体の列が生成されます。
- `project` – 生成された DynamicFrame で指定された型の値のみを保持して、潜在的なあいまいさを解決します。例えば、ChoiceType 列のデータが `int` または `string` となる可能性

がある場合、`project:string` アクションを指定すると、生成された `DynamicFrame` から `string` 型以外の値が削除されます。

`path` で配列を識別する場合は、あいまいさを避けるために配列名の後に空の角括弧を置きます。例えば、使用しているデータが次のように構造化されているとします。

```
"myList": [  
  { "price": 100.00 },  
  { "price": "$100.00" }  
]
```

`path` を `"myList[].price"` に設定し、`action` を `"cast:double"` に設定すると、文字列バージョンではなく、数値バージョンの料金を選択できます。

- `choice - specs` パラメータが `None` の場合のデフォルトの解決アクション。specs パラメータが `None` ではない場合、空の文字列以外に設定することはできません。

この引数では、上記の `specs` アクションに加えて、次のアクションもサポートされています。

- `MATCH_CATALOG` - 各 `ChoiceType` について、指定した Data Catalog テーブル内の対応する型へのキャストを試みます。
- `database` - `MATCH_CATALOG` 選択肢で使用する AWS Glue Data Catalog データベース (`MATCH_CATALOG` に必要)。
- `table_name` - `MATCH_CATALOG` アクションで使用する AWS Glue Data Catalog のテーブル名 (`MATCH_CATALOG` に必要)。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` - エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` - エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

```
apply(cls, *args, **kwargs)
```

継承元は `GlueTransform` [apply](#)。

```
name(cls)
```

継承元は `GlueTransform` [name](#)。

describeArgs(cls)

継承元は GlueTransform [describeArgs](#)。

describeReturn(cls)

継承元は GlueTransform [describeReturn](#)。

describeTransform(cls)

継承元は GlueTransform [describeTransform](#)。

describeErrors(cls)

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

SelectFields クラス

この SelectFields クラスは、既存の DynamicFrame から新しい DynamicFrame を作成し、指定したフィールドのみを保持します。SelectFields は、SQL SELECT ステートメントと同様の機能を提供します。

例

[DynamicFrame.select_fields\(\)](#) メソッドを使用して、DynamicFrame からフィールドを選択することをお勧めします。コード例については、「[例: select_fields を使用して、選択したフィールドで新しい DynamicFrame を作成する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)

- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, paths, transformation_ctx = "", info = "", stageThreshold = 0, totalThreshold = 0)
```

DynamicFrame のフィールド (ノード) を取得します。

- frame - フィールドを選択するための DynamicFrame (必須)。
- paths - 選択するフィールドへの完全パスのリスト (必須)。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。
- info - 変換のエラーに関連付けられた文字列 (オプション)。
- stageThreshold - エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- totalThreshold - エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

指定したフィールドのみを含む新しい DynamicFrame を返します。

```
apply(cls, *args, **kwargs)
```

継承元は GlueTransform [apply](#)。

```
name(cls)
```

継承元は GlueTransform [name](#)。

```
describeArgs(cls)
```

継承元は GlueTransform [describeArgs](#)。

```
describeReturn(cls)
```

継承元は GlueTransform [describeReturn](#)。

```
describeTransform(cls)
```

継承元は GlueTransform [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

SelectFromCollection クラス

`DynamicFrameCollection` で `DynamicFrame` を 1 つ選択します。

例

この例では、`SelectFromCollection` を使用して、`DynamicFrame` を `DynamicFrameCollection` から選択します。

データセットの例

この例では、`split_rows_collection` と呼ばれる `DynamicFrameCollection` から 2 つの `DynamicFrames` を選択します。次のリストに示しているのは、`split_rows_collection` のキーです。

```
dict_keys(['high', 'low'])
```

コードの例

```
# Example: Use SelectFromCollection to select
# DynamicFrames from a DynamicFrameCollection

from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.transforms import SelectFromCollection

# Create GlueContext
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

# Select frames and inspect entries
frame_low = SelectFromCollection.apply(dfc=split_rows_collection, key="low")
frame_low.toDF().show()

frame_high = SelectFromCollection.apply(dfc=split_rows_collection, key="high")
```

```
frame_high.toDF().show()
```

出力

```
+---+-----+-----+-----+-----+
| id|index|contact_details.val.type|contact_details.val.value|
+---+-----+-----+-----+-----+
| 1|  0|          fax|          202-225-3307|
| 1|  1|          phone|          202-225-5731|
| 2|  0|          fax|          202-225-3307|
| 2|  1|          phone|          202-225-5731|
| 3|  0|          fax|          202-225-3307|
| 3|  1|          phone|          202-225-5731|
| 4|  0|          fax|          202-225-3307|
| 4|  1|          phone|          202-225-5731|
| 5|  0|          fax|          202-225-3307|
| 5|  1|          phone|          202-225-5731|
| 6|  0|          fax|          202-225-3307|
| 6|  1|          phone|          202-225-5731|
| 7|  0|          fax|          202-225-3307|
| 7|  1|          phone|          202-225-5731|
| 8|  0|          fax|          202-225-3307|
| 8|  1|          phone|          202-225-5731|
| 9|  0|          fax|          202-225-3307|
| 9|  1|          phone|          202-225-5731|
| 10| 0|          fax|          202-225-6328|
| 10| 1|          phone|          202-225-4576|
+---+-----+-----+-----+-----+
```

only showing top 20 rows

```
+---+-----+-----+-----+-----+
| id|index|contact_details.val.type|contact_details.val.value|
+---+-----+-----+-----+-----+
| 11|  0|          fax|          202-225-6328|
| 11|  1|          phone|          202-225-4576|
| 11|  2|        twitter|        RepTrentFranks|
| 12|  0|          fax|          202-225-6328|
| 12|  1|          phone|          202-225-4576|
| 12|  2|        twitter|        RepTrentFranks|
| 13|  0|          fax|          202-225-6328|
| 13|  1|          phone|          202-225-4576|
| 13|  2|        twitter|        RepTrentFranks|
| 14|  0|          fax|          202-225-6328|
+---+-----+-----+-----+-----+
```

```

| 14| 1| phone| 202-225-4576|
| 14| 2| twitter| RepTrentFranks|
| 15| 0| fax| 202-225-6328|
| 15| 1| phone| 202-225-4576|
| 15| 2| twitter| RepTrentFranks|
| 16| 0| fax| 202-225-6328|
| 16| 1| phone| 202-225-4576|
| 16| 2| twitter| RepTrentFranks|
| 17| 0| fax| 202-225-6328|
| 17| 1| phone| 202-225-4576|
+---+-----+-----+-----+
only showing top 20 rows

```

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(dfc, key, transformation_ctx = "")
```

DynamicFrameCollection から DynamicFrame を 1 つ取得します。

- `dfc` – 選択すべき DynamicFrame を含む DynamicFrameCollection (必須)。
- `key` – 選択する DynamicFrame のキー (必須)。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。

```
apply(cls, *args, **kwargs)
```

継承元は GlueTransform [apply](#)。

name(cls)

継承元は GlueTransform [name](#)。

describeArgs(cls)

継承元は GlueTransform [describeArgs](#)。

describeReturn(cls)

継承元は GlueTransform [describeReturn](#)。

describeTransform(cls)

継承元は GlueTransform [describeTransform](#)。

describeErrors(cls)

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

Simplify_ddb_json class

特に DynamoDB JSON 構造にある DynamicFrame でネスト化された列を単純化し、新しい単純化された DynamicFrame を返します。

例

特に DynamoDB JSON 構造にある DynamicFrame でネスト化された列を単純化するには、DynamicFrame.simplify_ddb_json() メソッドを使用することをお勧めします。コード例については、「[例: simplify_ddb_json を使用して DynamoDB JSON の単純化を呼び出す](#)」を参照してください。

スピゴットクラス

AWS Glue ジョブで実行された変換が確認しやすくなるように、サンプルレコードを指定した送信先に書き込みます。

例

[DynamicFrame.spigot\(\)](#) メソッドを使用して、レコードのサブセットを DynamicFrame から指定された送信先に書き込むことをお勧めします。コード例については、「[例: spigot を使用して DynamicFrame のサンプルフィールドを Amazon S3 に書き込む](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, path, options, transformation_ctx = "")
```

変換中に指定した場所にサンプルレコードを書き込みます。

- frame - スピゴットする DynamicFrame (必須)。
- path - 書き込み先へのパス (必須)。
- options - オプションを指定する JSON のキーと値のペア (オプション)。`"topk"` オプションは、最初の k レコードを書き込む必要があることを指定します。`"prob"` オプションは、指定されたレコードを選択する確率を (10 進数で) 指定します。これを使用して、書き込むレコードを選択します。
- transformation_ctx - 状態情報を識別するために使用される一意の文字列 (オプション)。

```
apply(cls, *args, **kwargs)
```

GlueTransform [apply](#) から継承されました。

```
name(cls)
```

GlueTransform [name](#) から継承されました。

describeArgs(cls)

GlueTransform [describeArgs](#) から継承されました。

describeReturn(cls)

GlueTransform [describeReturn](#) から継承されました。

describeTransform(cls)

GlueTransform [describeTransform](#) から継承されました。

describeErrors(cls)

GlueTransform [describeErrors](#) から継承されました。

describe(cls)

GlueTransform [説明](#) から継承されました。

SplitFields クラス

指定されたフィールドで DynamicFrame を 2 つに新しく分割します。

例

[DynamicFrame.split_fields\(\)](#) メソッドを使用して、DynamicFrame のフィールドを分割することをお勧めします。コード例については、「[例: split_fields を使用して、選択したフィールドを別の DynamicFrame に分割する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)

• [説明](#)

```
__call__(frame, paths, name1 = None, name2 = None, transformation_ctx = "", info = "",
stageThreshold = 0, totalThreshold = 0)
```

DynamicFrame の 1 つまたは複数のフィールドを新しい DynamicFrame に分割し、残りのフィールドを含む別の新しい DynamicFrame を作成します。

- `frame` – 2 つの新しい DynamicFrame に分割するためのソース (必須)。
- `paths` – 分割されるフィールドへの完全パスのリスト (必須)。
- `name1` – 分割されるフィールドを含む DynamicFrame に割り当てる名前 (オプション)。名前が指定されていない場合、ソースフレームの名前に「1」を付加した名前が使用されます。
- `name2` – 指定されたフィールドの分割後に残されるフィールドを含む、DynamicFrame に割り当てる名前 (オプション)。名前が指定されていない場合、ソースフレームの名前に「2」を付加した名前が使用されます。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` – エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` – エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

```
apply(cls, *args, **kwargs)
```

継承元は `GlueTransform` [apply](#)。

```
name(cls)
```

継承元は `GlueTransform` [name](#)。

```
describeArgs(cls)
```

継承元は `GlueTransform` [describeArgs](#)。

```
describeReturn(cls)
```

継承元は `GlueTransform` [describeReturn](#)。

describeTransform(cls)

継承元は GlueTransform [describeTransform](#)。

describeErrors(cls)

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

SplitRows クラス

2 つの DynamicFrames を含む DynamicFrameCollection を作成します。1 つの DynamicFrame には分割するよう指定された行のみが含まれ、もう片方には残りの行すべてが含まれています。

例

[DynamicFrame.split_rows\(\)](#) メソッドを使用して、DynamicFrame の行を分割することをお勧めします。コード例については、「[例: split_rows を使用して、DynamicFrame 内の行を分割する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(frame, comparison_dict, name1="frame1", name2="frame2", transformation_ctx = "", info = None, stageThreshold = 0, totalThreshold = 0)`

DynamicFrame の 1 つ以上の行を、新しい DynamicFrame に分割します。

- `frame` – 2 つの新しい `DynamicFrame` に分割するためのソース (必須)。
- `comparison_dict` – 列へのフルパスを示すためのキーと、列の値の比較対象である値に対しコンパレータをマッピングする別のディクショナリを示すための値を持つ、ディクショナリ。たとえば、`{"age": {">": 10, "<": 20}}` は、"age" の値が 10 ~ 20 の行を分割します。ただし "age" の値がこの範囲外の行は除外されます (必須)。
- `name1` – 分割される行を含む `DynamicFrame` に割り当てる名前 (省略可能)。
- `name2` – 指定された行が分割された後に残る行を含む `DynamicFrame` に割り当てる名前 (省略可能)。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` – エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` – エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

describe(cls)

継承元は `GlueTransform` [説明](#)。

Unbox クラス

`DynamicFrame` の文字列フィールドをアンボックス (再フォーマット) します。

例

[DynamicFrame.unbox\(\)](#) メソッドを使用して、`DynamicFrame` のフィールドをアンボックスすることをお勧めします。コード例については、「[例: unbox を使用して、文字列フィールドを構造体にアンボックスする](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, path, format, transformation_ctx = "", info="", stageThreshold=0, totalThreshold=0,
**options)
```

`DynamicFrame` の文字列フィールドをアンボックスします。

- `frame` – フィールドをアンボックスする `DynamicFrame` (必須)。
- `path` – アンボックスする `StringNode` への完全パス (必須)。
- `format` – 形式の仕様 (オプション)。Amazon S3 や、複数の形式をサポートする AWS Glue 接続の場合に使用します。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。

- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` - エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` - エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `separator` - 区切りトークン (省略可能)。
- `escaper` - エスケープトークン (省略可能)。
- `skipFirst` - データの最初の行をスキップする必要がある場合は `True`、スキップしない場合は `False` (省略可能)。
- `withSchema` - アンボックスされるデータのスキーマを含む文字列 (オプション)。これは常に `StructType.json` を使用して作成する必要があります。
- `withHeader` - 解凍されるデータにヘッダーが含まれている場合は `True`、そうでない場合は `False` (省略可能)。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

describe(cls)

継承元は `GlueTransform` [説明](#)。

UnnestFrame クラス

`DynamicFrame` をネスト解除し、ネストされたオブジェクトを最上位の要素にフラット化して、配列オブジェクトの結合キーを生成します。

例

[DynamicFrame.unnest\(\)](#) メソッドを使用して、`DynamicFrame` のネストされた構造をフラット化することをお勧めします。コード例については、「[例: unnest を使用して、ネストされたフィールドを最上位フィールドに変換する](#)」を参照してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(frame, transformation_ctx = "", info="", stageThreshold=0, totalThreshold=0)
```

`DynamicFrame` をネスト解除し、ネストされたオブジェクトを最上位の要素にフラット化して、配列オブジェクトの結合キーを生成します。

- `frame` - ネスト解除する `DynamicFrame` (必須)。
- `transformation_ctx` - 状態情報を識別するために使用される一意の文字列 (オプション)。
- `info` - 変換のエラーに関連付けられた文字列 (オプション)。
- `stageThreshold` - エラーを出力する前に、変換で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。
- `totalThreshold` - エラーの出力を処理する前に、全体で発生する可能性のあるエラーの最大数 (オプション)。デフォルト値は 0 です。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

`FlagDuplicatesInColumn` クラス

`FlagDuplicatesInColumn` 変換は、各行に指定された値を持つ新しい列を返します。この列は、行のソース列の値がソース列の前の行の値と一致するかどうかを示します。一致が見つかったら、重複としてフラグが付けられます。最初の出現は、以前の行と一致しないため、フラグは付けられません。

例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsgluedi.transforms import *

sc = SparkContext()
```

```
spark = SparkSession(sc)

datasource1 = spark.read.json("s3://${BUCKET}/json/zips/raw/data")

try:
    df_output = column.FlagDuplicatesInColumn.apply(
        data_frame=datasource1,
        spark_context=sc,
        source_column="city",
        target_column="flag_col",
        true_string="True",
        false_string="False"
    )
except:
    print("Unexpected Error happened ")
    raise
```

出力

FlagDuplicatesInColumn 変換により、新しい列「flag_col」が「df_output」に追加されます DataFrame。この列には、対応する行の「city」列に重複する値があるかどうかを示す文字列値が含まれます。行に重複する「city」値がある場合、「flag_col」には「true_string」値「True」が含まれます。行に一意的「city」値がある場合、「flag_col」には「false_string」値「False」が含まれます。

結果の「df_output DataFrame」には、元の「datasource1」のすべての列と DataFrame、重複する「city」値を示す追加の「flag_col」列が含まれます。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(spark_context, data_frame, source_column, target_column,  
true_string=DEFAULT_TRUE_STRING, false_string=DEFAULT_FALSE_STRING)
```

FlagDuplicatesInColumn 変換は、各行に指定された値を持つ新しい列を返します。この列は、行のソース列の値がソース列の前の行の値と一致するかどうかを示します。一致が見つかったら、重複としてフラグが付けられます。最初の出現は、以前の行と一致しないため、フラグは付けられません。

- source_column - ソース列の名前。
- target_column - ターゲット列の名前。
- true_string - ソース列の値がその列の以前の値と重複した場合にターゲット列に挿入される文字列。
- false_string - ソース列の値がその列の以前の値と異なる場合にターゲット列に挿入される文字列。

```
apply(cls, *args, **kwargs)
```

継承元は GlueTransform [apply](#)。

```
name(cls)
```

継承元は GlueTransform [name](#)。

```
describeArgs(cls)
```

継承元は GlueTransform [describeArgs](#)。

```
describeReturn(cls)
```

継承元は GlueTransform [describeReturn](#)。

```
describeTransform(cls)
```

継承元は GlueTransform [describeTransform](#)。

```
describeErrors(cls)
```

継承元は GlueTransform [describeErrors](#)。

```
describe(cls)
```

継承元は GlueTransform [説明](#)。

FormatPhoneNumber クラス

FormatPhoneNumber 変換は、電話番号文字列がフォーマットされた値に変換される列を返します。

例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsglue.transforms import *

sc = SparkContext()
spark = SparkSession(sc)

input_df = spark.createDataFrame(
    [
        ("408-341-5669",),
        ("4083415669",)
    ],
    ["phone"],
)

try:
    df_output = column_formatting.FormatPhoneNumber.apply(
        data_frame=input_df,
        spark_context=sc,
        source_column="phone",
        default_region="US"
    )
    df_output.show()
except:
    print("Unexpected Error happened ")
    raise
```

出力

出力は次のようになります。

```
...
+-----+
| phone |
+-----+
```

```
| (408) 341-5669 |
| (408) 341-5669 |
+-----+
...
```

FormatPhoneNumber 変換では、「source_column」を「phone」として、「default_region」を「US」として受け取ります。

変換は、最初の形式に関係なく、両方の電話番号を標準の米国形式 `(408) 341-5669` に正常にフォーマットします。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(spark_context, data_frame, source_column, phone_number_format=None,
default_region=None, default_region_column=None)
```

FormatPhoneNumber 変換は、電話番号文字列がフォーマットされた値に変換される列を返します。

- source_column - 既存の列の名前。
- phone_number_format - 電話番号を変換する形式。形式を指定しない場合、デフォルトは E.164 であり、これは、国際認識標準電話番号形式です。有効な値には次のようなものがあります。
 - E164 (E の後の期間は省略)
- default_region - 電話番号自体に国コードが存在しない場合に、電話番号の地域を指定する 2 つまたは 3 つの大文字で構成される有効なリージョンコード。最大で、defaultRegion または defaultRegionColumn を指定できます。
- default_region_column - アドバンスドデータ型 の列の名前 Country。指定された列のリージョンコードは、電話番号自体に国コードが存在しない場合、電話番号の国コードを決定するため

に使用されます。最大で、defaultRegionまたはのいずれかdefaultRegionColumnを指定できます。

apply(cls, *args, **kwargs)

継承元は GlueTransform [apply](#)。

name(cls)

継承元は GlueTransform [name](#)。

describeArgs(cls)

継承元は GlueTransform [describeArgs](#)。

describeReturn(cls)

継承元は GlueTransform [describeReturn](#)。

describeTransform(cls)

継承元は GlueTransform [describeTransform](#)。

describeErrors(cls)

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

FormatCase クラス

FormatCase 変換は、列の各文字列を指定されたケースタイプに変更します。

例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsgluedi.transforms import *

sc = SparkContext()
spark = SparkSession(sc)
```

```
datasource1 = spark.read.json("s3://${BUCKET}/json/zips/raw/data")

try:
    df_output = data_cleaning.FormatCase.apply(
        data_frame=datasource1,
        spark_context=sc,
        source_column="city",
        case_type="LOWER"
    )
except:
    print("Unexpected Error happened ")
    raise
```

出力

FormatCase 変換は、「city」列の値を「case_type=LOWER」パラメータに基づいて小文字に変換します。結果の DataFrame 「df_output」には、元の「datasource1」のすべての列が含まれますが DataFrame、「city」列の値は小文字になります。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(spark_context, data_frame, source_column, case_type)`

FormatCase 変換は、列の各文字列を指定されたケースタイプに変更します。

- `source_column` - 既存の列の名前。
- `case_type` - サポートされているケースタイプは CAPITAL、LOWER、UPPER、です SENTENCE。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

FillWithMode クラス

FillWithMode 変換は、指定した電話番号形式に従って列をフォーマットします。値の一部が同じであるタイブレーカーロジックを指定することもできます。例えば、次の値を考えてみましょう。1 2 2 3 3 4

modeType が MINIMUM の場合 FillWithMode、 はモード値として 2 を返します。modeType が の場合 MAXIMUM、モードは 3 です。 の場合 AVERAGE、モードは 2.5 です。

例

```
from awsglue.context import *
from pyspark.sql import SparkSession
from awsgluedi.transforms import *
```

```
sc = SparkContext()
spark = SparkSession(sc)

input_df = spark.createDataFrame(
    [
        (105.111, 13.12),
        (1055.123, 13.12),
        (None, 13.12),
        (13.12, 13.12),
        (None, 13.12),
    ],
    ["source_column_1", "source_column_2"],
)

try:
    df_output = data_quality.FillWithMode.apply(
        data_frame=input_df,
        spark_context=sc,
        source_column="source_column_1",
        mode_type="MAXIMUM"
    )
    df_output.show()
except:
    print("Unexpected Error happened ")
    raise
```

出力

指定されたコードの出力は次のとおりです。

```
...
+-----+-----+
|source_column_1|source_column_2|
+-----+-----+
| 105.111| 13.12|
| 1055.123| 13.12|
| 1055.123| 13.12|
| 13.12| 13.12|
| 1055.123| 13.12|
+-----+-----+
...
```

「awsglue.data_quality」モジュールからのFillWithMode変換は、「input_df」に適用されず DataFrame。source_column_1 列の「null」値を、その列の Null 以外の値からの最大値 (「mode_type"MAXIMUM」) に置き換えます。

この場合、source_column_1列の最大値は`1055.123`です。したがって、の「null」値はsource_column_1、出力「df_output」で「1055.123 DataFrame」に置き換えられます。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(spark_context, data_frame, source_column, mode_type)`

FillWithMode 変換は、列内の文字列の大文字と小文字をフォーマットします。

- source_column - 既存の列の名前。
- mode_type - データ内のタイ値を解決する方法。この値は、MINIMUM、NONE、AVERAGEまたはのいずれかである必要がありますMAXIMUM。

`apply(cls, *args, **kwargs)`

継承元は GlueTransform [apply](#)。

`name(cls)`

継承元は GlueTransform [name](#)。

`describeArgs(cls)`

継承元は GlueTransform [describeArgs](#)。

describeReturn(cls)

継承元は GlueTransform [describeReturn](#)。

describeTransform(cls)

継承元は GlueTransform [describeTransform](#)。

describeErrors(cls)

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

FlagDuplicateRows クラス

FlagDuplicateRows 変換は、各行に指定された値を持つ新しい列を返します。この列は、その行がデータセット内の以前の行と完全に一致するかどうかを示します。一致が見つかったら、重複としてフラグが付けられます。最初の出現は、前の行と一致しないため、フラグは付けられません。

例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsglue.transforms import *

sc = SparkContext()
spark = SparkSession(sc)

input_df = spark.createDataFrame(
    [
        (105.111, 13.12),
        (13.12, 13.12),
        (None, 13.12),
        (13.12, 13.12),
        (None, 13.12),
    ],
    ["source_column_1", "source_column_2"],
)

try:
```

```
df_output = data_quality.FlagDuplicateRows.apply(  
    data_frame=input_df,  
    spark_context=sc,  
    target_column="flag_row",  
    true_string="True",  
    false_string="False",  
    target_index=1  
)  
except:  
    print("Unexpected Error happened ")  
    raise
```

出力

出力は、列に基づいて、行が重複しているかどうかflag_rowを示す追加のsource_column_1列 PySpark DataFrame を持つ になります。結果の `df_output DataFrame` には、次の行が含まれます。

```
...  
+-----+-----+-----+  
|source_column_1|source_column_2|flag_row|  
+-----+-----+-----+  
| 105.111| 13.12| False|  
| 13.12| 13.12| True|  
| null| 13.12| True|  
| 13.12| 13.12| True|  
| null| 13.12| True|  
+-----+-----+-----+  
...
```

flag_row 列は、行が重複しているかどうかを示します。「true_string」は「True」に設定され、「false_string」は「False」に設定されます。「target_index」は 1 に設定されます。つまり、flag_row列は出力の 2 番目の位置 (インデックス 1) に挿入されます DataFrame。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)

- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(spark_context, data_frame, target_column, true_string=DEFAULT_TRUE_STRING,  
false_string=DEFAULT_FALSE_STRING, target_index=None)
```

FlagDuplicateRows 変換は、各行に指定された値を持つ新しい列を返します。この列は、その行がデータセット内の以前の行と完全に一致するかどうかを示します。一致が見つかり、重複としてフラグが付けられます。最初の出現は、前の行と一致しないため、フラグは付けられません。

- true_string - 行が前の行と一致する場合に挿入される値。
- false_string - 行が一意である場合に挿入される値。
- target_column - データセットに挿入される新しい列の名前。

```
apply(cls, *args, **kwargs)
```

継承元は GlueTransform [apply](#)。

```
name(cls)
```

継承元は GlueTransform [name](#)。

```
describeArgs(cls)
```

継承元は GlueTransform [describeArgs](#)。

```
describeReturn(cls)
```

継承元は GlueTransform [describeReturn](#)。

```
describeTransform(cls)
```

継承元は GlueTransform [describeTransform](#)。

```
describeErrors(cls)
```

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は `GlueTransform` [説明](#)。

`RemoveDuplicates` クラス

選択したソース列で重複した値が発生した場合、`RemoveDuplicates`変換は行全体を削除します。

例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsgluedi.transforms import *

sc = SparkContext()
spark = SparkSession(sc)

input_df = spark.createDataFrame(
    [
        (105.111, 13.12),
        (13.12, 13.12),
        (None, 13.12),
        (13.12, 13.12),
        (None, 13.12),
    ],
    ["source_column_1", "source_column_2"],
)

try:
    df_output = data_quality.RemoveDuplicates.apply(
        data_frame=input_df,
        spark_context=sc,
        source_column="source_column_1"
    )
except:
    print("Unexpected Error happened ")
    raise
```

出力

出力は PySpark DataFrame になり、`source_column_1`列に基づいて重複が削除されます。結果の ``df_output DataFrame`` には、次の行が含まれます。

```
...
+-----+-----+
|source_column_1|source_column_2|
+-----+-----+
| 105.111| 13.12|
| 13.12| 13.12|
| null| 13.12|
+-----+-----+
...
```

source_column_1 値が `13.12` および `null` の行は、source_column_1列に基づいて重複が削除されるため DataFrame、出力に 1 回だけ表示されることに注意してください。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(spark_context, data_frame, source_column)`

選択したソース列で重複した値が発生した場合、RemoveDuplicates変換は行全体を削除します。

- source_column - 既存の列の名前。

`apply(cls, *args, **kwargs)`

継承元は GlueTransform [apply](#)。

`name(cls)`

継承元は GlueTransform [name](#)。

describeArgs(cls)

継承元は GlueTransform [describeArgs](#)。

describeReturn(cls)

継承元は GlueTransform [describeReturn](#)。

describeTransform(cls)

継承元は GlueTransform [describeTransform](#)。

describeErrors(cls)

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

MonthName クラス

MonthName 変換は、日付を表す文字列から、月の名前を含む新しい列を作成します。

例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsgluedi.transforms import *

sc = SparkContext()
spark = SparkSession(sc)

spark.conf.set("spark.sql.legacy.timeParserPolicy", "LEGACY")

input_df = spark.createDataFrame(
    [
        ("20-2018-12",),
        ("2018-20-12",),
        ("20182012",),
        ("12202018",),
        ("20122018",),
        ("20-12-2018",),
    ]
```

```

        ("12/20/2018",),
        ("02/02/02",),
        ("02 02 2009",),
        ("02/02/2009",),
        ("August/02/2009",),
        ("02/june/2009",),
        ("02/2020/june",),
        ("2013-02-21 06:35:45.658505",),
        ("August 02 2009",),
        ("2013/02/21",),
        (None,),
    ],
    ["column_1"],
)

try:
    df_output = datetime_functions.MonthName.apply(
        data_frame=input_df,
        spark_context=sc,
        source_column="column_1",
        target_column="target_column"
    )
    df_output.show()
except:
    print("Unexpected Error happened ")
    raise

```

出力

出力は次のようになります。

```

...
+-----+-----+
| column_1|target_column|
+-----+-----+
|20-2018-12 | December |
|2018-20-12 | null |
| 20182012| null |
| 12202018| null |
| 20122018| null |
|20-12-2018 | December |
|12/20/2018 | December |
| 02/02/02 | February |

```

```
|02 02 2009 | February |
|02/02/2009 | February |
|August/02/2009| August |
|02/june/2009| null |
|02/2020/june| null |
|2013-02-21 06:35:45.658505| February |
|August 02 2009| August |
| 2013/02/21| February |
| null | null |
+-----+-----+
...

```

MonthName 変換では、「source_column」は「column_1」、「target_column」は「target_column」になります。「column_1」列の日付/時刻文字列から月名を抽出し、「target_column」列に配置します。日付/時刻文字列が認識されない形式であるか、解析できない場合、「target_column」値は「null」に設定されます。

変換は、「20-12-2018」、「」、「」、12/20/201802/02/20092013-02-21 06:35:45.658505」、「2009年8月2日」など、さまざまな日付/時刻形式から月名を正常に抽出します。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(spark_context, data_frame, target_column, source_column=None, value=None)`

MonthName 変換は、日付を表す文字列から、月の名前を含む新しい列を作成します。

- source_column - 既存の列の名前。
- value - 評価する文字列。

- `target_column` – 新しく作成された列の名前。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

IsEven クラス

IsEven 変換は、ソース列または値が偶数であるかどうかを示すブール値を新しい列に返します。ソース列または値が 10 進数の場合、結果は `false` です。

例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsgluedi.transforms import *
```

```
sc = SparkContext()
spark = SparkSession(sc)

input_df = spark.createDataFrame(
    [(5,), (0,), (-1,), (2,), (None,)],
    ["source_column"],
)

try:
    df_output = math_functions.IsEven.apply(
        data_frame=input_df,
        spark_context=sc,
        source_column="source_column",
        target_column="target_column",
        value=None,
        true_string="Even",
        false_string="Not even",
    )
    df_output.show()
except:
    print("Unexpected Error happened ")
    raise
```

出力

出力は次のようになります。

```
...
+-----+-----+
|source_column|target_column|
+-----+-----+
| 5| Not even|
| 0| Even|
| -1| Not even|
| 2| Even|
| null| null|
+-----+-----+
...
```

IsEven 変換では、「source_column」を「source_column」として、「target_column」を「target_column」として受け取ります。「source_column」の値が偶数かどうかをチェックします。値が偶数の場合、「target_column」の値を「true_string」「Even」に設定します。値が奇数の場

合、「target_column」の値を「false_string」「Not even」に設定します。「source_column」値が「null」の場合、「target_column」値は「null」に設定されます。

変換は偶数 (0 と 2) を正しく識別し、「target_column」値を「Even」に設定します。奇数 (5 と -1) の場合、「target_column」の値を「Not even」に設定します。`"source_column"` の `null` 値の場合、`"target_column"` 値は `null` に設定されます。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(spark_context, data_frame, target_column, source_column=None,  
true_string=DEFAULT_TRUE_STRING, false_string=DEFAULT_FALSE_STRING, value=None)
```

IsEven 変換は、ソース列または値が偶数であるかどうかを示すブール値を新しい列に返します。ソース列または値が 10 進数の場合、結果は false です。

- source_column - 既存の列の名前。
- target_column - 作成する新しい列の名前。
- true_string - 値が偶数であるかどうかを示す文字列。
- false_string - 値が偶数ではないかどうかを示す文字列。

```
apply(cls, *args, **kwargs)
```

継承元は GlueTransform [apply](#)。

```
name(cls)
```

継承元は GlueTransform [name](#)。

describeArgs(cls)

継承元は GlueTransform [describeArgs](#)。

describeReturn(cls)

継承元は GlueTransform [describeReturn](#)。

describeTransform(cls)

継承元は GlueTransform [describeTransform](#)。

describeErrors(cls)

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

CryptographicHash クラス

CryptographicHash 変換は、アルゴリズムを列のハッシュ値に適用します。

例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsgluedi.transforms import *

secret = "${SECRET}"
sc = SparkContext()
spark = SparkSession(sc)

input_df = spark.createDataFrame(
    [
        (1, "1234560000"),
        (2, "1234560001"),
        (3, "1234560002"),
        (4, "1234560003"),
        (5, "1234560004"),
        (6, "1234560005"),
```

```

        (7, "1234560006"),
        (8, "1234560007"),
        (9, "1234560008"),
        (10, "1234560009"),
    ],
    ["id", "phone"],
)

try:
    df_output = pii.CryptographicHash.apply(
        data_frame=input_df,
        spark_context=sc,
        source_columns=["id", "phone"],
        secret_id=secret,
        algorithm="HMAC_SHA256",
        output_format="BASE64",
    )
    df_output.show()
except:
    print("Unexpected Error happened ")
    raise

```

出力

出力は次のようになります。

```

...
+---+-----+-----+-----+
| id| phone | id_hashed | phone_hashed |
+---+-----+-----+-----+
| 1| 1234560000 | QUI1zXTJiXmfIb... | juDBAmiRnn03g... |
| 2| 1234560001 | ZAUWiZ3dVTzCo... | vC81gUqBVDMNQ... |
| 3| 1234560002 | ZP4VvZWkqYifu... | K13QAKgswYpzB... |
| 4| 1234560003 | 3u8v03wQ8EQfj... | CPBzK1P8PZZkV... |
| 5| 1234560004 | eWkQJk4zA0Izx... | aLf7+mHcXqbLs... |
| 6| 1234560005 | xtI9fZCJZCvsa... | dy2DFgdYWmr0p... |
| 7| 1234560006 | iW9hew7jnHu0f... | wwFGMCOEv6o0v... |
| 8| 1234560007 | H9V1pqvgkFhfS... | g9WKhagIXy9ht... |
| 9| 1234560008 | xDhEuHaxAUbU5... | b3uQLKPY+Q5vU... |
| 10| 1234560009 | GRN6nFXkxk349... | VJdsKt8VbxBbt... |
+---+-----+-----+-----+
...

```

変換は、指定されたアルゴリズムとシークレットキーを使用して「id」列と「phone」列の値の暗号化ハッシュを計算し、Base64形式でハッシュをエンコードします。結果の`df_output DataFrame`には、元の`input_df`のすべての列に加えて DataFrame、計算されたハッシュを持つ追加の`id_hashed`列と`phone_hashed`列が含まれます。

方法

- [__call__](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(spark_context, data_frame, source_columns, secret_id, algorithm=None, secret_version=None, create_secret_if_missing=False, output_format=None, entity_type_filter=None)
```

CryptographicHash 変換は、アルゴリズムを列のハッシュ値に適用します。

- `source_columns` – 既存の列の配列。
- `secret_id` – Secrets Manager シークレットキーの ARN。ソース列をハッシュするためにハッシュベースのメッセージ認証コード (HMAC) プレフィックスアルゴリズムで使用されるキー。
- `secret_version` – オプション。デフォルトは最新のシークレットバージョンです。
- `entity_type_filter` – エンティティタイプのオプションの配列。検出された PII のみをフリーテキスト列で暗号化するために使用できます。
- `create_secret_if_missing` – オプションのブール値。true の場合、 は発信者に代わってシークレットを作成しようとします。
- `algorithm` – データのハッシュに使用されるアルゴリズム。有効な列挙値: MD5, SHA1, SHA256, SHA512, HMAC_MD5, HMAC_SHA1, HMAC_SHA256, HMAC_SHA512。

```
apply(cls, *args, **kwargs)
```

継承元は `GlueTransform` [apply](#)。

name(cls)

継承元は GlueTransform [name](#)。

describeArgs(cls)

継承元は GlueTransform [describeArgs](#)。

describeReturn(cls)

継承元は GlueTransform [describeReturn](#)。

describeTransform(cls)

継承元は GlueTransform [describeTransform](#)。

describeErrors(cls)

継承元は GlueTransform [describeErrors](#)。

describe(cls)

継承元は GlueTransform [説明](#)。

復号クラス

Decrypt 変換は AWS Glue 内で復号化されます。AWS Encryption SDK を使用して、AWS Glue の外部でデータを復号することもできます。指定された KMS キー ARN が列の暗号化に使用されたものと一致しない場合、復号オペレーションは失敗します。

例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsgluedi.transforms import *

kms = "${KMS}"
sc = SparkContext()
spark = SparkSession(sc)

input_df = spark.createDataFrame(
    [
```

```
(1, "1234560000"),
(2, "1234560001"),
(3, "1234560002"),
(4, "1234560003"),
(5, "1234560004"),
(6, "1234560005"),
(7, "1234560006"),
(8, "1234560007"),
(9, "1234560008"),
(10, "1234560009"),
],
["id", "phone"],
)

try:
    df_encrypt = pii.Encrypt.apply(
        data_frame=input_df,
        spark_context=sc,
        source_columns=["phone"],
        kms_key_arn=kms
    )
    df_decrypt = pii.Decrypt.apply(
        data_frame=df_encrypt,
        spark_context=sc,
        source_columns=["phone"],
        kms_key_arn=kms
    )
    df_decrypt.show()
except:
    print("Unexpected Error happened ")
    raise
```

出力

出力は、元の「id」列と復号された「phone」列 PySpark DataFrame を持つ になります。

```
...
+---+-----+
| id| phone|
+---+-----+
| 1| 1234560000|
| 2| 1234560001|
| 3| 1234560002|
```

```
| 4| 1234560003|
| 5| 1234560004|
| 6| 1234560005|
| 7| 1234560006|
| 8| 1234560007|
| 9| 1234560008|
| 10| 1234560009|
+---+-----+
```
```

Encrypt 変換では、`source\_columns` を `["phone"]` として、`kms\_key\_arn` を `\${KMS}` 環境変数の値として受け取ります。変換は、指定された KMS キーを使用して「phone」列の値を暗号化します。次に、暗号化された DataFrame 「df\_encrypt」が「awsglue.pii」モジュールからDecrypt変換に渡されます。`source\_columns` を `["phone"]` として、`kms\_key\_arn` を `\${KMS}` 環境変数の値として受け取ります。変換は、同じ KMS キーを使用して、「phone」列の暗号化された値を復号します。結果の「df\_decrypt DataFrame」には、元の「id」列と復号された「phone」列が含まれています。

## 方法

- [\\_\\_call\\_\\_](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(spark_context, data_frame, source_columns, kms_key_arn)`

Decrypt 変換は AWS Glue 内で復号化されます。AWS Encryption SDK を使用して、AWS Glue の外部でデータを復号することもできます。指定された KMS キー ARN が列の暗号化に使用されたものと一致しない場合、復号オペレーションは失敗します。

- `source_columns` – 既存の列の配列。
- `kms_key_arn` – ソース列の復号に使用する AWS Key Management Service キーのキー ARN。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

## 暗号化クラス

Encrypt 変換は、Key Management Service AWS キーを使用してソース列を暗号化します。Encrypt 変換は、セルあたり最大 128 MiB を暗号化できます。復号時に形式を保持しようとします。データ型を保持するには、データ型メタデータを 1KB 未満にシリアル化する必要があります。それ以外の場合は、`preserve_data_type`パラメータを `false` に設定する必要があります。データ型のメタデータは、暗号化コンテキストのプレーンテキストで保存されます。

## 例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsgluedi.transforms import *
```

```
kms = "${KMS}"
sc = SparkContext()
spark = SparkSession(sc)

input_df = spark.createDataFrame(
 [
 (1, "1234560000"),
 (2, "1234560001"),
 (3, "1234560002"),
 (4, "1234560003"),
 (5, "1234560004"),
 (6, "1234560005"),
 (7, "1234560006"),
 (8, "1234560007"),
 (9, "1234560008"),
 (10, "1234560009"),
],
 ["id", "phone"],
)

try:
 df_encrypt = pii.Encrypt.apply(
 data_frame=input_df,
 spark_context=sc,
 source_columns=["phone"],
 kms_key_arn=kms
)
except:
 print("Unexpected Error happened ")
 raise
```

## 出力

出力は、元の PySpark DataFrame 「id」 列と 「phone」 列の暗号化された値を含む追加の列を持つになります。

```
...
+---+-----+-----+
| id| phone | phone_encrypted |
+---+-----+-----+
| 1| 1234560000| EncryptedData1234...abc |
```

```

2	1234560001	EncryptedData5678...def
3	1234560002	EncryptedData9012...ghi
4	1234560003	EncryptedData3456...jkl
5	1234560004	EncryptedData7890...mno
6	1234560005	EncryptedData1234...pqr
7	1234560006	EncryptedData5678...stu
8	1234560007	EncryptedData9012...vwx
9	1234560008	EncryptedData3456...yz0
10	1234560009	EncryptedData7890...123
+---+-----+-----+
...

```

Encrypt 変換では、`source\_columns` を `["phone"]` として、`kms\_key\_arn` を `\${KMS}` 環境変数の値として受け取ります。変換は、指定された KMS キーを使用して「phone」列の値を暗号化します。結果の「df\_encrypt DataFrame」には、元の「id」列、元の「phone」列、および「phone」列の暗号化された値を含む「phone\_encrypted」という名前の追加列が含まれています。

## 方法

- [\\_\\_call\\_\\_](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

```
__call__(spark_context, data_frame, source_columns, kms_key_arn, entity_type_filter=None,
preserve_data_type=None)
```

Encrypt 変換は、Key Management Service AWS キーを使用してソース列を暗号化します。

- `source_columns` – 既存の列の配列。
- `kms_key_arn` – ソース列の暗号化に使用する AWS Key Management Service キーのキー ARN。
- `entity_type_filter` – エンティティタイプのオプションの配列。検出された PII のみをフリーテキスト列で暗号化するために使用できます。

- `preserve_data_type` – オプションのブール値。デフォルトは `true` です。 `false` の場合、データ型は保存されません。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

IntToIp クラス

`IntToIp` 変換は、ソース列の整数値または他の値を、次にターゲット列の対応する IPv4 値に変換し、その結果を新しい列で返します。

例

```
from pyspark.context import SparkContext
from pyspark.sql import SparkSession
from awsgluedi.transforms import *
```

```
sc = SparkContext()
spark = SparkSession(sc)

input_df = spark.createDataFrame(
 [
 (3221225473,),
 (0,),
 (1,),
 (100,),
 (168430090,),
 (4294967295,),
 (4294967294,),
 (4294967296,),
 (-1,),
 (None,)
],
 ["source_column_int"],
)

try:
 df_output = web_functions.IntToIp.apply(
 data_frame=input_df,
 spark_context=sc,
 source_column="source_column_int",
 target_column="target_column",
 value=None
)
 df_output.show()
except:
 print("Unexpected Error happened ")
 raise
```

## 出力

出力は次のようになります。

```
...
+-----+-----+
|source_column_int|target_column|
+-----+-----+
3221225473	192.0.0.1
0	0.0.0.0
1	0.0.0.1
```

```

100	0.0.0.100
168430090	10.0.0.10
4294967295	255.255.255.255
4294967294	255.255.255.254
4294967296	null
-1	null
null	null
+-----+-----+
...

```

IntToIp.apply 変換では、「source\_column」を「source\_column\_int」として受け取り、「target\_column」を「target\_column」として受け取り、「source\_column\_int」列の整数値を対応する IPv4 アドレス表現に変換して、結果を「target\_column」列に保存します。

IPv4 アドレス (0 ~ 4294967295) の範囲内の有効な整数値の場合、変換は IPv4 アドレス表現 (例: 192.0.0.1、0.0.0、10.0.0.10、255.255.255.255) に正常に変換します。

有効範囲外の整数値 (例: 4294967296、-1) の場合、「target\_column」値は「null」に設定されます。「source\_column\_int」列の「null」値の場合、「target\_column」値も「null」に設定されます。

## 方法

- [\\_\\_call\\_\\_](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(spark_context, data_frame, target_column, source_column=None, value=None)`

IntToIp 変換は、ソース列の整数値または他の値を、次にターゲット列の対応する IPv4 値に変換し、その結果を新しい列で返します。

- sourceColumn - 既存の列の名前。
- value - 評価する文字列。

- `targetColumn` – 作成する新しい列の名前。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

## IpToInt クラス

IpToInt 変換は、ソース列の Internet Protocol version 4 (IPv4) 値または他の値をターゲット列の対応する整数値に変換し、新しい列に結果を返します。

### 例

AWS Glue 4.0 以降では、を使用してジョブ引数を作成または更新します。key: `--enable-glue-di-transforms`, value: `true`

```
from pyspark.context import SparkContext
```

```
from awsgluedi.transforms import *

sc = SparkContext()

input_df = spark.createDataFrame(
 [
 ("192.0.0.1",),
 ("10.10.10.10",),
 ("1.2.3.4",),
 ("1.2.3.6",),
 ("http://12.13.14.15",),
 ("https://16.17.18.19",),
 ("1.2.3.4",),
 (None,),
 ("abc",),
 ("abc.abc.abc.abc",),
 ("321.123.123.123",),
 ("244.4.4.4",),
 ("255.255.255.255",),
],
 ["source_column_ip"],
)

df_output = web_functions.IpToInt.apply(
 data_frame=input_df,
 spark_context=sc,
 source_column="source_column_ip",
 target_column="target_column",
 value=None
)
df_output.show()
```

## 出力

出力は次のようになります。

```
...
+-----+-----+
|source_column_ip| target_column|
+-----+-----+
192.0.0.1	3221225473
10.10.10.10	168427722
1.2.3.4	16909060
```

```

1.2.3.6	16909062
http://12.13.14.15	null
https://16.17.18.19	null
1.2.3.4	16909060
null	null
abc	null
abc.abc.abc.abc	null
321.123.123.123	null
244.4.4.4	4102444804
255.255.255.255	4294967295
+-----+-----+
...

```

IpToInt 変換では、「source\_column」を「source\_column\_ip」として受け取り、「target\_column」を「target\_column」として受け取り、「source\_column\_ip」列の有効な IPv4 アドレス文字列を対応する 32 ビット整数表現に変換して、結果を「target\_column」列に格納します。

有効な IPv4 アドレス文字列 (例: "192.0.0.1"、"10.10.10.10"、"1.2.3.4") の場合、変換はそれらを整数表現 (例: 3221225473、168427722、16909060) に正常に変換します。有効な IPv4 アドレスではない文字列 (URL、URLs 「abc」のような非 IP 文字列、「abc.abc.abc.abc」のような無効な IP 形式など) の場合、「target\_column」値は「null」に設定されます。「source\_column\_ip」列の「null」値の場合、「target\_column」値も「null」に設定されます。

## 方法

- [\\_\\_call\\_\\_](#)
- [apply](#)
- [name](#)
- [describeArgs](#)
- [describeReturn](#)
- [describeTransform](#)
- [describeErrors](#)
- [説明](#)

`__call__(spark_context, data_frame, target_column, source_column=None, value=None)`

IpToInt 変換は、ソース列の Internet Protocol version 4 (IPv4) 値または他の値をターゲット列の対応する整数値に変換し、新しい列に結果を返します。

- `sourceColumn` - 既存の列の名前。
- `value` - 評価する文字列。
- `targetColumn` - 作成する新しい列の名前。

`apply(cls, *args, **kwargs)`

継承元は `GlueTransform` [apply](#)。

`name(cls)`

継承元は `GlueTransform` [name](#)。

`describeArgs(cls)`

継承元は `GlueTransform` [describeArgs](#)。

`describeReturn(cls)`

継承元は `GlueTransform` [describeReturn](#)。

`describeTransform(cls)`

継承元は `GlueTransform` [describeTransform](#)。

`describeErrors(cls)`

継承元は `GlueTransform` [describeErrors](#)。

`describe(cls)`

継承元は `GlueTransform` [説明](#)。

## データ統合変換

AWS Glue 4.0 以降では、`enable-glue-di-transforms` を使用してジョブ引数を作成または更新しますkey: `--enable-glue-di-transforms`, value: `true`。

ジョブスクリプトの例 :

```
from pyspark.context import SparkContext
```

```
from awsgluedi.transforms import *
sc = SparkContext()

input_df = spark.createDataFrame(
 [(5,), (0,), (-1,), (2,), (None,)],
 ["source_column"],
)

try:
 df_output = math_functions.IsEven.apply(
 data_frame=input_df,
 spark_context=sc,
 source_column="source_column",
 target_column="target_column",
 value=None,
 true_string="Even",
 false_string="Not even",
)
 df_output.show()
except:
 print("Unexpected Error happened ")
 raise
```

## ノートブックを使用したセッションの例

```
%idle_timeout 2880
%glue_version 4.0
%worker_type G.1X
%number_of_workers 5
%region eu-west-1
```

```
%%configure
{
 "--enable-glue-di-transforms": "true"
}
```

```
from pyspark.context import SparkContext
from awsgluedi.transforms import *

sc = SparkContext()
```

```
input_df = spark.createDataFrame(
 [(5,), (0,), (-1,), (2,), (None,)],
 ["source_column"],
)

try:
 df_output = math_functions.IsEven.apply(
 data_frame=input_df,
 spark_context=sc,
 source_column="source_column",
 target_column="target_column",
 value=None,
 true_string="Even",
 false_string="Not even",
)
 df_output.show()
except:
 print("Unexpected Error happened ")
 raise
```

## を使用したセッションの例 AWS CLI

```
aws glue create-session --default-arguments "--enable-glue-di-transforms=true"
```

## DI 変換 :

- [FlagDuplicatesInColumn クラス](#)
- [FormatPhoneNumber クラス](#)
- [FormatCase クラス](#)
- [FillWithMode クラス](#)
- [FlagDuplicateRows クラス](#)
- [RemoveDuplicates クラス](#)
- [MonthName クラス](#)
- [IsEven クラス](#)
- [CryptographicHash クラス](#)
- [復号クラス](#)
- [暗号化クラス](#)

- [IntToIp クラス](#)
- [IpToInt クラス](#)

Maven: プラグインを Spark アプリケーションにバンドルする

Spark アプリケーションをローカルで開発pom.xmlしながら Maven にプラグインの依存関係を追加することで、変換の依存関係を Spark アプリケーションと Spark ディストリビューション (バージョン 3.3) にバンドルできます。

```
<repositories>
 ...
 <repository>
 <id>aws-glue-etl-artifacts</id>
 <url>https://aws-glue-etl-artifacts.s3.amazonaws.com/release/ </url>
 </repository>
</repositories>
...
<dependency>
 <groupId>com.amazonaws</groupId>
 <artifactId>AWSGlueTransforms</artifactId>
 <version>4.0.0</version>
</dependency>
```

別の方法として、次のように AWS Glue Maven アーティファクトから直接バイナリをダウンロードし、Spark アプリケーションに含めることもできます。

```
#!/bin/bash
sudo wget -v https://aws-glue-etl-artifacts.s3.amazonaws.com/release/com.amazonaws/AWSGlueTransforms/4.0.0/AWSGlueTransforms-4.0.0.jar -P /usr/lib/spark/jars/
```

## Scala での AWS Glue ETL スクリプトのプログラミング

AWS Glue 用の Scala コード例とユーティリティは、GitHub ウェブサイトの [AWS Glue サンプルリポジトリ](#)で見つかります。

AWS Glue では、抽出、変換、ロード (ETL) ジョブを実行する PySpark Scala 言語のスクリプトのために、拡張機能がサポートされています。以下のセクションでは、AWS Glue Scala ライブラリと AWS Glue API を ETL スクリプトで使用する方法について説明します。また、ライブラリのリファレンス資料を提供します。

## 目次

- [Scala を使用した AWS Glue ETL スクリプトのプログラミング](#)
  - [開発エンドポイントでの Jupyter Notebook を使用した Scala ETL プログラムのテスト](#)
  - [Scala REPL での Scala ETL プログラムのテスト](#)
- [Scala スクリプトの例 - ストリーミング ETL](#)
- [AWS Glue Scala ライブラリの API](#)
  - [com.amazonaws.services.glue](#)
  - [com.amazonaws.services.glue.ml](#)
  - [com.amazonaws.services.glue.dq](#)
  - [com.amazonaws.services.glue.types](#)
  - [com.amazonaws.services.glue.util](#)
- [AWS Glue Scala ChoiceOption API](#)
  - [ChoiceOption 特性](#)
  - [ChoiceOption オブジェクト](#)
    - [Def apply](#)
  - [ケースクラス ChoiceOptionWithResolver](#)
  - [ケースクラス MatchCatalogSchemaChoiceOption](#)
- [Abstract DataSink クラス](#)
  - [Def writeDynamicFrame](#)
  - [Def pyWriteDynamicFrame](#)
  - [Def writeDataFrame](#)
  - [Def pyWriteDataFrame](#)
  - [Def setCatalogInfo](#)
  - [Def supportsFormat](#)
  - [Def setFormat](#)
  - [Def withFormat](#)
  - [Def setAccumulableSize](#)
  - [Def getOutputErrorRecordsAccumulable](#)
  - [Def errorsAsDynamicFrame](#)
- [DataSink オブジェクト](#)

- [Def recordMetrics](#)
- [AWS Glue Scala DataSource トレイト](#)
- [AWS Glue Scala DynamicFrame API](#)
  - [AWS Glue Scala DynamicFrame クラス](#)
    - [Val errorsCount](#)
    - [Def applyMapping](#)
    - [Def assertErrorThreshold](#)
    - [Def count](#)
    - [Def dropField](#)
    - [Def dropFields](#)
    - [Def dropNulls](#)
    - [Def errorsAsDynamicFrame](#)
    - [Def filter](#)
    - [Def getName](#)
    - [Def getNumPartitions](#)
    - [Def getSchemalfComputed](#)
    - [Def isSchemaComputed](#)
    - [Def javaToPython](#)
    - [Def join](#)
    - [Def map](#)
    - [Def mergeDynamicFrames](#)
    - [Def printSchema](#)
    - [Def recomputeSchema](#)
    - [Def relationalize](#)
    - [Def renameField](#)
    - [Def repartition](#)
    - [Def resolveChoice](#)
    - [Def schema](#)
    - [Def selectField](#)
    - [Def selectFields](#)

- [Def show](#)
- [Def simplifyDDBJson](#)
- [Def spigot](#)
- [Def splitFields](#)
- [Def splitRows](#)
- [Def stageErrorsCount](#)
- [Def toDF](#)
- [Def unbox](#)
- [Def unnest](#)
- [Def unnestDDBJson](#)
- [Def withFrameSchema](#)
- [Def withName](#)
- [Def withTransformationContext](#)
- [DynamicFrame オブジェクト](#)
  - [Def apply](#)
  - [Def emptyDynamicFrame](#)
  - [Def fromPythonRDD](#)
  - [Def ignoreErrors](#)
  - [Def inlineErrors](#)
  - [Def newFrameWithErrors](#)
- [AWS Glue Scala DynamicRecord クラス](#)
  - [def addField](#)
  - [def dropField](#)
  - [def setError](#)
  - [def isError](#)
  - [def getError](#)
  - [def clearError](#)
  - [def write](#)
  - [def readFields](#)
  - [def clone](#)

- [def schema](#)
- [def getRoot](#)
- [def toJson](#)
- [def getFieldNode](#)
- [def getField](#)
- [def hashCode](#)
- [def equals](#)
- [DynamicRecord オブジェクト](#)
  - [def apply](#)
- [RecordTraverser 特性](#)
- [AWS Glue スカラ API GlueContext](#)
  - [デフ addIngestionTime-カラム](#)
  - [def createDataFrame FromOptions](#)
  - [forEachBatch](#)
  - [def getCatalogSink](#)
  - [デフ getCatalogSource](#)
  - [def getJDBCSink](#)
  - [def getSink](#)
  - [getSinkWithdef フォーマット](#)
  - [def getSource](#)
  - [def フォーマット getSourceWith](#)
  - [def getSparkSession](#)
  - [def startTransaction](#)
  - [def commitTrac](#)
  - [def cancelTransaction](#)
  - [def this](#)
  - [def this](#)
  - [def this](#)
- [MappingSpec](#)
  - [MappingSpec ケースクラス](#)

- [MappingSpec オブジェクト](#)
- [Val orderingByTarget](#)
- [Def apply](#)
- [Def apply](#)
- [Def apply](#)
- [AWS Glue Scala ResolveSpec API](#)
  - [ResolveSpec オブジェクト](#)
    - [def](#)
    - [def](#)
  - [ResolveSpec ケースクラス](#)
    - [ResolveSpec def メソッド](#)
- [AWS Glue Scala ArrayNode API](#)
  - [ArrayNode ケースクラス](#)
    - [ArrayNode def メソッド](#)
- [AWS Glue Scala BinaryNode API](#)
  - [BinaryNode ケースクラス](#)
    - [BinaryNode val フィールド](#)
    - [BinaryNode def メソッド](#)
- [AWS Glue Scala BooleanNode API](#)
  - [BooleanNode ケースクラス](#)
    - [BooleanNode val フィールド](#)
    - [BooleanNode def メソッド](#)
- [AWS Glue Scala ByteNode API](#)
  - [ByteNode ケースクラス](#)
    - [ByteNode val フィールド](#)
    - [ByteNode def メソッド](#)
- [AWS Glue Scala DateNode API](#)
  - [DateNode ケースクラス](#)
    - [DateNode val フィールド](#)
    - [DateNode def メソッド](#)

- [AWS Glue Scala DecimalNode API](#)
  - [DecimalNode ケースクラス](#)
    - [DecimalNode val フィールド](#)
    - [DecimalNode def メソッド](#)
- [AWS Glue Scala DoubleNode API](#)
  - [DoubleNode ケースクラス](#)
    - [DoubleNode val フィールド](#)
    - [DoubleNode def メソッド](#)
- [AWS Glue Scala DynamicNode API](#)
  - [DynamicNode クラス](#)
    - [DynamicNode def メソッド](#)
  - [DynamicNode オブジェクト](#)
    - [DynamicNode def メソッド](#)
- [EvaluateDataQuality クラス](#)
  - [def apply](#)
  - [例](#)
- [AWS Glue Scala FloatNode API](#)
  - [FloatNode ケースクラス](#)
    - [FloatNode val フィールド](#)
    - [FloatNode def メソッド](#)
- [FillMissingValues クラス](#)
  - [def apply](#)
- [FindMatches クラス](#)
  - [def apply](#)
- [FindIncrementalMatches クラス](#)
  - [def apply](#)
- [AWS Glue Scala IntegerNode API](#)
  - [IntegerNode ケースクラス](#)
    - [IntegerNode val フィールド](#)
    - [IntegerNode def メソッド](#)

- [AWS Glue Scala LongNode API](#)
  - [LongNode ケースクラス](#)
    - [LongNode val フィールド](#)
    - [LongNode def メソッド](#)
- [AWS Glue Scala MapLikeNode API](#)
  - [MapLikeNode クラス](#)
    - [MapLikeNode def メソッド](#)
- [AWS Glue Scala MapNode API](#)
  - [MapNode ケースクラス](#)
    - [MapNode def メソッド](#)
- [AWS Glue Scala NullNode API](#)
  - [NullNode クラス](#)
  - [NullNode ケースオブジェクト](#)
- [AWS Glue Scala ObjectNode API](#)
  - [ObjectNode オブジェクト](#)
    - [ObjectNode def メソッド](#)
  - [ObjectNode ケースクラス](#)
    - [ObjectNode def メソッド](#)
- [AWS Glue Scala ScalarNode API](#)
  - [ScalarNode クラス](#)
    - [ScalarNode def メソッド](#)
  - [ScalarNode オブジェクト](#)
    - [ScalarNode def メソッド](#)
- [AWS Glue Scala ShortNode API](#)
  - [ShortNode ケースクラス](#)
    - [ShortNode val フィールド](#)
    - [ShortNode def メソッド](#)
- [AWS Glue Scala StringNode API](#)
  - [StringNode ケースクラス](#)
    - [StringNode val フィールド](#)

- [StringNode def メソッド](#)
- [AWS Glue Scala TimestampNode API](#)
  - [TimestampNode ケースクラス](#)
    - [TimestampNode val フィールド](#)
    - [TimestampNode def メソッド](#)
- [AWS Glue Scala GlueArgParser API](#)
  - [GlueArgParser オブジェクト](#)
    - [GlueArgParser def メソッド](#)
- [AWS Glue Scala Job API](#)
  - [Job オブジェクト](#)
    - [Job def メソッド](#)

## Scala を使用した AWS Glue ETL スクリプトのプログラミング

AWS Glue コンソールを使用して Scala ETL (抽出、変換、ロード) プログラムを自動的に生成し、必要に応じて変更してからジョブに割り当てることができます。または、独自のプログラムを最初から記述することもできます。詳細については、「[での Spark ジョブのジョブプロパティの設定 AWS Glue](#)」を参照してください。次に AWS Glue はサーバーで Scala プログラムをコンパイルし、その後、関連付けられているジョブを実行します。

プログラムをエラーなしでコンパイルして正常に実行するには、ジョブで実行する前にプログラムを REPL (Read-Eval-Print Loop) や Jupyter Notebook の開発エンドポイントにロードし、テストすることが重要です。コンパイルプロセスはサーバーで実行されるため、そこで問題が発生してもよく確認できません。

開発エンドポイントでの Jupyter Notebook を使用した Scala ETL プログラムのテスト

AWS Glue 開発エンドポイントで Scala プログラムをテストするには、「[開発エンドポイントの追加](#)」の説明に従って開発エンドポイントを設定します。

次に、コンピュータでローカルに実行されているか、Amazon EC2 ノートブックサーバーでリモートに実行されている Jupyter Notebook に開発エンドポイントを接続します。Jupyter Notebook のローカルバージョンをインストールするには、「[チュートリアル： JupyterLab の Jupyter Notebook](#)」の手順に従います。

Scala コードの実行とノートブックでの PySpark コードの実行が唯一異なる点は、ノートブックでは以下のコードを使用して各パラグラフを開始する必要があることです。

```
%spark
```

これにより、ノートブックサーバーがデフォルトで Spark インタープリタの PySpark フレーバーになるのを防止できます。

## Scala REPL での Scala ETL プログラムのテスト

AWS Glue Scala REPL を使用して開発エンドポイントで Scala プログラムをテストできます。

「[チュートリアル: SageMaker ノートブックを使用する](#)」の手順に従います。ただし、SSH-to-REPL コマンドの最後では、`-t gluepyspark` を `-t glue-spark-shell` に置き換えます。これにより、AWS Glue Scala REPL が呼び出されます。

完了後に REPL を閉じるには、「`sys.exit`」と入力します。

## Scala スクリプトの例 - ストリーミング ETL

### Example

次のスクリプト例では、Amazon Kinesis Data Streams に接続し、Data Catalog のスキーマを使用してデータストリームを解析した上で、そのストリームを Amazon S3 にある静的データセットに結合します。結合された結果は Parquet 形式で Amazon S3 に出力されます。

```
// This script connects to an Amazon Kinesis stream, uses a schema from the data
// catalog to parse the stream,
// joins the stream to a static dataset on Amazon S3, and outputs the joined results to
// Amazon S3 in parquet format.
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import java.util.Calendar
import org.apache.spark.SparkContext
import org.apache.spark.sql.Dataset
import org.apache.spark.sql.Row
import org.apache.spark.sql.SaveMode
import org.apache.spark.sql.Session
import org.apache.spark.sql.functions.from_json
import org.apache.spark.sql.streaming.Trigger
import scala.collection.JavaConverters._

object streamJoiner {
 def main(sysArgs: Array[String]) {
 val spark: SparkContext = new SparkContext()
```

```

val glueContext: GlueContext = new GlueContext(spark)
val sparkSession: SparkSession = glueContext.getSparkSession
import sparkSession.implicits._
// @params: [JOB_NAME]
val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
Job.init(args("JOB_NAME"), glueContext, args.asJava)

val staticData = sparkSession.read // read() returns type DataFrameReader
 .format("csv")
 .option("header", "true")
 .load("s3://awsexamplebucket-streaming-demo2/inputs/productsStatic.csv") //
load() returns a DataFrame

val datasource0 = sparkSession.readStream // readstream() returns type
DataStreamReader
 .format("kinesis")
 .option("streamName", "stream-join-demo")
 .option("endpointUrl", "https://kinesis.us-east-1.amazonaws.com")
 .option("startingPosition", "TRIM_HORIZON")
 .load // load() returns a DataFrame

val selectfields1 = datasource0.select(from_json($"data".cast("string"),
glueContext.getCatalogSchemaAsSparkSchema("stream-demos", "stream-join-demo2")) as
"data").select("data.*")

val datasink2 = selectfields1.writeStream.foreachBatch { (dataFrame: Dataset[Row],
batchId: Long) => { //foreachBatch() returns type DataStreamWriter
 val joined = dataFrame.join(staticData, "product_id")
 val year: Int = Calendar.getInstance().get(Calendar.YEAR)
 val month :Int = Calendar.getInstance().get(Calendar.MONTH) + 1
 val day: Int = Calendar.getInstance().get(Calendar.DATE)
 val hour: Int = Calendar.getInstance().get(Calendar.HOUR_OF_DAY)

 if (dataFrame.count() > 0) {
 joined.write // joined.write returns type
DataFrameWriter
 .mode(SaveMode.Append)
 .format("parquet")
 .option("quote", " ")
 .save("s3://awsexamplebucket-streaming-demo2/output/" + "/year=" +
"%04d".format(year) + "/month=" + "%02d".format(month) + "/day=" + "%02d".format(day)
+ "/hour=" + "%02d".format(hour) + "/"")
 }
}
}

```

```
 } // end foreachBatch()
 .trigger(Trigger.ProcessingTime("100 seconds"))
 .option("checkpointLocation", "s3://awsexamplebucket-streaming-demo2/
checkpoint/")
 .start().awaitTermination() // start() returns type StreamingQuery
 Job.commit()
 }
}
```

## AWS Glue Scala ライブラリの API

AWS Glue では、抽出、変換、ロード (ETL) ジョブのスクリプト化に使用する PySpark Scala 言語の拡張機能がサポートされています。以下のセクションでは、AWS Glue Scala ライブラリの API について説明します。

com.amazonaws.services.glue

Scala ライブラリの com.amazonaws.services.glueAWS Glue パッケージには、以下の API が含まれています。

- [ChoiceOption](#)
- [DataSink](#)
- [DataSource 特性](#)
- [DynamicFrame](#)
- [DynamicRecord](#)
- [GlueContext](#)
- [MappingSpec](#)
- [ResolveSpec](#)

com.amazonaws.services.glue.ml

AWS Glue Scala ライブラリの com.amazonaws.services.glue.ml パッケージには、以下の API が含まれています。

- [FillMissingValues](#)
- [FindIncrementalMatches](#)
- [FindMatches](#)

com.amazonaws.services.glue.dq

AWS Glue Scala ライブラリの com.amazonaws.services.glue.dq パッケージには、以下の API が含まれています。

- [EvaluateDataQuality](#)

com.amazonaws.services.glue.types

Scala ライブラリの com.amazonaws.services.glue.typesAWS Glue パッケージには、以下の API が含まれています。

- [ArrayNode](#)
- [BinaryNode](#)
- [BooleanNode](#)
- [ByteNode](#)
- [DateNode](#)
- [DecimalNode](#)
- [DoubleNode](#)
- [DynamicNode](#)
- [FloatNode](#)
- [IntegerNode](#)
- [LongNode](#)
- [MapLikeNode](#)
- [MapNode](#)
- [NullNode](#)
- [ObjectNode](#)
- [ScalarNode](#)
- [ShortNode](#)
- [StringNode](#)
- [TimestampNode](#)

com.amazonaws.services.glue.util

Scala ライブラリの com.amazonaws.services.glue.utilAWS Glue パッケージには、以下の API が含まれています。

- [GlueArgParser](#)
- [ジョブ](#)

AWS Glue Scala ChoiceOption API

トピック

- [ChoiceOption 特性](#)
- [ChoiceOption オブジェクト](#)
- [ケースクラス ChoiceOptionWithResolver](#)
- [ケースクラス MatchCatalogSchemaChoiceOption](#)

パッケージ: com.amazonaws.services.glue

ChoiceOption 特性

```
trait ChoiceOption extends Serializable
```

ChoiceOption オブジェクト

ChoiceOption

```
object ChoiceOption
```

DynamicFrame のすべての ChoiceType ノードに適用可能な選択肢を解決するための一般的な戦略。

- val CAST
- val MAKE\_COLS
- val MAKE\_STRUCT
- val MATCH\_CATALOG

- `val PROJECT`

## Def apply

```
def apply(choice: String): ChoiceOption
```

## ケースクラス ChoiceOptionWithResolver

```
case class ChoiceOptionWithResolver(name: String, choiceResolver: ChoiceResolver)
 extends ChoiceOption {}
```

## ケースクラス MatchCatalogSchemaChoiceOption

```
case class MatchCatalogSchemaChoiceOption() extends ChoiceOption {}
```

## Abstract DataSink クラス

### トピック

- [Def writeDynamicFrame](#)
- [Def pyWriteDynamicFrame](#)
- [Def writeDataFrame](#)
- [Def pyWriteDataFrame](#)
- [Def setCatalogInfo](#)
- [Def supportsFormat](#)
- [Def setFormat](#)
- [Def withFormat](#)
- [Def setAccumulableSize](#)
- [Def getOutputErrorRecordsAccumulable](#)
- [Def errorsAsDynamicFrame](#)
- [DataSink オブジェクト](#)

パッケージ: `com.amazonaws.services.glue`

```
abstract class DataSink
```

ライターは DataSource に似ています。DataSink は、DynamicFrame が書き込まれる送信先と形式をカプセル化します。

### Def writeDynamicFrame

```
def writeDynamicFrame(frame : DynamicFrame,
 callSite : CallSite = CallSite("Not provided", "")
) : DynamicFrame
```

### Def pyWriteDynamicFrame

```
def pyWriteDynamicFrame(frame : DynamicFrame,
 site : String = "Not provided",
 info : String = "")
```

### Def writeDataFrame

```
def writeDataFrame(frame: DataFrame,
 glueContext: GlueContext,
 callSite: CallSite = CallSite("Not provided", ""))
): DataFrame
```

### Def pyWriteDataFrame

```
def pyWriteDataFrame(frame: DataFrame,
 glueContext: GlueContext,
 site: String = "Not provided",
 info: String = ""
): DataFrame
```

### Def setCatalogInfo

```
def setCatalogInfo(catalogDatabase: String,
 catalogTableName : String,
```

```
catalogId : String = "")
```

### Def supportsFormat

```
def supportsFormat(format : String) : Boolean
```

### Def setFormat

```
def setFormat(format : String,
 options : JsonOptions
) : Unit
```

### Def withFormat

```
def withFormat(format : String,
 options : JsonOptions = JsonOptions.empty
) : DataSink
```

### Def setAccumulableSize

```
def setAccumulableSize(size : Int) : Unit
```

### Def getOutputErrorRecordsAccumulable

```
def getOutputErrorRecordsAccumulable : Accumulable[List[OutputError], OutputError]
```

### Def errorsAsDynamicFrame

```
def errorsAsDynamicFrame : DynamicFrame
```

### DataSink オブジェクト

```
object DataSink
```

## Def recordMetrics

```
def recordMetrics(frame : DynamicFrame,
 ctxt : String
) : DynamicFrame
```

## AWS Glue Scala DataSource トレイト

パッケージ: com.amazonaws.services.glue

DynamicFrame を生成するための高水準インターフェイス。

```
trait DataSource {

 def getDynamicFrame : DynamicFrame

 def getDynamicFrame(minPartitions : Int,
 targetPartitions : Int
) : DynamicFrame
 def getDataFrame : DataFrame

 /** @param num: the number of records for sampling.
 * @param options: optional parameters to control sampling behavior. Current
 available parameter for Amazon S3 sources in options:
 * 1. maxSamplePartitions: the maximum number of partitions the sampling will
 read.
 * 2. maxSampleFilesPerPartition: the maximum number of files the sampling will
 read in one partition.
 */
 def getSampleDynamicFrame(num:Int, options: JsonOptions = JsonOptions.empty):
 DynamicFrame

 def glueContext : GlueContext

 def setFormat(format : String,
 options : String
) : Unit

 def setFormat(format : String,
 options : JsonOptions
) : Unit

 def supportsFormat(format : String) : Boolean
```

```
def withFormat(format : String,
 options : JsonObject = JsonObject.empty
) : DataSource
}
```

## AWS Glue Scala DynamicFrame API

パッケージ: `com.amazonaws.services.glue`

### 目次

- [AWS Glue Scala DynamicFrame クラス](#)

- [Val errorsCount](#)
- [Def applyMapping](#)
- [Def assertErrorThreshold](#)
- [Def count](#)
- [Def dropField](#)
- [Def dropFields](#)
- [Def dropNulls](#)
- [Def errorsAsDynamicFrame](#)
- [Def filter](#)
- [Def getName](#)
- [Def getNumPartitions](#)
- [Def getSchemaIfComputed](#)
- [Def isSchemaComputed](#)
- [Def javaToPython](#)
- [Def join](#)
- [Def map](#)
- [Def mergeDynamicFrames](#)
- [Def printSchema](#)
- [Def recomputeSchema](#)
- [Def relationalize](#)
- [Def renameField](#)
- [Def repartition](#)

- [Def resolveChoice](#)
- [Def schema](#)
- [Def selectField](#)
- [Def selectFields](#)
- [Def show](#)
- [Def simplifyDDBJson](#)
- [Def spigot](#)
- [Def splitFields](#)
- [Def splitRows](#)
- [Def stageErrorsCount](#)
- [Def toDF](#)
- [Def unbox](#)
- [Def unnest](#)
- [Def unnestDDBJson](#)
- [Def withFrameSchema](#)
- [Def withName](#)
- [Def withTransformationContext](#)
- [DynamicFrame オブジェクト](#)
  - [Def apply](#)
  - [Def emptyDynamicFrame](#)
  - [Def fromPythonRDD](#)
  - [Def ignoreErrors](#)
  - [Def inlineErrors](#)
  - [Def newFrameWithErrors](#)

## AWS Glue Scala DynamicFrame クラス

パッケージ: com.amazonaws.services.glue

```
class DynamicFrame extends Serializable with Logging (
 val glueContext : GlueContext,
 _records : RDD[DynamicRecord],
 val name : String = s"",
```

```
val transformationContext : String = DynamicFrame.UNDEFINED,
callSite : CallSite = CallSite("Not provided", ""),
stageThreshold : Long = 0,
totalThreshold : Long = 0,
prevErrors : => Long = 0,
errorExpr : => Unit = {})
```

DynamicFrame は、自己記述型の [DynamicRecord](#) オブジェクトの分散コレクションです。

DynamicFrame は、ETL (抽出、変換、ロード) オペレーションの柔軟なデータモデルを提供するように設計されています。これらのオブジェクトを作成するのにスキーマは必要なく、乱雑または不整合な値や型を持つデータの読み取りと変換に使用できます。スキーマは、スキーマを必要とするオペレーションでオンデマンドで計算できます。

DynamicFrame は、データクリーニングと ETL 用の広範な変換を提供します。また、既存のコードと統合するための SparkSQL DataFrames との相互変換や、DataFrames が提供する多くの分析オペレーションをサポートしています。

以下のパラメータは、DynamicFrame を生成する AWS Glue 変換の多くで共有されます。

- `transformationContext` - この DynamicFrame の識別子。実行間で保持されるジョブのブックマーク状態のキーとして、`transformationContext` が使用されます。
- `callSite` - エラーレポートのコンテキスト情報を提供します。これらの値は、Python から呼び出すときに、自動的に設定されます。
- `stageThreshold` - この DynamicFrame の計算から例外がスローされるまでのエラーレコードの最大許容数。以前の DynamicFrame にあるレコードは除きます。
- `totalThreshold` - 例外がスローされるまでの合計エラーレコードの最大数。以前のフレームのレコードも含まれます。

## Val errorsCount

```
val errorsCount
```

この DynamicFrame のエラーレコードの数。以前のオペレーションのエラーも含まれます。

## Def applyMapping

```
def applyMapping(mappings : Seq[Product4[String, String, String, String]],
 caseSensitive : Boolean = true,
```

```

transformationContext : String = "",
callSite : CallSite = CallSite("Not provided", ""),
stageThreshold : Long = 0,
totalThreshold : Long = 0
) : DynamicFrame

```

- mappings - 新しい DynamicFrame を生成するための一連のマッピング。
- caseSensitive - ソース列で大文字と小文字を区別して扱うかどうかを指定。このパラメータを false に設定すると、AWS Glue Data Catalog のように大文字と小文字を区別しないストアと統合する場合に役立ちます。

マッピングのシーケンスに基づく選択列、プロジェクト列、およびキャスト列。

各マッピングは、ソース列/タイプとターゲット列/タイプで構成されます。マッピングは、4 タプル (source\_path、source\_type、 target\_path、target\_type) として指定するか、同じ情報を含む [MappingSpec](#) オブジェクトとして指定できます。

マッピングは、シンプルなプロジェクションやキャストに使用できるだけでなく、パスのコンポーネントを "." (ピリオド) で区切ることでフィールドのネストまたはネスト解除にも使用できます。

例えば、以下のスキーマを含む DynamicFrame があるとします。

```

{{{
 root
 |-- name: string
 |-- age: int
 |-- address: struct
 | |-- state: string
 | |-- zip: int
 }}}

```

以下の呼び出しを行って、state フィールドと zip フィールドをネスト解除できます。

```

{{{
 df.applyMapping(
 Seq(("name", "string", "name", "string"),
 ("age", "int", "age", "int"),
 ("address.state", "string", "state", "string"),
 ("address.zip", "int", "zip", "int")))
 }}}

```

その結果、スキーマは以下ようになります。

```
{{{
 root
 |-- name: string
 |-- age: int
 |-- state: string
 |-- zip: int
}}}
```

`applyMapping` を使用して列を再ネストすることもできます。以下の例では、前の変換を反転し、ターゲットに `address` という名前の構造体を作成します。

```
{{{
 df.applyMapping(
 Seq(("name", "string", "name", "string"),
 ("age", "int", "age", "int"),
 ("state", "string", "address.state", "string"),
 ("zip", "int", "address.zip", "int")))
}}}
```

"." (ピリオド) 文字を含むフィールド名は、バッククォート (` `) で囲むことで使用できます。

#### Note

現在、`applyMapping` メソッドを使用して、配列の下にネストされた列をマップすることはできません。

## Def `assertErrorThreshold`

```
def assertErrorThreshold : Unit
```

このアクションでは、計算を適用し、エラーレコード数が `stageThreshold` と `totalThreshold` を下回っていることを確認します。いずれかの条件が失敗すると、例外をスローします。

## Def `count`

```
lazy
def count
```

この `DynamicFrame` の要素数を返します。

### Def dropField

```
def dropField(path : String,
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided", ""),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

指定した列が削除された新しい `DynamicFrame` を返します。

### Def dropFields

```
def dropFields(fieldNames : Seq[String], // The column names to drop.
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided", ""),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

指定した列が削除された新しい `DynamicFrame` を返します。

このメソッドを使用して、ネストされた列 (配列内の列など) を削除できますが、特定の配列要素を削除することはできません。

### Def dropNulls

```
def dropNulls(transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided", ""),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0)
```

すべての `NULL` 列が削除された新しい `DynamicFrame` を返します。

#### Note

`NullType` タイプの列のみが削除されます。他の列にある個別の `null` 値は削除または変更されません。

## Def errorsAsDynamicFrame

```
def errorsAsDynamicFrame
```

この DynamicFrame からのエラーレコードを含む新しい DynamicFrame を返します。

## Def filter

```
def filter(f : DynamicRecord => Boolean,
 errorMsg : String = "",
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided"),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

"f" 関数が true を返すレコードのみを含む新しい DynamicFrame を生成します。フィルター関数 "f" は、入力レコードを変更しないものとします。

## Def getName

```
def getName : String
```

この DynamicFrame の名前を返します。

## Def getNumPartitions

```
def getNumPartitions
```

この DynamicFrame のパーティション数を返します。

## Def getSchemaIfComputed

```
def getSchemaIfComputed : Option[Schema]
```

計算済みのスキーマを返します。スキーマが計算済みでない場合は、データをスキャンしません。

## Def isSchemaComputed

```
def isSchemaComputed : Boolean
```

この `DynamicFrame` に対してスキーマが計算された場合は `true` を返し、それ以外の場合は `false` を返します。このメソッドが `false` を返した場合、`schema` メソッドを呼び出すには、この `DynamicFrame` のレコードを再び渡す必要があります。

### Def javaToPython

```
def javaToPython : JavaRDD[Array[Byte]]
```

### Def join

```
def join(keys1 : Seq[String],
 keys2 : Seq[String],
 frame2 : DynamicFrame,
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided", ""),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

- `keys1` - 結合に使用するこの `DynamicFrame` の列。
- `keys2` - 結合に使用する `frame2` の列。 `keys1` と同じ長さにする必要があります。
- `frame2` - 結合の対象になる `DynamicFrame`。

指定したキーを使用して `frame2` との等結合を行った結果を返します。

### Def map

```
def map(f : DynamicRecord => DynamicRecord,
 errorMsg : String = "",
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided", ""),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

指定した関数 "f" をこの `DynamicFrame` の各レコードに適用することで生成された新しい `DynamicFrame` を返します。

このメソッドは、指定した関数を適用する前に各レコードをコピーするため、レコードを安全に変更できます。特定のレコードでマッピング関数から例外がスローされた場合、そのレコードはエラーとしてマークされ、スタックトレースがエラーレコードの列として保存されます。

## Def mergeDynamicFrames

```
def mergeDynamicFrames(stageDynamicFrame: DynamicFrame, primaryKeys: Seq[String],
 transformationContext: String = "",
 options: JsonOptions = JsonOptions.empty, callSite: CallSite =
 CallSite("Not provided"),
 stageThreshold: Long = 0, totalThreshold: Long = 0):
 DynamicFrame
```

- `stageDynamicFrame` – マージするステージング `DynamicFrame`。
- `primaryKeys` – ソースおよびステージング `DynamicFrame` からのレコードを照合するプライマリキーフィールドのリスト。
- `transformationContext` – 現在の変換に関するメタデータを取得するために使用される一意の文字列 (オプション)。
- `options` – この変換に関する追加情報を提供する JSON の名前と値のペアを示す文字列。
- `callSite` – エラーレポートのコンテキスト情報を提供するために使用します。
- `stageThreshold` – A Long。指定された変換で処理がエラーアウトする必要があるエラーの数。
- `totalThreshold` – A Long。この変換までで処理がエラーアウトする必要があるエラーの合計数。

レコードを識別するために、この `DynamicFrame` を指定されたプライマリキーに基づくステージング `DynamicFrame` とマージします。重複レコード (同じプライマリキーを持つレコード) は重複除外されません。ステージングフレームに一致するレコードがない場合、すべてのレコード (重複を含む) がソースから保持されます。ステージングフレームに一致するレコードがある場合、ステージングフレームのレコードによって、AWS Glue のソースのレコードが上書きされます。

以下の場合、返される `DynamicFrame` にはレコード A が含まれます。

1. A がソースフレームとステージングフレームの両方に存在する場合、ステージングフレームの A が返されます。
2. A がソーステーブルにあり、`A.primaryKeys` が `stagingDynamicFrame` がない場合 (これは、ステージングテーブルで A が更新されていないことを意味します)。

ソースフレームとステージングフレームが、同じスキーマを持つ必要はありません。

### Example

```
val mergedFrame: DynamicFrame = srcFrame.mergeDynamicFrames(stageFrame, Seq("id1", "id2"))
```

### Def printSchema

```
def printSchema : Unit
```

この DynamicFrame のスキーマを、人間が判読できる形式で stdout に出力します。

### Def recomputeSchema

```
def recomputeSchema : Schema
```

スキーマの再計算を強制します。これには、データのスキャンが必要ですが、現在のスキーマの一部のフィールドがデータにない場合、スキーマが "強化" される場合があります。

再計算されたスキーマを返します。

### Def relationalize

```
def relationalize(rootTableName : String,
 stagingPath : String,
 options : JsonObject = JsonObject.empty,
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided"),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : Seq[DynamicFrame]
```

- `rootTableName` – 出力の基本 DynamicFrame に使用する名前。ピボット配列によって作成される DynamicFrame は、この名前をプレフィックスとして使用します。
- `stagingPath` – 中間データを書き込む Amazon Simple Storage Service (Amazon S3) のパス。
- `options` – Relationalize のオプションと設定。現在使用されていません。

すべてのネストされた構造をフラット化し、配列を個別のテーブルにピボットします。

このオペレーションでは、リレーショナルデータベースに取り込むための深くネストされたデータを準備できます。ネストされた構造体は、[ネスト解除](#) 変換と同じ方法でフラット化されます。さらに、配列は個別のテーブルにピボットされ、各配列要素が行になります。例えば、以下のデータを含む `DynamicFrame` があるとします。

```
{"name": "Nancy", "age": 47, "friends": ["Fred", "Lakshmi"]}
{"name": "Stephanie", "age": 28, "friends": ["Yao", "Phil", "Alvin"]}
{"name": "Nathan", "age": 54, "friends": ["Nicolai", "Karen"]}
```

以下のコードを実行します。

```
{{{
 df.relationalize("people", "s3:/my_bucket/my_path", JsonOptions.empty)
}}}
```

これにより、2つのテーブルが生成されます。最初のテーブルは "people" という名前で、内容は以下のとおりです。

```
{{{
 {"name": "Nancy", "age": 47, "friends": 1}
 {"name": "Stephanie", "age": 28, "friends": 2}
 {"name": "Nathan", "age": 54, "friends": 3}
}}}
```

ここで、友人の配列は自動生成された結合キーに置き換えられています。別のテーブルは `people.friends` という名前で、以下の内容で作成されます。

```
{{{
 {"id": 1, "index": 0, "val": "Fred"}
 {"id": 1, "index": 1, "val": "Lakshmi"}
 {"id": 2, "index": 0, "val": "Yao"}
 {"id": 2, "index": 1, "val": "Phil"}
 {"id": 2, "index": 2, "val": "Alvin"}
 {"id": 3, "index": 0, "val": "Nicolai"}
 {"id": 3, "index": 1, "val": "Karen"}
}}}
```

このテーブルで、"id" は配列要素の元のレコードを識別する結合キーです。"index" は元の配列内の位置を参照します。"val" は実際の配列エントリです。

`relationalize` メソッドは、このプロセスをすべての配列に再帰的に適用することで作成した `DynamicFrame` のシーケンスを返します。

#### Note

AWS Glue ライブラリは、新しいテーブルの結合キーを自動的に生成します。ジョブの実行間で結合キーが必ず一意になるように、ジョブのブックマークを有効にする必要があります。

## Def renameField

```
def renameField(oldName : String,
 newName : String,
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided", ""),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

- `oldName` – 列の元の名前。
- `newName` – 列の新しい名前。

指定したフィールドの名前が変更された新しい `DynamicFrame` を返します。

このメソッドを使用して、ネストされたフィールドの名前を変更できます。例えば、以下のコードはアドレス構造体内の `state` の名前を `state_code` に変更します。

```
{{{
 df.renameField("address.state", "address.state_code")
}}}
```

## Def repartition

```
def repartition(numPartitions : Int,
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided", ""),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

```
) : DynamicFrame
```

numPartitions パーティションを含む新しい DynamicFrame を返します。

## Def resolveChoice

```
def resolveChoice(specs : Seq[Product2[String, String]] = Seq.empty[ResolveSpec],
 choiceOption : Option[ChoiceOption] = None,
 database : Option[String] = None,
 tableName : Option[String] = None,
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided", ""),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

- choiceOption - 仕様シーケンスに列挙されていない ChoiceType 列に適用するアクション。
- database - match\_catalog アクションで使用する Data Catalog データベース。
- tableName - match\_catalog アクションで使用する Data Catalog テーブル。

1 つ以上の ChoiceType をより限定されたタイプに置き換えて新しい DynamicFrame を返します。

resolveChoice を使用するには 2 つの方法があります。最初の方法では、特定の列のシーケンスと解決方法を指定します。これらは (列、アクション) ペアで構成されたタプルとして指定します。

以下のアクションを指定できます。

- cast:type - すべての値を指定した型にキャストしようとします。
- make\_cols - それぞれの異なるタイプを columnName\_type という名前の列に変換します。
- make\_struct - 列を各区別型のキーを持つ構造体に変換します。
- project:type - 指定した型の値のみを保持します。

resolveChoice のもう 1 つのモードでは、すべての ChoiceType に対して単一の解決策を指定します。このモードは、実行前に ChoiceType の完全なリストが不明な場合に使用できます。このモードでは、上記のアクションに加えて、以下のアクションもサポートされています。

- `match_catalogChoiceType` – 指定したカタログテーブルの対応するタイプへの各のキャストを試行します。

例:

`user.id` 列を解決するために `int` にキャストし、`address` フィールドで構造体のみを保持します。

```
{{{
 df.resolveChoice(specs = Seq(("user.id", "cast:int"), ("address", "project:struct")))
}}}
```

すべての `ChoiceType` を解決するために、各選択肢を別個の列に変換します。

```
{{{
 df.resolveChoice(choiceOption = Some(ChoiceOption("make_cols")))
}}}
```

すべての `ChoiceType` を解決するために、指定したカタログテーブルのタイプにキャストします。

```
{{{
 df.resolveChoice(choiceOption = Some(ChoiceOption("match_catalog")),
 database = Some("my_database"),
 tableName = Some("my_table"))
}}}
```

## Def schema

```
def schema : Schema
```

この `DynamicFrame` のスキーマを返します。

返されたスキーマには必ず、この `DynamicFrame` のレコードにあるすべてのフィールドが含まれます。しかし、その他のフィールドが含まれる場合がまれにあります。[ネスト解除](#) メソッドを使用して、この `DynamicFrame` のレコードに基づいてスキーマを "強化" できます。

## Def selectField

```
def selectField(fieldName : String,
```

```
transformationContext : String = "",
callSite : CallSite = CallSite("Not provided", ""),
stageThreshold : Long = 0,
totalThreshold : Long = 0
) : DynamicFrame
```

単一のフィールドを DynamicFrame として返します。

### Def selectFields

```
def selectFields(paths : Seq[String],
transformationContext : String = "",
callSite : CallSite = CallSite("Not provided", ""),
stageThreshold : Long = 0,
totalThreshold : Long = 0
) : DynamicFrame
```

- paths - 選択する列名のシーケンス。

指定した列を含む新しい DynamicFrame を返します。

#### Note

selectFields メソッドでは、最上位の列のみを選択できます。[applyMapping](#) メソッドでは、ネストされた列を選択できます。

### Def show

```
def show(numRows : Int = 20) : Unit
```

- numRows - 出力する行の数。

この DynamicFrame の行を JSON 形式で出力します。

### Def simplifyDDBJson

DynamoDB は AWS Glue DynamoDB エクスポートコネクタを使用してエクスポートし、特殊なネスト構造の JSON ファイルが生成されます。詳細については、「[データオブジェクト](#)」を参照して

ください。simplifyDDBJsonこのタイプのデータの DynamicFrame にあるネスト化された列を単純化し、新しい単純化された DynamicFrame を返します。リストタイプに複数タイプまたは Map タイプが含まれている場合、リストの要素は単純化されません。このメソッドは、DynamoDB エクスポート JSON 形式のデータのみをサポートします。unnest が他の種類のデータにも同様の変更を加えることを考慮してください。

```
def simplifyDDBJson() : DynamicFrame
```

このメソッドにはパラメータはありません。

## 入力例

DynamoDB エクスポートによって生成された次のスキーマを考えてみましょう。

```
root
|-- Item: struct
| |-- parentMap: struct
| | |-- M: struct
| | | |-- childMap: struct
| | | | |-- M: struct
| | | | | |-- appName: struct
| | | | | | |-- S: string
| | | | | | |-- packageName: struct
| | | | | | | |-- S: string
| | | | | | | |-- updatedAt: struct
| | | | | | | | |-- N: string
| | |-- strings: struct
| | | |-- SS: array
| | | | |-- element: string
| | |-- numbers: struct
| | | |-- NS: array
| | | | |-- element: string
| | |-- binaries: struct
| | | |-- BS: array
| | | | |-- element: string
| | |-- isDDBJson: struct
| | | |-- BOOL: boolean
| | |-- nullValue: struct
| | | |-- NULL: boolean
```

## コードの例

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.DynamoDbDataSink
import org.apache.spark.SparkContextimport scala.collection.JavaConverters._

object GlueApp {

 def main(sysArgs: Array[String]): Unit = {
 val glueContext = new GlueContext(SparkContext.getOrCreate())
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)

 val dynamicFrame = glueContext.getSourceWithFormat(
 connectionType = "dynamodb",
 options = JsonOptions(Map(
 "dynamodb.export" -> "ddb",
 "dynamodb.tableArn" -> "ddbTableARN",
 "dynamodb.s3.bucket" -> "exportBucketLocation",
 "dynamodb.s3.prefix" -> "exportBucketPrefix",
 "dynamodb.s3.bucketOwner" -> "exportBucketAccountID",
))
).getDynamicFrame()

 val simplified = dynamicFrame.simplifyDDBJson()
 simplified.printSchema()

 Job.commit()
 }
}
```

## 出力例

この `simplifyDDBJson` 変換により、以下のように簡略化されます。

```
root
|-- parentMap: struct
| |-- childMap: struct
| | |-- appName: string
| | |-- packageName: string
| | |-- updatedAt: string
```

```
|-- strings: array
| |-- element: string
|-- numbers: array
| |-- element: string
|-- binaries: array
| |-- element: string
|-- isDDBJson: boolean
|-- nullValue: null
```

## Def spigot

```
def spigot(path : String,
 options : JsonOptions = new JsonOptions("{}"),
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided"),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

同じレコードを返すが、副作用としてレコードのサブセットを書き出すパススルー変換。

- path – 出力を `s3://bucket//path` 形式で書き込む先の Amazon S3 のパス。
- options - サンプルング動作を記述するオプションの `JsonOptions` マップ。

このレコードと同じレコードを含む `DynamicFrame` を返します。

デフォルトでは、path で指定した場所に任意の 100 レコードを書き込みます。この動作は、options マップを使用してカスタマイズできます。有効なキーは以下のとおりです。

- topk – 書き出されるレコードの総数を指定します。デフォルトは 100 です。
- prob – 各レコードを含める確率を (小数として) 指定します。デフォルトは1です。

例えば、以下の呼び出しでは、データセットをサンプルングするために、20 パーセントの確率で各レコードを選択し、200 レコードを書き出した後で停止します。

```
{{{
 df.spigot("s3://my_bucket/my_path", JsonOptions(Map("topk" -> 200, "prob" ->
 0.2)))
}}}
```

## Def splitFields

```
def splitFields(paths : Seq[String],
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided", ""),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : Seq[DynamicFrame]
```

- paths - 最初の DynamicFrame に含めるパス。

2つの DynamicFrame のシーケンスを返します。DynamicFrame の1つ目には指定したパスが含まれ、2つ目には他のすべての列が含まれます。

### 例

この例では、persons Glue Data Catalog の legislators データベースの AWS テーブルから作成した DynamicFrame を2つに分割し、指定したフィールドを1つ目の DynamicFrame に、残りのフィールドを2つ目の DynamicFrame に配置しています。次に、結果から最初の DynamicFrame を選択します。

```
val InputFrame = glueContext.getCatalogSource(database="legislators",
 tableName="persons",
 transformationContext="InputFrame").getDynamicFrame()

val SplitField_collection = InputFrame.splitFields(paths=Seq("family_name", "name",
 "links.note",
 "links.url", "gender", "image", "identifiers.scheme", "identifiers.identifier",
 "other_names.lang",
 "other_names.note", "other_names.name"), transformationContext="SplitField_collection")

val ResultFrame = SplitField_collection(0)
```

## Def splitRows

```
def splitRows(paths : Seq[String],
 values : Seq[Any],
 operators : Seq[String],
 transformationContext : String,
 callSite : CallSite,
 stageThreshold : Long,
```

```
totalThreshold : Long
) : Seq[DynamicFrame]
```

列と定数を比較する述語に基づいて行を分割します。

- `paths` - 比較に使用する列。
- `values` - 比較に使用する定数値。
- `operators` - 比較に使用する演算子。

2つの `DynamicFrame` のシーケンスを返します。1つ目には述語が `true` の行を含め、2つ目には述語が `false` の行を含めます。

述語を指定するには3つのシーケンスを使用します。"paths"には、ネストされている可能性がある列の名前を含め、"values"には、比較する定数値を含め、"operators"には、比較に使用する演算子を含めます。3つすべてのシーケンスを同じ長さにする必要があります。n演算子では、n番目の列をn番目の値と比較します。

各演算子は `"!="`、`"="`、`"<="`、`"<"`、`">="`、`">"` のいずれかにする必要があります。

以下の呼び出しの例では、`DynamicFrame` を分割し、1つ目の出力フレームには米国の65才を超える人々のレコード、2つ目には他のすべてのレコードを含めています。

```
{{{
 df.splitRows(Seq("age", "address.country"), Seq(65, "USA"), Seq(">=", "="))
}}}
```

### Def stageErrorsCount

```
def stageErrorsCount
```

この `DynamicFrame` の計算中に生じたエラーレコードの数を返します。この `DynamicFrame` に入力として渡した以前のオペレーションのエラーは含まれません。

### Def toDF

```
def toDF(specs : Seq[ResolveSpec] = Seq.empty[ResolveSpec]) : DataFrame
```

この `DynamicFrame` を、同じスキーマとレコードを含む Apache Spark SQL `DataFrame` に変換します。

**Note**

DataFrame は ChoiceType をサポートしていないため、このメソッドは ChoiceType 列を StructType に自動的に変換します。選択肢の解決の詳細とオプションについては、「[resolveChoice](#)」を参照してください。

**Def unbox**

```
def unbox(path : String,
 format : String,
 optionString : String = "{}",
 transformationContext : String = "",
 callSite : CallSite = CallSite("Not provided"),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

- path – 解析する列。文字列またはバイナリにする必要があります。
- format – 解析に使用する形式。
- optionString – 形式に渡すオプション (CSV 区切り記号など)。

指定した形式に従って、埋め込まれた文字列またはバイナリ列を解析します。解析された列は、元の列名で構造体の下にネストされます。

例えば、JSON 列が埋め込まれた CSV ファイルがあるとします。

```
name, age, address
Sally, 36, {"state": "NE", "city": "Omaha"}
...
```

最初の解析の後、以下のスキーマを含む DynamicFrame を取得します。

```
{{
 root
 |-- name: string
 |-- age: int
 |-- address: string
}}
```

アドレス列で `unbox` を呼び出して、特定のコンポーネントを解析できます。

```
{{{
 df.unbox("address", "json")
}}}
```

これにより、以下のスキーマを含む `DynamicFrame` を取得します。

```
{{{
 root
 |-- name: string
 |-- age: int
 |-- address: struct
 | |-- state: string
 | |-- city: string
}}}
```

## Def unnest

```
def unnest(transformationContext : String = "",
 callSite : CallSite = CallSite("Not Provided"),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0
) : DynamicFrame
```

すべてのネストされた構造体が平坦化された新しい `DynamicFrame` を返します。名前は "." (ピリオド) 文字を使用して生成されます。

例えば、以下のスキーマを含む `DynamicFrame` があるとします。

```
{{{
 root
 |-- name: string
 |-- age: int
 |-- address: struct
 | |-- state: string
 | |-- city: string
}}}
```

次の呼び出しでアドレス構造体がネスト解除されます。

```

{{{
 df.unnest()
}}}

```

その結果、スキーマは以下のようになります。

```

{{{
 root
 |-- name: string
 |-- age: int
 |-- address.state: string
 |-- address.city: string
}}}

```

このメソッドはまた、配列内でネストされた構造体をネスト解除しません。しかしこれまでの慣習上、そのようなフィールドの名前の前には、外側の配列の名前と ".val" が付けられます。

#### Def unnestDDBJson

```

unnestDDBJson(transformationContext : String = "",
 callSite : CallSite = CallSite("Not Provided"),
 stageThreshold : Long = 0,
 totalThreshold : Long = 0): DynamicFrame

```

特に DynamoDB JSON 構造体内にある DynamicFrame でネスト化された列をネスト解除すると、新しくネスト解除された DynamicFrame が返されます。構造体型の配列のある列は、ネスト解除されません。これは、通常の unnest 変換とは異なる特殊なタイプのネスト解除変換で、データが DynamoDB JSON 構造体に格納されている必要があることに注意してください。詳細については、「[DynamoDB JSON](#)」を参照してください。

例えば、DynamoDB JSON 構造体のあるエクスポートを読み取るスキーマは次のようになります。

```

root
|-- Item: struct
| |-- ColA: struct
| | |-- S: string
| |-- ColB: struct
| | |-- S: string
| |-- ColC: struct
| | |-- N: string

```

```
| |-- ColD: struct
| | |-- L: array
| | | |-- element: null
```

`unnestDDBJson()` 変換は、以下のように変換します。

```
root
|-- ColA: string
|-- ColB: string
|-- ColC: string
|-- ColD: array
| |-- element: null
```

次のコード例は、AWS Glue DynamoDB エクスポートコネクタを使用し、DynamoDB JSON `unnest` を呼び出し、パーティションの数を表示します。

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.DynamoDbDataSink
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._

object GlueApp {

 def main(sysArgs: Array[String]): Unit = {
 val glueContext = new GlueContext(SparkContext.getOrCreate())
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)

 val dynamicFrame = glueContext.getSourceWithFormat(
 connectionType = "dynamodb",
 options = JsonOptions(Map(
 "dynamodb.export" -> "ddb",
 "dynamodb.tableArn" -> "<test_source>",
 "dynamodb.s3.bucket" -> "<bucket name>",
 "dynamodb.s3.prefix" -> "<bucket prefix>",
 "dynamodb.s3.bucketOwner" -> "<account_id of bucket>",
))
).getDynamicFrame()
```

```
val unnested = dynamicFrame.unnestDDBJson()
print(unnested.getNumPartitions())

Job.commit()
}
}
```

## Def withFrameSchema

```
def withFrameSchema(getSchema : () => Schema) : DynamicFrame
```

- getSchema – 使用するスキーマを返す関数。高価である可能性が高い計算を延期するために、パラメータがゼロである関数として指定します。

この DynamicFrame のスキーマを、指定した値に設定します。高価なスキーマの再計算を回避するために主に内部で使用されます。渡すスキーマには、データ内に存在するすべての列を含める必要があります。

## Def withName

```
def withName(name : String) : DynamicFrame
```

- name - 使用する新しい名前。

この DynamicFrame のコピーを新しい名前で見返します。

## Def withTransformationContext

```
def withTransformationContext(ctx : String) : DynamicFrame
```

この DynamicFrame のコピーを、指定した変換コンテキストで見返します。

## DynamicFrame オブジェクト

パッケージ: com.amazonaws.services.glue

```
object DynamicFrame
```

## Def apply

```
def apply(df : DataFrame,
 glueContext : GlueContext
) : DynamicFrame
```

## Def emptyDynamicFrame

```
def emptyDynamicFrame(glueContext : GlueContext) : DynamicFrame
```

## Def fromPythonRDD

```
def fromPythonRDD(rdd : JavaRDD[Array[Byte]],
 glueContext : GlueContext
) : DynamicFrame
```

## Def ignoreErrors

```
def ignoreErrors(fn : DynamicRecord => DynamicRecord) : DynamicRecord
```

## Def inlineErrors

```
def inlineErrors(msg : String,
 callSite : CallSite
) : (DynamicRecord => DynamicRecord)
```

## Def newFrameWithErrors

```
def newFrameWithErrors(prevFrame : DynamicFrame,
 rdd : RDD[DynamicRecord],
 name : String = "",
 transformationContext : String = "",
 callSite : CallSite,
 stageThreshold : Long,
 totalThreshold : Long
) : DynamicFrame
```

## AWS Glue Scala DynamicRecord クラス

### トピック

- [def addField](#)
- [def dropField](#)
- [def setError](#)
- [def isError](#)
- [def getError](#)
- [def clearError](#)
- [def write](#)
- [def readFields](#)
- [def clone](#)
- [def schema](#)
- [def getRoot](#)
- [def toJson](#)
- [def getFieldNode](#)
- [def getField](#)
- [def hashCode](#)
- [def equals](#)
- [DynamicRecord オブジェクト](#)
- [RecordTraverser 特性](#)

パッケージ: com.amazonaws.services.glue

```
class DynamicRecord extends Serializable with Writable with Cloneable
```

DynamicRecord は、処理対象のデータセット内のデータ行を表す自己記述型のデータ構造体です。自己記述型とは、DynamicRecord が表す行のスキーマを、レコード自体を検査することで取得できるという意味です。DynamicRecord は Apache Spark の Row に似ています。

def addField

```
def addField(path : String,
```

```
dynamicNode : DynamicNode
) : Unit
```

指定したパスに [DynamicNode](#) を追加します。

- path - 追加するフィールドのパス。
- dynamicNode - 指定したパスに追加する [DynamicNode](#)。

def dropField

```
def dropField(path: String, underRename: Boolean = false): Option[DynamicNode]
```

指定したパスに配列がない場合は、指定したパスから [DynamicNode](#) を削除し、削除したノードを返します。

- path - 削除するフィールドへのパス。
- underRenamedropField – 名前変更の変換の一部として が呼び出された場合は true、それ以外の場合は false (デフォルトは false)。

scala.Option Option ([DynamicNode](#)) を返します。

def setError

```
def setError(error : Error)
```

error パラメータの指定に従って、このレコードをエラーレコードとして設定します。

戻り値は DynamicRecord。

def isError

```
def isError
```

このレコードがエラーレコードであるかどうかを確認します。

def getError

```
def getError
```

レコードがエラーレコードである場合、Error を受け取ります。このレコードがエラーレコードである場合は `scala.Some Some (エラー)` を返し、それ以外の場合は `scala.None` を返します。

def clearError

```
def clearError
```

Error を `scala.None.None` に設定します。

def write

```
override def write(out : DataOutput) : Unit
```

def readFields

```
override def readFields(in : DataInput) : Unit
```

def clone

```
override def clone : DynamicRecord
```

このレコードを新しい `DynamicRecord` に複製し、返します。

def schema

```
def schema
```

レコードを検査して `Schema` を取得します。

def getRoot

```
def getRoot : ObjectNode
```

レコードのルート `ObjectNode` を取得します。

def toJson

```
def toJson : String
```

レコードの JSON 文字列を取得します。

### def getFieldNode

```
def getFieldNode(path : String) : Option[DynamicNode]
```

DynamicNode のオプションとして指定した path でフィールドの値を取得します。

フィールドが存在する場合は scala.Some Some([DynamicNode](#)) を返し、それ以外の場合は scala.None.None を返します。

### def getField

```
def getField(path : String) : Option[Any]
```

DynamicNode のオプションとして指定した path でフィールドの値を取得します。

scala.Some Some (値) を返します。

### def hashCode

```
override def hashCode : Int
```

### def equals

```
override def equals(other : Any)
```

### DynamicRecord オブジェクト

```
object DynamicRecord
```

### def apply

```
def apply(row : Row,
 schema : SparkStructType)
```

メソッドを適用して Apache Spark SQL Row を [DynamicRecord](#) に変換します。

- row – Spark SQL Row。
- schema – 行の Schema。

戻り値は DynamicRecord。

### RecordTraverser 特性

```
trait RecordTraverser {
 def nullValue(): Unit
 def byteValue(value: Byte): Unit
 def binaryValue(value: Array[Byte]): Unit
 def booleanValue(value: Boolean): Unit
 def shortValue(value: Short) : Unit
 def intValue(value: Int) : Unit
 def longValue(value: Long) : Unit
 def floatValue(value: Float): Unit
 def doubleValue(value: Double): Unit
 def decimalValue(value: BigDecimal): Unit
 def stringValue(value: String): Unit
 def dateValue(value: Date): Unit
 def timestampValue(value: Timestamp): Unit
 def objectStart(length: Int): Unit
 def objectKey(key: String): Unit
 def objectEnd(): Unit
 def mapStart(length: Int): Unit
 def mapKey(key: String): Unit
 def mapEnd(): Unit
 def arrayStart(length: Int): Unit
 def arrayEnd(): Unit
}
```

### AWS Glueスカラ API GlueContext

パッケージ: com.amazonaws.services.glue

```
class GlueContext extends SQLContext(sc) (
 @transient val sc : SparkContext,
 val defaultSourcePartitioner : PartitioningStrategy)
```

GlueContext は、Amazon Simple Storage Service (Amazon S3) や、AWS Glue Data Catalog、JDBC その他との間で、[DynamicFrame](#) の読み取りと書き込みを行うためのエントリポ

イントです。このクラスが提供するユーティリティ関数によって [DataSource 特性](#) オブジェクトと [DataSink](#) オブジェクトが作成され、これらのオブジェクトを使用して DynamicFrame を読み書きできます。

また、ソースから作成されたパーティションの数がパーティションの最小しきい値 (デフォルトは 10) 未満の場合は、DynamicFrame で GlueContext を使用してパーティションのターゲット数 (デフォルトは 20) を設定することもできます。

### デフ addIngestionTime-カラム

```
def addIngestionTimeColumns(
 df : DataFrame,
 timeGranularity : String = "") : DataFrame
```

#### 入力 DataFrame への取り込み時間列

(`ingest_year`、`ingest_month`、`ingest_day`、`ingest_hour`、`ingest_minute`) を追加します。Amazon S3 のデータカタログテーブルをターゲットとして指定する場合、この関数は、AWS Glue により生成されたスクリプト内で自動的に生成されます。この関数は、出力テーブル上で、取り込み時間列があるパーティションを自動的に更新します。これにより、入力データに明示的な取り込み時間列を指定しなくても、取り込み時間において出力データの自動的なパーティション化が行えます。

- `dataFrame` – 取り込み時間列の追加先である `dataFrame`。
- `timeGranularity` – 時間列の詳細度。有効な値は「`day`」、「`hour`」、および「`minute`」です。例えば、関数に対し「`hour`」が渡された場合、元の `dataFrame` は「`ingest_year`」、「`ingest_month`」、「`ingest_day`」に加え「`ingest_hour`」の時間列を持つこととなります。

時間の詳細度列を追加した後、そのデータフレームを返します。

例：

```
glueContext.addIngestionTimeColumns(dataFrame, "hour")
```

### def createDataFrame FromOptions

```
def createDataFrameFromOptions(connectionType : String,
 connectionOptions : JsonOptions,
 transformationContext : String = "",
```

```
format : String = null,
formatOptions : JsonOptions = JsonOptions.empty
) : DataSource
```

指定された接続と形式で作成された DataFrame を返します。この関数は AWS Glue ストリーミングソースでのみ使用してください。

- `connectionType` – ストリーミング接続タイプ。有効な値は、`kinesis` および `kafka` です。
- `connectionOptions` – 接続オプション。これは、Kinesis と Kafka で異なります。各ストリーミングデータソースのすべての接続オプションの一覧は、[AWS Glue for Spark での ETL の接続タイプとオプション](#) で確認いただけます。ストリーミング接続オプションについては、以下の違いに注意してください。
  - Kinesis ストリーミングのソースには `streamARN`、`startingPosition`、`inferSchema`、および `classification` が必要です。
  - Kafka ストリーミングのソースには `connectionName`、`topicName`、`startingOffsets`、`inferSchema`、および `classification` が必要です。
- `transformationContext` – 使用する変換コンテキスト (オプション)。
- `format` – 形式の仕様 (オプション)。Amazon S3 や、複数のフォーマットをサポートする AWS Glue 接続の場合に使用します。サポートされる形式については、「[AWS Glue for Spark での入出力のデータ形式に関するオプション](#)」を参照してください。
- `formatOptions` – 指定された形式に関するオプション。サポートされる形式オプションについては、「[データ形式に関するオプション](#)」を参照してください。

Amazon Kinesis ストリーミングソースの例:

```
val data_frame_datasource0 =
glueContext.createDataFrameFromOptions(transformationContext = "datasource0",
connectionType = "kinesis",
connectionOptions = JsonOptions("""{"streamName": "example_stream", "startingPosition":
"TRIM_HORIZON", "inferSchema": "true", "classification": "json"}"""))
```

Kafka ストリーミングソースの例 :

```
val data_frame_datasource0 =
glueContext.createDataFrameFromOptions(transformationContext = "datasource0",
connectionType = "kafka",
```

```
connectionOptions = JsonOptions("""{"connectionName": "example_connection",
 "topicName": "example_topic", "startingPosition": "earliest", "inferSchema": "false",
 "classification": "json", "schema": "`column1` STRING, `column2` STRING}"""))
```

forEachBatch

### forEachBatch(frame, batch\_function, options)

ストリーミングソースから読み取られるすべてのマイクロバッチに渡される、batch\_function を適用します。

- frame— DataFrame 現在のマイクロバッチを含む。
- batch\_function – すべてのマイクロバッチに適用される関数。
- options – マイクロバッチの処理方法に関する情報を保持している、キーと値のペアの集合。以下のような必須オプションがあります。
  - windowSize – 各バッチの処理にかかる時間。
  - checkpointLocation – ストリーミング ETL ジョブ用に、チェックポイントが格納される場所。
  - batchMaxRetries – 失敗した場合にこのバッチを再試行する最大回数。デフォルト値は 3 です。このオプションは、Glue バージョン 2.0 以降でのみ設定可能です。

例:

```
glueContext.forEachBatch(data_frame_datasource0, (dataFrame: Dataset[Row], batchId:
Long) =>
 {
 if (dataFrame.count() > 0)
 {
 val datasource0 = DynamicFrame(glueContext.addIngestionTimeColumns(dataFrame,
"hour"), glueContext)
 // @type: DataSink
 // @args: [database = "tempdb", table_name = "fromoptionsoutput",
stream_batch_time = "100 seconds",
 // stream_checkpoint_location = "s3://from-options-testing-eu-central-1/
fromOptionsOutput/checkpoint/",
 // transformation_ctx = "datasink1"]
 // @return: datasink1
 // @inputs: [frame = datasource0]
 val options_datasink1 = JsonOptions(
```

```

 Map("partitionKeys" -> Seq("ingest_year", "ingest_month", "ingest_day",
"ingest_hour"),
 "enableUpdateCatalog" -> true))
 val datasink1 = glueContext.getCatalogSink(
 database = "tempdb",
 tableName = "fromoptionsoutput",
 redshiftTmpDir = "",
 transformationContext = "datasink1",
 additionalOptions = options_datasink1).writeDynamicFrame(datasource0)
 }
}, JsonOptions("""{"windowSize" : "100 seconds",
 "checkpointLocation" : "s3://from-options-testing-eu-central-1/
fromOptionsOutput/checkpoint/"})"""))

```

## def getCatalogSink

```

def getCatalogSink(database : String,
 tableName : String,
 redshiftTmpDir : String = "",
 transformationContext : String = ""
 additionalOptions: JsonOptions = JsonOptions.empty,
 catalogId: String = null
) : DataSink

```

Data Catalog で定義されているテーブル内の指定した場所に行う [DataSink](#) を作成します。

- database – Data Catalog 内のデータベース名。
- tableName – Data Catalog 内のテーブル名。
- redshiftTmpDir – 特定のデータシンクで使用する一時的なステージングディレクトリ。デフォルトでは空に設定されます。
- transformationContext – ジョブのブックマークで使用されるシンクに関連付けられた変換コンテキスト。デフォルトでは空に設定されます。
- additionalOptions – AWS Glue で使用する追加のオプション。
- catalogId – 現在アクセスされている Data Catalog のカタログ ID (アカウント ID)。null の場合は、呼び出し元のアカウント ID のデフォルトが使用されます。

戻り値は DataSink。

## デフ getCatalogSource

```
def getCatalogSource(database : String,
 tableName : String,
 redshiftTmpDir : String = "",
 transformationContext : String = ""
 pushDownPredicate : String = " "
 additionalOptions: JsonOptions = JsonOptions.empty,
 catalogId: String = null
) : DataSource
```

Data Catalog のテーブル定義からデータを読み取る [DataSource 特性](#) を作成します。

- database – Data Catalog 内のデータベース名。
- tableName – Data Catalog 内のテーブル名。
- redshiftTmpDir – 特定のデータシンクで使用する一時的なステージングディレクトリ。デフォルトでは空に設定されます。
- transformationContext – ジョブのブックマークで使用されるシンクに関連付けられた変換コンテキスト。デフォルトでは空に設定されます。
- pushDownPredicate – データセットのすべてのファイルをリストアップして読み取る必要がないフィルタパーティション。詳しくは、「[プッシュダウン述語を使用した事前フィルタ処理](#)」を参照してください。
- additionalOptions - オプションの名前と値のペアのコレクション。[AWS Glue for Spark での ETL の接続タイプとオプション](#) でリストされている使用可能なオプション (endpointUrl、streamName、bootstrap.servers、security.protocol、topicName、class および delimiter を除く)。別のオプションとして、catalogPartitionPredicate もサポートされています。

catalogPartitionPredicate – カタログ式を渡して、インデックス列に基づいたフィルタリングができます。これにより、フィルタリングをサーバー側で処理できます。詳細については、「[AWS Glue パーティションインデックス](#)」を参照してください。push\_down\_predicate と catalogPartitionPredicate では、異なる構文が使用されることに注意してください。前者では Spark SQL の標準構文を使用し、後者では JSQL パーサーを使用します。

- catalogId – 現在アクセスされている Data Catalog のカタログ ID (アカウント ID)。null の場合は、呼び出し元のアカウント ID のデフォルトが使用されます。

戻り値は DataSource。

## ストリーミングソースの例

```
val data_frame_datasource0 = glueContext.getCatalogSource(
 database = "tempdb",
 tableName = "test-stream-input",
 redshiftTmpDir = "",
 transformationContext = "datasource0",
 additionalOptions = JsonOptions("""{
 "startingPosition": "TRIM_HORIZON", "inferSchema": "false"}""")
).getDataFrame()
```

## def getJDBCSink

```
def getJDBCSink(catalogConnection : String,
 options : JsonOptions,
 redshiftTmpDir : String = "",
 transformationContext : String = "",
 catalogId: String = null
) : DataSink
```

Data Catalog 内の Connection オブジェクトに指定されている、JDBC データベースに対し書き込みを行う [DataSink](#) を作成します。Connection オブジェクトには、URL、ユーザー名、パスワード、VPC、サブネット、セキュリティグループなど、JDBC シンクに接続するための情報がありません。

- catalogConnection – 書き込み先の JDBC URL を含む Data Catalog の接続名。
- options – JDBC データストアへの書き込みに必要な追加情報を提供する、名前/値ペアを示す JSON 形式の文字列。これには、以下が含まれます。
  - dbtable (必須) - JDBC テーブルの名前。データベース内でスキーマをサポートする JDBC データストアの場合、schema.table-name を指定します。スキーマを指定しない場合、デフォルトの「パブリック」スキーマが使用されます。次の例では、test\_db データベースの test という名前のスキーマおよび test\_table という名前のテーブルを指すオプションパラメータを示します。

```
options = JsonOptions("""{"dbtable": "test.test_table", "database": "test_db"}""")
```

- database (必須) - JDBC データベースの名前。
- その他のオプションはすべて SparkSQL JDBC ライターに直接渡されます。詳細については、「[Redshift data source for Spark](#)」を参照してください。

- `redshiftTmpDir` – 特定のデータシンクで使用する一時的なステージングディレクトリ。デフォルトでは空に設定されます。
- `transformationContext` – ジョブのブックマークで使用するシンクに関連付けられた変換コンテキスト。デフォルトでは空に設定されます。
- `catalogId` – 現在アクセスされている Data Catalog のカタログ ID (アカウント ID)。null の場合は、呼び出し元のアカウント ID のデフォルトが使用されます。

コード例:

```
getJDBCSink(catalogConnection = "my-connection-name", options =
 JsonOptions("""{"dbtable": "my-jdbc-table", "database": "my-jdbc-db"}"""),
 redshiftTmpDir = "", transformationContext = "datasink4")
```

戻り値は `DataSink`。

`def getSink`

```
def getSink(connectionType : String,
 connectionOptions : JsonOptions,
 transformationContext : String = ""
) : DataSink
```

Amazon Simple Storage Service (Amazon S3)、JDBC、AWS Glue データカタログ、Apache Kafka や Amazon Kinesis データストリームなどのデステイネーションにデータを書き込むデータを作成します。 [DataSink](#)

- `connectionType` - 接続のタイプ。 [the section called “接続パラメータ”](#) を参照してください。
- `connectionOptions` – データシンクとの接続を確立するための追加情報を提供する JSON 形式の名前と値のペアの文字列。 [the section called “接続パラメータ”](#) を参照してください。
- `transformationContext` – ジョブのブックマークで使用するシンクに関連付けられた変換コンテキスト。デフォルトでは空に設定されます。

戻り値は `DataSink`。

`getSinkWithdef` フォーマット

```
def getSinkWithFormat(connectionType : String,
```

```
options : JsonOptions,
transformationContext : String = "",
format : String = null,
formatOptions : JsonOptions = JsonOptions.empty
) : DataSink
```

Amazon S3、JDBC、Data Catalog、Apache Kafka、Amazon Kinesis Data Streams などの書き込み先にデータを書き込む [DataSink](#) を作成します。書き込み先に書き込むデータの形式も設定します。

- `connectionType` - 接続のタイプ。 [the section called “接続パラメータ”](#) を参照してください。
- `options` - データシンクとの接続を確立するための追加情報を提供する JSON 形式の名前/値ペアの文字列。 [the section called “接続パラメータ”](#) を参照してください。
- `transformationContext` - ジョブのブックマークで使用されるシンクに関連付けられた変換コンテキスト。デフォルトでは空に設定されます。
- `format` - 書き込み先に書き込むデータの形式。
- `formatOptions` — 書き込み先でデータをフォーマットするための追加オプションを提供する JSON 形式の名前と値のペアの文字列。 [データ形式に関するオプション](#) を参照してください。

戻り値は `DataSink`。

`def getSource`

```
def getSource(connectionType : String,
 connectionOptions : JsonOptions,
 transformationContext : String = ""
 pushDownPredicate
) : DataSource
```

Amazon S3、[DataSource 特性](#) JDBC、AWS Glue データカタログなどのソースからデータを読み取るデータを作成します。Kafka および Kinesis ストリーミングデータソースもサポートしています。

- `connectionType` - データソースの型。 [the section called “接続パラメータ”](#) を参照してください。
- `connectionOptions` - データソースとの接続を確立するための追加情報を提供する JSON 形式の名前と値のペアの文字列。詳細については、「[the section called “接続パラメータ”](#)」を参照してください。

Kinesis ストリーミングソースには、streamARN、startingPosition、inferSchema、および classification の接続オプションが必要です。

Kafka ストリーミングソースに

は、connectionName、topicName、startingOffsets、inferSchema、および classification の接続オプションが必要です。

- transformationContext – ジョブのブックマークで使用されるシンクに関連付けられた変換コンテキスト。デフォルトでは空に設定されます。
- pushDownPredicate – パーティション列に関する述語です。

戻り値は DataSource。

Amazon Kinesis ストリーミングソースの例:

```
val kinesisOptions = jsonOptions()
data_frame_datasource0 = glueContext.getSource("kinesis",
 kinesisOptions).getDataFrame()

private def jsonOptions(): JsonOptions = {
 new JsonOptions(
 s""""{"streamARN": "arn:aws:kinesis:eu-central-1:123456789012:stream/
fromOptionsStream",
 |"startingPosition": "TRIM_HORIZON",
 |"inferSchema": "true",
 |"classification": "json"}"""".stripMargin)
}
```

Kafka ストリーミングソースの例 :

```
val kafkaOptions = jsonOptions()
val data_frame_datasource0 = glueContext.getSource("kafka",
 kafkaOptions).getDataFrame()

private def jsonOptions(): JsonOptions = {
 new JsonOptions(
 s""""{"connectionName": "ConfluentKafka",
 |"topicName": "kafka-auth-topic",
 |"startingOffsets": "earliest",
 |"inferSchema": "true",
 |"classification": "json"}"""".stripMargin)
```

```
}
```

## def フォーマット getSourceWith

```
def getSourceWithFormat(connectionType : String,
 options : JsonOptions,
 transformationContext : String = "",
 format : String = null,
 formatOptions : JsonOptions = JsonOptions.empty
) : DataSource
```

Amazon S3、JDBC、AWS Glue データカタログなどのソースからデータを読み取り、ソースに保存されているデータの形式を設定するを作成します。 [DataSource 特性](#)

- `connectionType` – データソースの型。 [the section called “接続パラメータ”](#) を参照してください。
- `options` – データソースとの接続を確立するための追加情報を提供する、名前と値のペアを示す JSON 形式の文字列。 [the section called “接続パラメータ”](#) を参照してください。
- `transformationContext` – ジョブのブックマークで使用するシンクに関連付けられている変換コンテキスト。デフォルトでは空に設定されます。
- `format` – ソースに保存されるデータの形式。 `connectionType` が "s3" のときは、 `format` を指定することもできます。 "avro"、 "csv"、 "grokLog"、 "ion"、 "json"、 "xml"、 "parquet"、 "orc" のいずれかになります。
- `formatOptions` – ソースでデータを解析するための追加オプションを提供する、名前と値ペアを示す JSON 形式の文字列。 [データ形式に関するオプション](#) を参照してください。

戻り値は DataSource。

## 例

Amazon S3 のカンマ区切り値 (CSV) DynamicFrame ファイルであるデータソースからを作成します。

```
val datasource0 = glueContext.getSourceWithFormat(
 connectionType="s3",
 options =JsonOptions(s""""{"paths": ["s3://csv/nycflights.csv"]}""""),
 transformationContext = "datasource0",
 format = "csv",
 formatOptions=JsonOptions(s""""{"withHeader":"true","separator": ",","}""""))
```

```
).getDynamicFrame()
```

JDBC 接続を使用して PostgreSQL DynamicFrame であるデータソースからを作成します。

```
val datasource0 = glueContext.getSourceWithFormat(
 connectionType="postgresql",
 options =JsonOptions(s""""{
 "url":"jdbc:postgresql://databasePostgres-1.rds.amazonaws.com:5432/testdb",
 "dbtable": "public.company",
 "redshiftTmpDir":"","",
 "user":"username",
 "password":"password123"
 }""""),
 transformationContext = "datasource0").getDynamicFrame()
```

JDBC 接続を使用して MySQL DynamicFrame であるデータソースからを作成します。

```
val datasource0 = glueContext.getSourceWithFormat(
 connectionType="mysql",
 options =JsonOptions(s""""{
 "url":"jdbc:mysql://databaseMysql-1.rds.amazonaws.com:3306/testdb",
 "dbtable": "athenatest_nycflights13_csv",
 "redshiftTmpDir":"","",
 "user":"username",
 "password":"password123"
 }""""),
 transformationContext = "datasource0").getDynamicFrame()
```

def getSparkSession

```
def getSparkSession : SparkSession
```

この GlueContext に関連付けられている SparkSession オブジェクトを取得します。  
SparkSession このオブジェクトを使用して、create from で使用するテーブルと UDF を登録しま  
すDataFrame。DynamicFrames

を返します。 SparkSession

def startTransaction

```
def startTransaction(readOnly: Boolean):String
```

新しいトランザクションの開始。Lake Formation [startTransaction](#) API を内部的に呼び出します。

- `readOnly` — (Boolean) このトランザクションを読み取り専用にするか、または読み取りおよび書き込みを行うかを示します。読み取り専用のトランザクション ID を使用した書き込みは拒否されます。読み取り専用トランザクションはコミットする必要はありません。

トランザクション ID を返します。

`def commitTrac`

```
def commitTransaction(transactionId: String, waitForCommit: Boolean): Boolean
```

指定されたトランザクションをコミットしようとしています。commitTransaction トランザクションのコミットが完了する前に戻る可能性があります。Lake Formation [startTransaction](#) API を内部的に呼び出します。

- `transactionId` — (文字列) コミットするトランザクション。
- `waitForCommit` — (ブール値) commitTransaction がすぐに戻るかどうか指定します。デフォルト値は True です。false の場合、commitTransaction はトランザクションがコミットされるまでポーリングし待機します。最大で 6 回の再試行でエクスポネンシャルバックオフを使用すると、待機時間は 1 分に制限されます。

コミットが完了したかどうかを示すブール値を返します。

`def cancelTransaction`

```
def cancelTransaction(transactionId: String): Unit
```

指定されたトランザクションをキャンセルしようとしています。Lake Formation [CancelTransaction](#) API を内部的に呼び出します。

- `transactionId` — (文字列) キャンセルするトランザクション。

戻り値は、トランザクションが以前にコミットされた場合は `TransactionCommittedException` 例外です。

## def this

```
def this(sc : SparkContext,
 minPartitions : Int,
 targetPartitions : Int)
```

指定した SparkContext、最小パーティション数、ターゲットパーティション数を使用して GlueContext オブジェクトを作成します。

- sc - SparkContext。
- minPartitions - 最小パーティション数。
- targetPartitions - ターゲットパーティション数。

戻り値は GlueContext。

## def this

```
def this(sc : SparkContext)
```

渡された SparkContext で GlueContext オブジェクトを作成します。最小限のパーティション数を 10 に設定し、ターゲットパーティション数を 20 に設定します。

- sc - SparkContext。

戻り値は GlueContext。

## def this

```
def this(sparkContext : JavaSparkContext)
```

渡された JavaSparkContext で GlueContext オブジェクトを作成します。最小限のパーティション数を 10 に設定し、ターゲットパーティション数を 20 に設定します。

- sparkContext - JavaSparkContext。

戻り値は GlueContext。

## MappingSpec

パッケージ: com.amazonaws.services.glue

### MappingSpec ケースクラス

```
case class MappingSpec(sourcePath: SchemaPath,
 sourceType: DataType,
 targetPath: SchemaPath,
 targetType: DataType
) extends Product4[String, String, String, String] {
 override def _1: String = sourcePath.toString
 override def _2: String = ExtendedTypeName.fromDataType(sourceType)
 override def _3: String = targetPath.toString
 override def _4: String = ExtendedTypeName.fromDataType(targetType)
}
```

- sourcePath - ソースフィールドの SchemaPath。
- sourceType - ソースフィールドの DataType。
- targetPath - ターゲットフィールドの SchemaPath。
- targetType - ターゲットフィールドの DataType。

MappingSpec は、ソースパスとソースデータ型から、ターゲットパスとターゲットデータ型へのマッピングを指定します。ソースフレームのソースパスの値は、ターゲットフレームのターゲットパスに表示されます。ソースデータ型はターゲットデータ型にキャストされます。

Product4 からの拡張であるため、applyMapping インターフェイスで Product4 をすべて処理できます。

### MappingSpec オブジェクト

```
object MappingSpec
```

MappingSpec オブジェクトには以下のメンバーがあります。

#### Val orderingByTarget

```
val orderingByTarget: Ordering[MappingSpec]
```

## Def apply

```
def apply(sourcePath : String,
 sourceType : DataType,
 targetPath : String,
 targetType : DataType
) : MappingSpec
```

MappingSpec を作成します。

- sourcePath - ソースパスの文字列表現。
- sourceType - ソース DataType。
- targetPath - ターゲットパスの文字列表現。
- targetType - ターゲット DataType。

戻り値は MappingSpec。

## Def apply

```
def apply(sourcePath : String,
 sourceTypeString : String,
 targetPath : String,
 targetTypeString : String
) : MappingSpec
```

MappingSpec を作成します。

- sourcePath - ソースパスの文字列表現。
- sourceType - ソースデータ型の文字列表現。
- targetPath - ターゲットパスの文字列表現。
- targetType - ターゲットデータ型の文字列表現。

MappingSpec を返します。

## Def apply

```
def apply(product : Product4[String, String, String, String]) : MappingSpec
```

MappingSpec を作成します。

- product - ソースパス、ソースデータ型、ターゲットパス、およびターゲットデータ型の Product4。

戻り値は MappingSpec。

AWS Glue Scala ResolveSpec API

トピック

- [ResolveSpec オブジェクト](#)
- [ResolveSpec ケースクラス](#)

パッケージ: com.amazonaws.services.glue

ResolveSpec オブジェクト

ResolveSpec

```
object ResolveSpec
```

def

```
def apply(path : String,
 action : String
) : ResolveSpec
```

ResolveSpec を作成します。

- path - 解決する必要がある選択肢フィールドの文字列表現。
- action — 解決のためのアクション。アクションは Project、KeepAsStruct、または Cast のいずれかになります。

戻り値は ResolveSpec。

def

```
def apply(product : Product2[String, String]) : ResolveSpec
```

ResolveSpec を作成します。

- product -: ソースパスおよび解決アクションの Product2。

戻り値は ResolveSpec。

ResolveSpec ケースクラス

```
case class ResolveSpec extends Product2[String, String] (
 path : SchemaPath,
 action : String)
```

ResolveSpec を作成します。

- path - 解決する必要がある選択枝フィールドの SchemaPath。
- action — 解決のためのアクション。アクションは Project、KeepAsStruct、または Cast のいずれかになります。

ResolveSpec def メソッド

```
def _1 : String
```

```
def _2 : String
```

AWS Glue Scala ArrayNode API

パッケージ: com.amazonaws.services.glue.types

ArrayNode ケースクラス

ArrayNode

```
case class ArrayNode extends DynamicNode (
 value : ArrayBuffer[DynamicNode])
```

ArrayNode def メソッド

```
def add(node : DynamicNode)
```

```
def clone
```

```
def equals(other : Any)
```

```
def get(index : Int) : Option[DynamicNode]
```

```
def getValue
```

```
def hashCode : Int
```

```
def isEmpty : Boolean
```

```
def nodeType
```

```
def remove(index : Int)
```

```
def this
```

```
def toIterator : Iterator[DynamicNode]
```

```
def toJson : String
```

```
def update(index : Int,
 node : DynamicNode)
```

## AWS Glue Scala BinaryNode API

パッケージ: `com.amazonaws.services.glue.types`

### BinaryNode ケースクラス

#### BinaryNode

```
case class BinaryNode extends ScalarNode(value, TypeCode.BINARY) (
 value : Array[Byte])
```

## BinaryNode val フィールド

- ordering

## BinaryNode def メソッド

```
def clone
```

```
def equals(other : Any)
```

```
def hashCode : Int
```

## AWS Glue Scala BooleanNode API

パッケージ: com.amazonaws.services.glue.types

## BooleanNode ケースクラス

### BooleanNode

```
case class BooleanNode extends ScalarNode(value, TypeCode.BOOLEAN) (
 value : Boolean)
```

## BooleanNode val フィールド

- ordering

## BooleanNode def メソッド

```
def equals(other : Any)
```

## AWS Glue Scala ByteNode API

パッケージ: com.amazonaws.services.glue.types

## ByteNode ケースクラス

### ByteNode

```
case class ByteNode extends ScalarNode(value, TypeCode.BYTE) (
 value : Byte)
```

```
value : Byte)
```

ByteNode val フィールド

- ordering

ByteNode def メソッド

```
def equals(other : Any)
```

AWS Glue Scala DateNode API

パッケージ: com.amazonaws.services.glue.types

DateNode ケースクラス

DateNode

```
case class DateNode extends ScalarNode(value, TypeCode.DATE) (
 value : Date)
```

DateNode val フィールド

- ordering

DateNode def メソッド

```
def equals(other : Any)
```

```
def this(value : Int)
```

AWS Glue Scala DecimalNode API

パッケージ: com.amazonaws.services.glue.types

DecimalNode ケースクラス

DecimalNode

```
case class DecimalNode extends ScalarNode(value, TypeCode.DECIMAL) (
 value : Decimal)
```

```
value : BigDecimal)
```

DecimalNode val フィールド

- ordering

DecimalNode def メソッド

```
def equals(other : Any)
```

```
def this(value : Decimal)
```

AWS Glue Scala DoubleNode API

パッケージ: com.amazonaws.services.glue.types

DoubleNode ケースクラス

DoubleNode

```
case class DoubleNode extends ScalarNode(value, TypeCode.DOUBLE) (
 value : Double)
```

DoubleNode val フィールド

- ordering

DoubleNode def メソッド

```
def equals(other : Any)
```

AWS Glue Scala DynamicNode API

トピック

- [DynamicNode クラス](#)
- [DynamicNode オブジェクト](#)

パッケージ: com.amazonaws.services.glue.types

## DynamicNode クラス

### DynamicNode

```
class DynamicNode extends Serializable with Cloneable
```

### DynamicNode def メソッド

```
def getValue : Any
```

平坦な値を取得し、現在のレコードにバインドします。

```
def nodeType : TypeCode
```

```
def toJson : String
```

### デバッグのメソッド:

```
def toRow(schema : Schema,
 options : Map[String, ResolveOption]
) : Row
```

```
def typeName : String
```

## DynamicNode オブジェクト

### DynamicNode

```
object DynamicNode
```

### DynamicNode def メソッド

```
def quote(field : String,
 useQuotes : Boolean
) : String
```

```
def quote(node : DynamicNode,
 useQuotes : Boolean
```

```
) : String
```

## EvaluateDataQuality クラス

AWS Glue Data Quality は AWS Glue のプレビューリリースに含まれ、変更される可能性があります。

パッケージ: com.amazonaws.services.glue.dq

```
object EvaluateDataQuality
```

### def apply

```
def apply(frame: DynamicFrame,
 ruleset: String,
 publishingOptions: JsonOptions = JsonOptions.empty): DynamicFrame
```

DynamicFrame に対してデータ品質ルールセットを評価し、評価の結果を含む新しい DynamicFrame を返します。AWS Glue Data Quality の詳細については、「[AWS Glue データ品質](#)」を参照してください。

- frame – データ品質を評価する DynamicFrame。
- ruleset – 文字列形式のデータ品質定義言語 (DQDL) ルールセット。DQDL の詳細については、[データ品質定義言語 \(DQDL\) リファレンス](#) のガイドを参照してください。
- publishingOptions – 評価結果とメトリクスを発行する次のオプションを指定するディクショナリ。
  - dataQualityEvaluationContext – AWS Glue が Amazon CloudWatch メトリクスとデータ品質結果を発行する名前空間を指定する文字列。集計されたメトリクスは CloudWatch に表示され、完全な結果は AWS Glue Studio インターフェイスに表示されます。
    - 必須: いいえ
    - デフォルト値: default\_context
  - enableDataQualityCloudWatchMetrics – データ品質評価の結果を CloudWatch に発行するかどうかを指定します。dataQualityEvaluationContext オプションを使用してメトリクスの名前空間を指定します。
    - 必須: いいえ

- デフォルト値: False
- `enableDataQualityResultsPublishing` – データ品質結果を AWS Glue Studio インターフェイスの [Data Quality] (データ品質) タブに表示するかどうかを指定します。
  - 必須: いいえ
  - デフォルト値: True
- `resultsS3Prefix` – AWS Glue がデータ品質評価結果を書き込める Amazon S3 ロケーションを指定します。
  - 必須: いいえ
  - デフォルト値: "" (空の文字列)

## 例

次のコード例は、`SelectFields` 変換を実行する前に `DynamicFrame` のデータ品質を評価する方法を示しています。変換を試行する前に、すべてのデータ品質ルールに合格していることを、このスクリプトで検証します。

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.MappingSpec
import com.amazonaws.services.glue.errors.CallSite
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._
import com.amazonaws.services.glue.dq.EvaluateDataQuality

object GlueApp {
 def main(sysArgs: Array[String]) {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 // @params: [JOB_NAME]
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)

 // Create DynamicFrame with data
 val Legislators_Area = glueContext.getCatalogSource(database="legislators",
 tableName="areas_json", transformationContext="S3bucket_node1").getDynamicFrame()

 // Define data quality ruleset
```

```
val DQ_Ruleset = ""
 Rules = [ColumnExists "id"]
""

// Evaluate data quality
val DQ_Results = EvaluateDataQuality.apply(frame=Legislators_Area,
ruleset=DQ_Ruleset, publishingOptions=JsonOptions("""{"dataQualityEvaluationContext":
"Legislators_Area", "enableDataQualityMetrics": "true",
"enableDataQualityResultsPublishing": "true"}""))
 assert(DQ_Results.filter(_.getField("Outcome").contains("Failed")).count == 0,
"Failing DQ rules for Legislators_Area caused the job to fail.")

// Script generated for node Select Fields
val SelectFields_Results = Legislators_Area.selectFields(paths=Seq("id", "name"),
transformationContext="Legislators_Area")

Job.commit()
}
}
```

## AWS Glue Scala FloatNode API

パッケージ: `com.amazonaws.services.glue.types`

### FloatNode ケースクラス

#### FloatNode

```
case class FloatNode extends ScalarNode(value, TypeCode.FLOAT) (
 value : Float)
```

#### FloatNode val フィールド

- `ordering`

#### FloatNode def メソッド

```
def equals(other : Any)
```

### FillMissingValues クラス

パッケージ: `com.amazonaws.services.glue.ml`

```
object FillMissingValues
```

## def apply

```
def apply(frame: DynamicFrame,
 missingValuesColumn: String,
 outputColumn: String = "",
 transformationContext: String = "",
 callSite: CallSite = CallSite("Not provided", ""),
 stageThreshold: Long = 0,
 totalThreshold: Long = 0): DynamicFrame
```

指定された列の動的フレームの欠落値を埋め、推定値を含む新しい列を持つ新しいフレームを返します。欠落した値がない行の場合、指定した列の値が新しい列に複製されます。

- `frame` — 欠落した値を埋めるための `DynamicFrame`。必須。
- `missingValuesColumn` — 欠落した値 (`null` 値および空の文字列) を含む列。必須。
- `outputColumn` — 値が欠落しているすべての行に推定値が挿入された新しい列の名前。デフォルトは、末尾に `"_filled"` が付加された `missingValuesColumn` の値です。
- `transformationContext` — 状態の情報を識別するために使用される一意の文字列 (オプション)。
- `callSite` — エラーレポートのコンテキスト情報を提供するために使用します (オプション)。
- `stageThreshold` — エラーで終了する前に変換で許容される、エラーの最大発生数 (オプション、デフォルト値は 0)。
- `totalThreshold` — エラーによる終了が処理される前に、全体で許容されるエラーの最大発生数 (オプション、デフォルト値は 0)。

値が欠落した行のための推定値を含む (追加された) 1 つの列、さらに他の行の現在値を含む、新しい動的フレームを返します。

## FindMatches クラス

パッケージ: `com.amazonaws.services.glue.ml`

```
object FindMatches
```

## def apply

```
def apply(frame: DynamicFrame,
 transformId: String,
 transformationContext: String = "",
 callSite: CallSite = CallSite("Not provided", ""),
 stageThreshold: Long = 0,
 totalThreshold: Long = 0,
 enforcedMatches: DynamicFrame = null): DynamicFrame,
computeMatchConfidenceScores: Boolean
```

入力フレームで一致するものを検索し、一致したグループごとに、一意の ID を使用する新しい列を持つ新しいフレームを返します。

- `frame` — 検索の対象となる `DynamicFrame`。必須。
- `transformId` — 入力フレームに適用する `FindMatches` 変換に関連付けられた一意の ID。必須。
- `transformationContext` — この `DynamicFrame` の識別子。 `transformationContext` は、実行全体で保持される、ジョブのブックマーク状態のキーとして使用されます。オプション。
- `callSite` — エラーレポートのコンテキスト情報を提供するために使用します。これらの値は、Python から呼び出すときに、自動的に設定されます。オプション。
- `stageThreshold` — この `DynamicFrame` の計算から例外がスローされるまでに許容される、エラーレコードの最大数 (以前の `DynamicFrame` に含まれているレコードは除きます)。オプション。デフォルト値は 0 です。
- `totalThreshold` — 例外がスローされるまでのエラーレコードの合計最大数。以前のフレームにあるレコードも含まれます。オプション。デフォルト値は 0 です。
- `enforcedMatches` — 強制的な一致のためのフレーム。オプション。デフォルト: `null`。
- `computeMatchConfidenceScores` — 一致するレコードの各グループの信頼スコアをコンピューティングするかどうかを示すブール値。オプション。デフォルトは `False` です。

一致するレコードの各グループに割り当てられた一意の識別子を含む、新しい動的フレームを返します。

### FindIncrementalMatches クラス

パッケージ: `com.amazonaws.services.glue.ml`

```
object FindIncrementalMatches
```

## def apply

```
apply(existingFrame: DynamicFrame,
 incrementalFrame: DynamicFrame,
 transformId: String,
 transformationContext: String = "",
 callSite: CallSite = CallSite("Not provided", ""),
 stageThreshold: Long = 0,
 totalThreshold: Long = 0,
 enforcedMatches: DynamicFrame = null): DynamicFrame,
computeMatchConfidenceScores: Boolean
```

既存ならびに増分フレームにわたって一致を検索し、一致グループごとに一意の ID を使用する列を含む、新しいフレームを返します。

- `existingframe` — グループごとに一致する ID が割り当てられている既存のフレーム。必須。
- `incrementalframe` — 既存のフレームに対する一致を検索するために使用される増分フレームです。必須。
- `transformId` — 入力フレームに適用する `FindIncrementalMatches` 変換に関連付けられた一意の ID。必須。
- `transformationContext` — この `DynamicFrame` の識別子。`transformationContext` は、実行全体で保持される、ジョブのブックマーク状態のキーとして使用されます。オプション。
- `callSite` — エラーレポートのコンテキスト情報を提供するために使用します。これらの値は、Python から呼び出すときに、自動的に設定されます。オプション。
- `stageThreshold` — この `DynamicFrame` の計算から例外がスローされるまでに許容される、エラーレコードの最大数 (以前の `DynamicFrame` に含まれているレコードは除きます)。オプション。デフォルト値は 0 です。
- `totalThreshold` — 例外がスローされるまでのエラーレコードの合計最大数。以前のフレームにあるレコードも含まれます。オプション。デフォルト値は 0 です。
- `enforcedMatches` — 強制的な一致のためのフレーム。オプション。デフォルト: `null`。
- `computeMatchConfidenceScores` — 一致するレコードの各グループの信頼スコアをコンピューティングするかどうかを示すブール値。オプション。デフォルトは `False` です。

一致するレコードの各グループに割り当てられた一意の識別子を含む、新しい動的フレームを返します。

## AWS Glue Scala IntegerNode API

パッケージ: com.amazonaws.services.glue.types

### IntegerNode ケースクラス

#### IntegerNode

```
case class IntegerNode extends ScalarNode(value, TypeCode.INT) (
 value : Int)
```

#### IntegerNode val フィールド

- ordering

#### IntegerNode def メソッド

```
def equals(other : Any)
```

## AWS Glue Scala LongNode API

パッケージ: com.amazonaws.services.glue.types

### LongNode ケースクラス

#### LongNode

```
case class LongNode extends ScalarNode(value, TypeCode.LONG) (
 value : Long)
```

#### LongNode val フィールド

- ordering

#### LongNode def メソッド

```
def equals(other : Any)
```

## AWS Glue Scala MapLikeNode API

パッケージ: com.amazonaws.services.glue.types

## MapLikeNode クラス

### MapLikeNode

```
class MapLikeNode extends DynamicNode (
 value : mutable.Map[String, DynamicNode])
```

### MapLikeNode def メソッド

```
def clear : Unit
```

```
def get(name : String) : Option[DynamicNode]
```

```
def getValue
```

```
def has(name : String) : Boolean
```

```
def isEmpty : Boolean
```

```
def put(name : String,
 node : DynamicNode
) : Option[DynamicNode]
```

```
def remove(name : String) : Option[DynamicNode]
```

```
def toIterator : Iterator[(String, DynamicNode)]
```

```
def toJson : String
```

```
def toJson(useQuotes : Boolean) : String
```

例: この JSON の場合:

```
{"foo": "bar"}
```

useQuotes == true であれば、toJson は {"foo": "bar"} を生成します。useQuotes == false であれば、toJson は {foo: bar} @return を生成します。

## AWS Glue Scala MapNode API

パッケージ: com.amazonaws.services.glue.types

### MapNode ケースクラス

#### MapNode

```
case class MapNode extends MapLikeNode(value) (
 value : mutable.Map[String, DynamicNode])
```

#### MapNode def メソッド

```
def clone
```

```
def equals(other : Any)
```

```
def hashCode : Int
```

```
def nodeType
```

```
def this
```

## AWS Glue Scala NullNode API

### トピック

- [NullNode クラス](#)
- [NullNode ケースオブジェクト](#)

パッケージ: com.amazonaws.services.glue.types

### NullNode クラス

#### NullNode

```
class NullNode
```

## NullNode ケースオブジェクト

### NullNode

```
case object NullNode extends NullNode
```

## AWS Glue Scala ObjectNode API

### トピック

- [ObjectNode オブジェクト](#)
- [ObjectNode ケースクラス](#)

パッケージ: `com.amazonaws.services.glue.types`

## ObjectNode オブジェクト

### ObjectNode

```
object ObjectNode
```

## ObjectNode def メソッド

```
def apply(frameKeys : Set[String],
 v1 : mutable.Map[String, DynamicNode],
 v2 : mutable.Map[String, DynamicNode],
 resolveWith : String
) : ObjectNode
```

## ObjectNode ケースクラス

### ObjectNode

```
case class ObjectNode extends MapLikeNode(value) (
 val value : mutable.Map[String, DynamicNode])
```

## ObjectNode def メソッド

```
def clone
```

```
def equals(other : Any)
```

```
def hashCode : Int
```

```
def nodeType
```

```
def this
```

## AWS Glue Scala ScalarNode API

### トピック

- [ScalarNode クラス](#)
- [ScalarNode オブジェクト](#)

パッケージ: `com.amazonaws.services.glue.types`

### ScalarNode クラス

#### ScalarNode

```
class ScalarNode extends DynamicNode (
 value : Any,
 scalarType : TypeCode)
```

### ScalarNode def メソッド

```
def compare(other : Any,
 operator : String
) : Boolean
```

```
def getValue
```

```
def hashCode : Int
```

```
def nodeType
```

```
def toJson
```

## ScalarNode オブジェクト

### ScalarNode

```
object ScalarNode
```

### ScalarNode def メソッド

```
def apply(v : Any) : DynamicNode
```

```
def compare(tv : Ordered[T],
 other : T,
 operator : String
) : Boolean
```

```
def compareAny(v : Any,
 y : Any,
 o : String)
```

```
def withEscapedSpecialCharacters(jsonToEscape : String) : String
```

## AWS Glue Scala ShortNode API

パッケージ: `com.amazonaws.services.glue.types`

### ShortNode ケースクラス

#### ShortNode

```
case class ShortNode extends ScalarNode(value, TypeCode.SHORT) (
 value : Short)
```

### ShortNode val フィールド

- `ordering`

## ShortNode def メソッド

```
def equals(other : Any)
```

## AWS Glue Scala StringNode API

パッケージ: com.amazonaws.services.glue.types

## StringNode ケースクラス

### StringNode

```
case class StringNode extends ScalarNode(value, TypeCode.STRING) (
 value : String)
```

## StringNode val フィールド

- ordering

## StringNode def メソッド

```
def equals(other : Any)
```

```
def this(value : UTF8String)
```

## AWS Glue Scala TimestampNode API

パッケージ: com.amazonaws.services.glue.types

## TimestampNode ケースクラス

### TimestampNode

```
case class TimestampNode extends ScalarNode(value, TypeCode.TIMESTAMP) (
 value : Timestamp)
```

## TimestampNode val フィールド

- ordering

## TimestampNode def メソッド

```
def equals(other : Any)
```

```
def this(value : Long)
```

## AWS Glue Scala GlueArgParser API

パッケージ: com.amazonaws.services.glue.util

## GlueArgParser オブジェクト

### GlueArgParser

```
object GlueArgParser
```

これは、AWSGlueDataplanePython パッケージ内の `utils.getResolvedOptions` の Python バージョンと厳密に一致しています。

## GlueArgParser def メソッド

```
def getResolvedOptions(args : Array[String],
 options : Array[String]
) : Map[String, String]
```

```
def initParser(userOptionsSet : mutable.Set[String]) : ArgumentParser
```

## Example ジョブに渡された引数を取得する

ジョブ引数を取得するために、`getResolvedOptions` メソッドを使用できます。aws\_region という名前のジョブ引数を取得する次の例を考えてみましょう。

```
val args = GlueArgParser.getResolvedOptions(sysArgs,
 Seq("JOB_NAME", "aws_region").toArray)
Job.init(args("JOB_NAME"), glueContext, args.asJava)
val region = args("aws_region")
println(region)
```

## AWS Glue Scala Job API

パッケージ: com.amazonaws.services.glue.util

### Job オブジェクト

#### ジョブ

```
object Job
```

### Job def メソッド

```
def commit
```

```
def init(jobName : String,
 glueContext : GlueContext,
 args : java.util.Map[String, String] = Map[String, String]()
) : this.type
```

```
def init(jobName : String,
 glueContext : GlueContext,
 endpoint : String,
 args : java.util.Map[String, String]
) : this.type
```

```
def isInitialized
```

```
def reset
```

```
def runId
```

## AWS Glue for Spark ETL スクリプトのプログラミングの機能と最適化

以下のセクションでは、言語を問わず、AWS Glue for Spark ETL (抽出、変換、ロード) プログラミングに適用される一般的な手法と値について説明します。

### トピック

- [AWS Glue for Spark での ETL の接続タイプとオプション](#)
- [AWS Glue for Spark での入出力のデータ形式に関するオプション](#)
- [Spark SQL ジョブ用の AWS Glue Data Catalog のサポート](#)
- [ジョブのブックマークを使用する](#)
- [AWS Glue Studio 外でのセンシティブデータ検出の使用](#)
- [AWS Glue ビジュアルジョブ API](#)

## AWS Glue for Spark での ETL の接続タイプとオプション

AWS Glue for Spark では、PySpark と Scala のさまざまなメソッドと変換で、`connectionType` パラメータを使用しながら接続タイプを指定します。`connectionOptions` または `options` パラメータを使用して接続オプションを指定します。

`connectionType` パラメータには、次の表に示す値を指定できます。各タイプの関連する `connectionOptions` (または `options`) パラメータ値については、以下のセクションで説明します。特に明記されていない限り、パラメータは、接続がソースまたはシンクとして使用されるときに適用されます。

接続オプションの設定と使用方法を示すサンプルコードについては、「各接続タイプのホームページ」を参照してください。

<code>connectionType</code>	接続先
<a href="#">dynamodb</a>	<a href="#">Amazon DynamoDB</a> データベース
<a href="#">kinesis</a>	<a href="#">Amazon Kinesis Data Streams</a>
<a href="#">s3</a>	<a href="#">Simple Storage Service (Amazon S3)</a>
<a href="#">documentdb</a>	<a href="#">Amazon DocumentDB (MongoDB 互換)</a> データベース
<a href="#">opensearch</a>	<a href="#">Amazon OpenSearch Service</a> 。
<a href="#">Redshift</a>	<a href="#">Amazon Redshift</a> データベース
<a href="#">kafka</a>	<a href="#">Kafka</a> または <a href="#">Amazon Managed Streaming for Apache Kafka</a>
<a href="#">azurecosmos</a>	Azure Cosmos for NoSQL。

connectionType	接続先
<a href="#">azuresql</a>	Azure SQL。
<a href="#">bigquery</a>	Google BigQuery。
<a href="#">mongodb</a>	<a href="#">MongoDB</a> データベース (MongoDB Atlas を含む)。
<a href="#">sqlserver</a>	Microsoft SQL Server データベース (「 <a href="#">JDBC 接続</a> 」を参照)
<a href="#">mysql</a>	<a href="#">MySQL</a> データベース (「 <a href="#">JDBC 接続</a> 」を参照)
<a href="#">oracle</a>	<a href="#">Oracle</a> データベース (「 <a href="#">JDBC 接続</a> 」を参照)
<a href="#">postgresql</a>	<a href="#">PostgreSQL</a> データベース (「 <a href="#">JDBC 接続</a> 」を参照)
<a href="#">saphana</a>	SAP HANA。
<a href="#">snowflake</a>	<a href="#">Snowflake</a> データレイク
<a href="#">teradata</a>	Teradata Vantage。
<a href="#">vertica</a>	Vertica。
<a href="#">custom.*</a>	Spark、Athena、または JDBC データストア (「 <a href="#">カスタムと AWS Marketplace での、connectionType の値</a> 」を参照)
<a href="#">marketplace.*</a>	Spark、Athena、または JDBC データストア (「 <a href="#">カスタムと AWS Marketplace での、connectionType の値</a> 」を参照)

## DynamoDB 接続

Spark 用の AWS Glue を使用して AWS Glue の DynamoDB 内のテーブルに対する読み込みと書き込みを行うことができます。AWS Glue ジョブにアタッチされている IAM 権限を使用して DynamoDB に接続します。AWS Glue は、別の AWS アカウントの DynamoDB テーブルに対するデータの書き込みをサポートしています。詳細については、「[the section called “DynamoDB テーブルへのクロスアカウントおよびクロスリージョンでのアクセス”](#)」を参照してください。

AWS Glue DynamoDB ETL コネクタに加えて、DynamoDB エクスポートコネクタを使用して DynamoDB から読み取ることができます。このコネクタは、DynamoDB

ExportTableToPointInTime リクエストを呼び出し、これを指定した Amazon S3 の場所に [DynamoDB JSON](#) 形式で保存します。次に、AWS Glue は、Amazon S3 のエクスポート場所からデータを読み取ることによって DynamicFrame オブジェクトを作成します。

DynamoDB のライターは、AWS Glue バージョン 1.0 以降で利用可能です。AWS Glue DynamoDB のエクスポートコネクタは、AWS Glue バージョン 2.0 以降で利用可能です。

DynamoDB の詳細については、[Amazon DynamoDB](#) のドキュメントを参照してください。

#### Note

DynamoDB ETL リーダーでは、フィルタまたはプッシュダウン述語は使用できません。

## DynamoDB 接続の設定

AWS Glue から DynamoDB に接続するには、AWS Glue ジョブに関連付けられた IAM ロールに DynamoDB と対話する権限を付与します。DynamoDB からの読み取りまたは書き込みに必要な権限の詳細については、「IAM ドキュメント」の [DynamoDB のアクション、リソース、および条件キー](#) を参照してください。

次の状況では、追加の設定が必要になる場合があります。

- DynamoDB エクスポートコネクタを使用するときは、ジョブが DynamoDB テーブルのエクスポートをリクエストできるように IAM を設定する必要があります。さらに、エクスポート用の Amazon S3 バケットを特定し、DynamoDB がそのバケットに書き込むための適切な権限と、AWS Glue ジョブがそこから読み取るための適切な権限を IAM に付与する必要があります。詳細については、「[DynamoDB でテーブルのエクスポートをリクエストする](#)」を参照してください。
- AWS Glue ジョブに特定の Amazon VPC 接続要件がある場合は、NETWORK AWS Glue 接続タイプを使用してネットワークオプションを指定します。DynamoDB へのアクセスは IAM によって承認されるため、AWS Glue DynamoDB 接続タイプは必要ありません。

## DynamoDB からの読み込みと書き込み

次のコード例では、DynamoDB テーブルからの読み取り (ETL コネクタ経由) と、DynamoDB のテーブルへの書き込み方法を示します。ここでは、あるテーブルから読み取りを行い、別のテーブルに対して書き込みを行っています。

## Python

```
import sys
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from awsglue.utils import getResolvedOptions

args = getResolvedOptions(sys.argv, ["JOB_NAME"])
glue_context= GlueContext(SparkContext.getOrCreate())
job = Job(glue_context)
job.init(args["JOB_NAME"], args)

dyf = glue_context.create_dynamic_frame.from_options(
 connection_type="dynamodb",
 connection_options={"dynamodb.input.tableName": test_source,
 "dynamodb.throughput.read.percent": "1.0",
 "dynamodb.splits": "100"
 }
)
print(dyf.getNumPartitions())

glue_context.write_dynamic_frame_from_options(
 frame=dyf,
 connection_type="dynamodb",
 connection_options={"dynamodb.output.tableName": test_sink,
 "dynamodb.throughput.write.percent": "1.0"
 }
)

job.commit()
```

## Scala

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.DynamoDbDataSink
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._
```

```
object GlueApp {

 def main(sysArgs: Array[String]): Unit = {
 val glueContext = new GlueContext(SparkContext.getOrCreate())
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)

 val dynamicFrame = glueContext.getSourceWithFormat(
 connectionType = "dynamodb",
 options = JsonOptions(Map(
 "dynamodb.input.tableName" -> test_source,
 "dynamodb.throughput.read.percent" -> "1.0",
 "dynamodb.splits" -> "100"
))
).getDynamicFrame()

 print(dynamicFrame.getNumPartitions())

 val dynamoDbSink: DynamoDbDataSink = glueContext.getSinkWithFormat(
 connectionType = "dynamodb",
 options = JsonOptions(Map(
 "dynamodb.output.tableName" -> test_sink,
 "dynamodb.throughput.write.percent" -> "1.0"
))
).asInstanceOf[DynamoDbDataSink]

 dynamoDbSink.writeDynamicFrame(dynamicFrame)

 Job.commit()
 }
}
```

## DynamoDB エクスポートコネクタを使用する

DynamoDB テーブルサイズが 80 GB を超える場合、エクスポートコネクタは ETL コネクタよりもパフォーマンスが向上します。さらに、エクスポートリクエストが AWS Glue ジョブで Spark プロセスの外部で実行される場合、[AWS Glue ジョブの自動スケーリング](#)を有効にして、エクスポートリクエスト中の DPU 使用量を節約できます。エクスポートコネクタでは、Spark エグゼキューターの並列処理のためのスプリット数や、DynamoDB スループットの読み取り率を設定する必要がありません。

**Note**

DynamoDBには、ExportTableToPointInTime リクエストを呼び出すための特定の要件があります。詳細については、「[DynamoDB でテーブルのエクスポートをリクエストする](#)」を参照してください。例えば、このコネクタを使用するには、テーブルでポイントインタイムリストア (PITR) を有効にする必要があります。DynamoDB コネクタは、Amazon S3 への DynamoDB エクスポートの AWS KMS 暗号化もサポートします。AWS Glue ジョブの設定でセキュリティ設定を指定することで、DynamoDB エクスポートの AWS KMS 暗号化が有効になります。KMS キーは、Amazon S3 バケットと同じリージョンにある必要があります。

DynamoDB エクスポートと Amazon S3 ストレージのコストに追加料金が適用されることに注意してください。Amazon S3 でエクスポートされたデータは、ジョブ実行の終了後も保持されるため、DynamoDB を追加でエクスポートしなくても再利用できます。このコネクタを使用するための要件として、テーブルに対してポイントインタイムリカバリ (PITR) が有効になっていることが必要です。

DynamoDB ETL コネクタまたはエクスポートコネクタは、DynamoDB ソースに適用されるフィルタまたはプッシュダウン述語をサポートしていません。

次のコード例は、( エクスポートコネクタを経由して ) 読み取り、パーティションの数を出力する方法を示しています。

**Python**

```
import sys
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from awsglue.utils import getResolvedOptions

args = getResolvedOptions(sys.argv, ["JOB_NAME"])
glue_context= GlueContext(SparkContext.getOrCreate())
job = Job(glue_context)
job.init(args["JOB_NAME"], args)

dyf = glue_context.create_dynamic_frame.from_options(
 connection_type="dynamodb",
 connection_options={
 "dynamodb.export": "ddb",
 "dynamodb.tableArn": test_source,
```

```
 "dynamodb.s3.bucket": bucket_name,
 "dynamodb.s3.prefix": bucket_prefix,
 "dynamodb.s3.bucketOwner": account_id_of_bucket,
 }
)
print(dyf.getNumPartitions())

job.commit()
```

## Scala

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.DynamoDbDataSink
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._

object GlueApp {

 def main(sysArgs: Array[String]): Unit = {
 val glueContext = new GlueContext(SparkContext.getOrCreate())
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)

 val dynamicFrame = glueContext.getSourceWithFormat(
 connectionType = "dynamodb",
 options = JsonOptions(Map(
 "dynamodb.export" -> "ddb",
 "dynamodb.tableArn" -> test_source,
 "dynamodb.s3.bucket" -> bucket_name,
 "dynamodb.s3.prefix" -> bucket_prefix,
 "dynamodb.s3.bucketOwner" -> account_id_of_bucket,
))
).getDynamicFrame()

 print(dynamicFrame.getNumPartitions())

 Job.commit()
 }
}
```

```
}
```

これらの例は、( エクスポートコネクタを介して ) 読み取りを行い、dynamodb 分類を持つ AWS Glue データカタログテーブルからパーティションの数を出力する方法を示しています。

## Python

```
import sys
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from awsglue.utils import getResolvedOptions

args = getResolvedOptions(sys.argv, ["JOB_NAME"])
glue_context= GlueContext(SparkContext.getOrCreate())
job = Job(glue_context)
job.init(args["JOB_NAME"], args)

dynamicFrame = glue_context.create_dynamic_frame.from_catalog(
 database=catalog_database,
 table_name=catalog_table_name,
 additional_options={
 "dynamodb.export": "ddb",
 "dynamodb.s3.bucket": s3_bucket,
 "dynamodb.s3.prefix": s3_bucket_prefix
 }
)
print(dynamicFrame.getNumPartitions())

job.commit()
```

## Scala

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.DynamoDbDataSink
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._
```

```
object GlueApp {

 def main(sysArgs: Array[String]): Unit = {
 val glueContext = new GlueContext(SparkContext.getOrCreate())
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)

 val dynamicFrame = glueContext.getCatalogSource(
 database = catalog_database,
 tableName = catalog_table_name,
 additionalOptions = JsonOptions(Map(
 "dynamodb.export" -> "ddb",
 "dynamodb.s3.bucket" -> s3_bucket,
 "dynamodb.s3.prefix" -> s3_bucket_prefix
))
).getDynamicFrame()
 print(dynamicFrame.getNumPartitions())
 }
}
```

## DynamoDB エクスポート JSON の使用を簡素化

AWS Glue DynamoDB エクスポートコネクタを使用した DynamoDB エクスポートでは、特殊なネスト構造の JSON ファイルが生成されます。詳細については、[データオブジェクト](#)を参照してください。AWS Glue は DynamicFrame 変換を提供します。これを使用して、このような構造をダウンストリームアプリケーションで使いやすい形式に変換できます。

変換は 2 つの方法のいずれかで起動できます。。DynamoDB から読み取るメソッドを呼び出すときに、値 "true" を使用して接続オプション "dynamodb.simplifyDDBJson" を設定できます。また、AWS Glue ライブラリで独立して利用可能なメソッドとして変換を呼び出すこともできます。

DynamoDB エクスポートによって生成された次のスキーマを考えてみましょう。

```
root
|-- Item: struct
| |-- parentMap: struct
| | |-- M: struct
| | | |-- childMap: struct
| | | | |-- M: struct
| | | | | |-- appName: struct
| | | | | | |-- S: string
| | | | | | |-- packageName: struct
```

```

| | | | | | |-- S: string
| | | | | | |-- updatedAt: struct
| | | | | | |-- N: string
| |-- strings: struct
| | |-- SS: array
| | | |-- element: string
| |-- numbers: struct
| | |-- NS: array
| | | |-- element: string
| |-- binaries: struct
| | |-- BS: array
| | | |-- element: string
| |-- isDDBJson: struct
| | |-- BOOL: boolean
| |-- nullValue: struct
| | |-- NULL: boolean

```

この `simplifyDDBJson` 変換により、以下のように簡略化されます。

```

root
|-- parentMap: struct
| |-- childMap: struct
| | |-- appName: string
| | |-- packageName: string
| | |-- updatedAt: string
|-- strings: array
| |-- element: string
|-- numbers: array
| |-- element: string
|-- binaries: array
| |-- element: string
|-- isDDBJson: boolean
|-- nullValue: null

```

#### Note

`simplifyDDBJson` は、AWS Glue バージョン 3.0 以降で有効です。 `unnestDDBJson` トランスフォームを使用して DynamoDB のエクスポート JSON を簡略化することもできます。ユーザーには `unnestDDBJson` から `simplifyDDBJson` への移行をお勧めします。

## DynamoDB オペレーションにおける並列処理の設定

パフォーマンスを向上させるために、DynamoDB コネクタで使用できる特定のパラメーターを調整できます。並列処理パラメーターを調整する際の目標は、プロビジョニングされた AWS Glue ワーカーを最大限に活用することです。そして、さらにパフォーマンスを向上させる必要がある場合は、DPU 数を増加してジョブをスケールアウトすることをお勧めします。

ETL コネクタを使用する場合、`dynamodb.splits` パラメーターを使用して DynamoDB 読み取りオペレーションの並列処理を変更できます。 エクスポートコネクタで読み込みをする場合は、Spark エグゼキュータの並列処理の分割数を構成する必要がありません。DynamoDB 書き込みオペレーションの並列処理は、`dynamodb.output.numParallelTasks` を使用して変更できます。

### DynamoDB ETL コネクタによる読み取り

`dynamodb.splits` ジョブ構成に設定されている最大ワーカー数と以下 `numSlots` の計算に基づいて計算することをお勧めします。自動スケーリングの場合、実際に利用可能なワーカー数は、その上限に基づいて変更される可能性があります。最大ワーカー数の設定についての詳細は、「[the section called “Spark ジョブプロパティの設定”](#)」の「[ワーカー数] (NumberOfWorkers)」を参照してください。

- `numExecutors = NumberOfWorkers - 1`

1 つのエグゼキューターが Spark ドライバー専用になっている場合、他のエグゼキューターはデータの処理に使用されます。

- `numSlotsPerExecutor =`

AWS Glue 3.0 and later versions

- WorkerType が G.1X の場合は 4
- WorkerType が G.2X の場合は 8
- WorkerType が G.4X の場合は 16
- WorkerType が G.8X の場合は 32

AWS Glue 2.0 and legacy versions

- WorkerType が G.1X の場合は 8
- WorkerType が G.2X の場合は 16

- `numSlots = numSlotsPerExecutor * numExecutors`

`dynamodb.splits` を使用可能なスロット数、`numSlots` に設定することをおすすめします。

## DynamoDB への書き込み

この `dynamodb.output.numParallelTasks` パラメータは、以下の計算を使用して Spark タスクごとの WCU を決定するために使用されます。

$$\text{permittedWcuPerTask} = ( \text{TableWCU} * \text{dynamodb.throughput.write.percent} ) / \text{dynamodb.output.numParallelTasks}$$

DynamoDB ライターは、設定が DynamoDB に書き込まれる Spark タスクの数を正確に表している場合に最適に機能します。場合によっては、書き込みパフォーマンスを向上させるために、デフォルトの計算をオーバーライドする必要がある場合があります。このパラメータを指定しない場合、Spark タスクごとに許可される WCU は、次の式により自動的に計算されます。

- `numPartitions = dynamicframe.getNumPartitions()`
- `numSlots` (このセクションすでに定義したとおり)
- `numParallelTasks = min(numPartitions, numSlots)`
- 例 1。DPU=10、WorkerType=Standard。入力 DynamicFrame には、100 個の RDD パーティションがあります。
  - `numPartitions = 100`
  - `numExecutors = (10 - 1) * 2 - 1 = 17`
  - `numSlots = 4 * 17 = 68`
  - `numParallelTasks = min(100, 68) = 68`
- 例 2。DPU=10、WorkerType=Standard。入力 DynamicFrame には、20 個の RDD パーティションがあります。
  - `numPartitions = 20`
  - `numExecutors = (10 - 1) * 2 - 1 = 17`
  - `numSlots = 4 * 17 = 68`
  - `numParallelTasks = min(20, 68) = 20`

### Note

従来の AWS Glue バージョンのジョブと Standard Worker を使用するジョブでは、スロット数を計算するために違った方法が必要となります。これらのジョブのパフォーマンスを調整

する必要がある場合は、サポートされている AWS Glue バージョンに移行することをおすすめします。

## DynamoDB 接続オプションのリファレンス

Amazon DynamoDB への接続を指定します。

接続オプションは、ソース接続とシンク接続とで異なります。

"connectionType": ソースとして ETL コネクタを使用する "dynamodb"

AWS Glue DynamoDB ETL コネクタを使用している場合は、"connectionType": "dynamodb" をソースとして次の接続オプションを使用します。

- "dynamodb.input.tableName": (必須) 読み取り元の DynamoDB テーブル。
- "dynamodb.throughput.read.percent": (オプション) 使用する読み込みキャパシティーユニット (RCU) の割合。デフォルトでは、"0.5" に設定されます。許容値は "0.1" から "1.5" (これらの値を含む) です。
  - 0.5 ではデフォルトの読み込み速度を表し、AWS Glue はテーブルの読み込み容量の半分を消費しようとすることを意味します。上記の値を 0.5 より上に設定すると、AWS Glue は読み取りのリクエストレートを増加させ、0.5 より低くした場合はそのリクエストレートを減少させます。(実際の読み込みレートは、DynamoDB テーブルに、統一ディストリビューションのキーがあるかどうかなどの要因によって変わります。)
  - DynamoDB テーブルがオンデマンドモードの場合、AWS Glue はテーブルの読み取り容量を 40000 として処理します。大きなテーブルをエクスポートする場合は、DynamoDB テーブルをオンデマンドモードに切り替えることをお勧めします。
- "dynamodb.splits": (オプション) 読み取り中にこの DynamoDB テーブルを分割するパーティションの数を定義します。デフォルトでは、"1" に設定されます。許容値は "1" から "1,000,000" (これらの値を含む) です。

1 は並列処理がないことを表します。より良いパフォーマンスを得るためには、より大きな値を (以下の式を使用して) 指定することを強くお勧めします。値を適切に設定する方法の詳細は、[「the section called “DynamoDB 並列処理”」](#)を参照してください。

- "dynamodb.sts.roleArn": (オプション) クロスアカウントアクセスのために引き受ける IAM ロールの ARN。このパラメータは、AWS Glue 1.0 以降で使用可能です。

- "dynamodb.sts.roleSessionName": (任意) STS セッション名。デフォルトでは、「glue-dynamodb-read-sts-session」に設定されています。このパラメータは、AWS Glue 1.0 以降で使用可能です。

"connectionType" : AWS Glue DynamoDB エクスポートコネクタをソースとする "dynamodb"

AWS Glue バージョン 2.0 以降でのみ使用可能な AWS Glue DynamoDB エクスポートコネクタを使用する場合は、ソースとして "connectionType": "dynamodb" で次の接続オプションを使用します。

- "dynamodb.export": (必須) 文字列値:
  - ddb に設定すると、AWS Glue DynamoDB エクスポートコネクタが有効になり、AWS Glue ジョブ中に新しい ExportTableToPointInTimeRequest が呼び出されます。dynamodb.s3.bucket と dynamodb.s3.prefix から渡された場所で新しいエクスポートが生成されます。
  - s3 に設定すると、AWS Glue DynamoDB エクスポートコネクタが有効になりますが、新しい DynamoDB エクスポートの作成はスキップされ、代わりに dynamodb.s3.bucket と dynamodb.s3.prefix がそのテーブルの過去のエクスポートの Amazon S3 ロケーションとして使用されます。
- "dynamodb.tableArn": (必須) 読み取り元の DynamoDB テーブル。
- "dynamodb.unnestDDBJson": (オプション) デフォルト: false。有効な値: ブール値 true に設定すると、エクスポートに存在する DynamoDB JSON 構造体のネスト解除の変換が実行されます。"dynamodb.unnestDDBJson" と "dynamodb.simplifyDDBJson" を同時に true に設定するとエラーになります。AWS Glue 3.0 以降のバージョンでは、DynamoDB マップタイプを簡略化するときの動作を改善するために "dynamodb.simplifyDDBJson" を使用することをおすすめします。詳細については、「[the section called “DynamoDB エクスポート JSON の使用を簡素化”](#)」を参照してください。
- "dynamodb.simplifyDDBJson": (オプション) デフォルト: false。有効な値: ブール値 true に設定すると、エクスポートに存在する DynamoDB JSON 構造のスキーマを簡素化するための変換を実行します。これは "dynamodb.unnestDDBJson" オプションと同じ目的ですが、DynamoDB テーブル内の DynamoDB マップタイプやネストされたマップタイプをより適切にサポートします。このオプションは、AWS Glue バージョン 3.0 以降でのみ有効です。"dynamodb.unnestDDBJson" と "dynamodb.simplifyDDBJson" を同時に true に設定するとエラーになります。詳細については、「[the section called “DynamoDB エクスポート JSON の使用を簡素化”](#)」を参照してください。

- "dynamodb.s3.bucket": (オプション) DynamoDB ExportTableToPointInTime プロセスが実行される Simple Storage Service (Amazon S3) バケットの場所を示します。エクスポートファイル形式は DynamoDB JSON です。
- "dynamodb.s3.prefix": (オプション) DynamoDB ExportTableToPointInTime の読み込みが保存される Simple Storage Service (Amazon S3) バケット内の Amazon S3 プレフィックスの場所を示します。dynamodb.s3.prefix と dynamodb.s3.bucket のどちらも指定しなかった場合、これらの値は AWS Glue ジョブの設定で指定された一時ディレクトリの場所がデフォルトとなります。詳細については、「[AWS Glue で使用される特別なパラメータ](#)」を参照してください。
- "dynamodb.s3.bucketOwner": クロスアカウント Simple Storage Service (Amazon S3) アクセスのために必要なバケット所有者を示します。
- "dynamodb.sts.roleArn": (オプション) DynamoDB テーブルのクロスアカウントアクセスおよび/またはクロスリージョンアクセスのために割り当てる IAM ロールの ARN。注: 同じ IAM ロールの ARN を使用して、ExportTableToPointInTime リクエストに指定された Simple Storage Service (Amazon S3) の場所にアクセスします。
- "dynamodb.sts.roleSessionName": (任意) STS セッション名。デフォルトでは、「glue-dynamodb-read-sts-session」に設定されています。
- "dynamodb.exportTime" (オプション) 有効な値: ISO-8601 インスタントを表す文字列。エクスポートが実行されるべき時点。
- "dynamodb.sts.region": (リージョンのエンドポイントを使用してリージョン間通話を行う場合は必須) 読み取りたい DynamoDB テーブルをホストするリージョン。

"connectionType": ETLコネクタをシンクとして使用する "dynamodb"

"connectionType": "dynamodb" をシンクとして、次の接続オプションを使用します。

- "dynamodb.output.tableName": (必須) 書き込み先の DynamoDB テーブル。
- "dynamodb.throughput.write.percent": (オプション) 使用する書き込みキャパシティユニット (WCU) の割合。デフォルトでは、"0.5" に設定されます。許容値は "0.1" から "1.5" (これらの値を含む) です。
- 0.5 ではデフォルトの読み込み速度を表し、AWS Glue はテーブルの書き込み容量の半分を消費しようとすることを意味します。上記の値を 0.5 より上に設定すると、AWS Glue は書き込みリクエストレートを増加させ、0.5 より低くした場合はそのリクエストレートを減少させます。(実際の書き込みレートは、統一されたキーのディストリビューションが、DynamoDB テーブルにあるかどうかなどの要因によって変わります)。

- DynamoDB テーブルがオンデマンドモードの場合、AWS Glue はテーブルの書き込み容量を 40000 として処理します。大きなテーブルをインポートする場合は、DynamoDB テーブルをオンデマンドモードに切り替えることをお勧めします。
- "dynamodb.output.numParallelTasks": (オプション) 同時に DynamoDB に書き込める並列タスクの数を定義します。Spark タスクごとに許容される WCU を計算するために使用されます。ほとんどの場合、AWS Glue はこの値の妥当なデフォルト値を計算します。 詳細については、「[the section called “DynamoDB 並列処理”](#)」を参照してください。
- "dynamodb.output.retry": (オプション) DynamoDB から ProvisionedThroughputExceededException が送られている場合の、再試行の実行回数を定義します。デフォルトでは、"10" に設定されています。
- "dynamodb.sts.roleArn": (オプション) クロスアカウントアクセスのために引き受ける IAM ロールの ARN。
- "dynamodb.sts.roleSessionName": (任意) STS セッション名。デフォルトでは、「glue-dynamodb-write-sts-session」に設定されています。

## DynamoDB テーブルへのクロスアカウントおよびクロスリージョンでのアクセス

AWS Glue ETL ジョブでは、DynamoDB テーブルに対する、クロスリージョンならびにクロスアカウントでのアクセスを共にサポートしています。AWS Glue ETL ジョブは、別の AWS アカウントの DynamoDB テーブルからのデータ読み取り、および別の AWS アカウントの DynamoDB テーブルへのデータ書き込みをサポートしています。加えて AWS Glue では、別のリージョンの DynamoDB テーブルからの読み取りと、別のリージョンの DynamoDB テーブルへの書き込みもサポートしています。このセクションでは、このアクセスを設定する手順と、スクリプトの例を示します。

このセクションに示す手順では、IAM のチュートリアルに則り、IAM ロールを作成しそのロールにアクセス許可を付与しています。このチュートリアルでは、ロールの引き受けについても解説されていますが、今回はその代わりに、AWS Glue でのロールの引き受けにはジョブ・スクリプトを使用します。また、このチュートリアルでは、クロスアカウントに関連した一般的な作業についても説明しています。詳しくは、IAM ユーザーガイドの「[チュートリアル: IAM ロールを使用した AWS アカウント間のアクセス権の委譲](#)」を参照してください。

### ロールを作成する

[チュートリアルステップ 1](#) に従い、アカウント A に IAM ロールを作成します。ロールのアクセス許可を定義する際に、AmazonDynamoDBReadOnlyAccess または AmazonDynamoDBFullAccess など既存のポリシーをアタッチして、ロールに対し DynamoDB の読み取り/書き込みを許

可します。次の例に、アクセス許可ポリシー AmazonDynamoDBFullAccess を使用する、DynamoDBCrossAccessRole という名前のロールを作成する方法を示します。

### ロールにアクセス許可を付与する

IAM ユーザーガイドにある[チュートリアル](#)のステップ 2 に従い、アカウント B が新しく作成されたロールに切り替えることを許可します。次の例では、ここで示されるステートメントを使用して新しいポリシーを作成しています。

```
{
 "Version": "2012-10-17",
 "Statement": {
 "Effect": "Allow",
 "Action": "sts:AssumeRole",
 "Resource": "<DynamoDBCrossAccessRole's ARN>"
 }
}
```

次に、DynamoDB へのアクセスに使用するグループ/ロール/ユーザーに、このポリシーをアタッチします。

### AWS Glue Job スクリプトの中でロールを引き受ける

この段階で、アカウント B にログインして、AWS Glue ジョブを作成できるようになっています。ジョブを作成するには、[での Spark ジョブのジョブプロパティの設定 AWS Glue](#) にある手順を参照してください。

ジョブスクリプト内で DynamoDBCrossAccessRole ロールを引き受けるには、`dynamodb.sts.roleArn` パラメーターを使用する必要があります。このロールを引き受けることで、アカウント B の DynamoDB にアクセスするために必要な、一時的な認証情報を取得できます。このためのサンプルスクリプトを示します。

リージョン間でのクロスアカウント読み取りの場合 (ETL コネクタ):

```
import sys
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from awsglue.utils import getResolvedOptions

args = getResolvedOptions(sys.argv, ["JOB_NAME"])
```

```
glue_context= GlueContext(SparkContext.getOrCreate())
job = Job(glue_context)
job.init(args["JOB_NAME"], args)

dyf = glue_context.create_dynamic_frame_from_options(
 connection_type="dynamodb",
 connection_options={
 "dynamodb.region": "us-east-1",
 "dynamodb.input.tableName": "test_source",
 "dynamodb.sts.roleArn": "<DynamoDBCrossAccessRole's ARN>"
 }
)
dyf.show()
job.commit()
```

リージョン間でのクロスアカウント読み取りの場合 (ETL コネクタ):

```
import sys
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
from awsglue.utils import getResolvedOptions

args = getResolvedOptions(sys.argv, ["JOB_NAME"])
glue_context= GlueContext(SparkContext.getOrCreate())
job = Job(glue_context)
job.init(args["JOB_NAME"], args)

dyf = glue_context.create_dynamic_frame_from_options(
 connection_type="dynamodb",
 connection_options={
 "dynamodb.export": "ddb",
 "dynamodb.tableArn": "<test_source ARN>",
 "dynamodb.sts.roleArn": "<DynamoDBCrossAccessRole's ARN>"
 }
)
dyf.show()
job.commit()
```

リージョンにまたがるクロスアカウント読み取りおよび書き込みの場合:

```
import sys
from pyspark.context import SparkContext
```

```
from awsglue.context import GlueContext
from awsglue.job import Job
from awsglue.utils import getResolvedOptions

args = getResolvedOptions(sys.argv, ["JOB_NAME"])
glue_context= GlueContext(SparkContext.getOrCreate())
job = Job(glue_context)
job.init(args["JOB_NAME"], args)

dyf = glue_context.create_dynamic_frame_from_options(
 connection_type="dynamodb",
 connection_options={
 "dynamodb.region": "us-east-1",
 "dynamodb.input.tableName": "test_source"
 }
)
dyf.show()

glue_context.write_dynamic_frame_from_options(
 frame=dyf,
 connection_type="dynamodb",
 connection_options={
 "dynamodb.region": "us-west-2",
 "dynamodb.output.tableName": "test_sink",
 "dynamodb.sts.roleArn": "<DynamoDBCrossAccessRole's ARN>"
 }
)

job.commit()
```

## Kinesis 接続

Amazon Kinesis Data Streams から読み取る、または Amazon Kinesis Data Streams に書き込むには、データカタログテーブルに格納されている情報を使用するか、データストリームに直接アクセスするための情報を指定します。Kinesis から Spark に情報を読み取り DataFrame、それを Glue に変換することができます。AWS DynamicFrameKinesis には JSON DynamicFrames 形式で書き込むことができます。データストリームに直接アクセスする場合は、これらのオプションを使用して、データストリームへのアクセス方法に関する情報を提供します。

`getCatalogSource` または `create_data_frame_from_catalog` を使用して Kinesis ストリーミングソースからレコードを消費する場合、ジョブはデータカタログデータベースとテーブル名の情報を持っており、それを使用して Kinesis ストリー

ミングソースから読み込むためのいくつかの基本パラメータを取得することができます。getSource、getSourceWithFormat、createDataFrameFromOptions、または create\_data\_frame\_from\_options を使用している場合、ここで説明する接続オプションを使用して、これらの基本パラメータを指定する必要があります。

Kinesis の接続オプションは、GlueContext クラス内の指定されたメソッドに対して以下の引数で指定することができます。

- Scala
  - connectionOptions: getSource、createDataFrameFromOptions、getSink で使用
  - additionalOptions: getCatalogSource、getCatalogSink で使用
  - options: getSourceWithFormat、getSinkWithFormat で使用
- Python
  - connection\_options:  
create\_data\_frame\_from\_options、write\_dynamic\_frame\_from\_options で使用
  - additional\_options:  
create\_data\_frame\_from\_catalog、write\_dynamic\_frame\_from\_catalog で使用
  - options: getSource、getSink で使用

ストリーミング ETL ジョブに関する注意事項と制限事項については、「[the section called “ストリーミング ETL に関する注意と制限”](#)」を参照してください。

## Kinesis の設定

AWS Glue Spark ジョブで Kinesis データストリームに接続するには、いくつかの前提条件が必要です。

- 読み取りの場合、AWS Glue ジョブには Kinesis データストリームへの読み取りアクセスレベル IAM 権限が必要です。
- 書き込みを行う場合、AWS Glue ジョブには Kinesis データストリームへの書き込みアクセスレベルの IAM 権限が必要です。

場合によっては、追加の前提条件を設定する必要があります。

- AWS Glue ジョブが (通常は他のデータセットに接続するための) 追加のネットワーク接続で設定されていて、その接続のいずれかが Amazon VPC ネットワークオプションを提供している場合、

ジョブは Amazon VPC 経由で通信するように指示されます。この場合、Amazon VPC を介して通信するように Kinesis データストリームを設定する必要があります。そのためには、Amazon VPC と Kinesis データストリームの間インターフェイス VPC エンドポイントを作成します。詳細については、「[インターフェイス VPC エンドポイントと Amazon Kinesis Data Streams の使用](#)」を参照してください。

- 別のアカウントで Amazon Kinesis Data Streams を指定する場合は、クロスアカウントアクセスを許可するようにロールとポリシーを設定する必要があります。詳細については、「[Example: Read From a Kinesis Stream in a Different Account](#)」を参照してください。

ストリーミング ETL ジョブの前提条件の詳細については、「[the section called “ストリーミング ETL ジョブ”](#)」を参照してください。

例: Kinesis ストリームからの読み込み

例: Kinesis ストリームからの読み込み

[the section called “forEachBatch”](#) と組み合わせて使用します。

Amazon Kinesis ストリーミングソースの例:

```
kinesis_options =
 { "streamARN": "arn:aws:kinesis:us-east-2:777788889999:stream/fromOptionsStream",
 "startingPosition": "TRIM_HORIZON",
 "inferSchema": "true",
 "classification": "json"
 }
data_frame_datasource0 =
 glueContext.create_data_frame.from_options(connection_type="kinesis",
 connection_options=kinesis_options)
```

例: Kinesis ストリームへの書き込み

例: Kinesis ストリームからの読み込み

[the section called “forEachBatch”](#) と組み合わせて使用します。

Amazon Kinesis ストリーミングソースの例:

```
kinesis_options =
 { "streamARN": "arn:aws:kinesis:us-east-2:777788889999:stream/fromOptionsStream",
```

```
 "startingPosition": "TRIM_HORIZON",
 "inferSchema": "true",
 "classification": "json"
 }
data_frame_datasource0 =
 glueContext.create_data_frame.from_options(connection_type="kinesis",
 connection_options=kinesis_options)
```

## Kinesis 接続オプションのリファレンス

Amazon Kinesis Data Streams への接続を指定します。

Kinesis ストリーミングデータソースには、次の接続オプションを使用します。

- "streamARN" (必須) 読み取り/書き込みに使用されます。Kinesis データストリームの ARN。
- "classification" (読み取りに必須) 読み取りで使用。レコード内のデータで使用されるファイル形式。データカタログを通じて提供されていない限り、必須です。
- "streamName" - (オプション) 読み取りで使用。読み取り対象/読み取り元の Kinesis データストリームの名前。endpointUrl で使用。
- "endpointUrl" - (オプション) 読み取りで使用。デフォルト: "https://kinesis.us-east-1.amazonaws.com"。Kinesis AWS ストリームのエンドポイント。特別なリージョンに接続する場合を除き、これを変更する必要はありません。
- "partitionKey" - (オプション) 書き込みに使用。レコードを作成する際に使用される Kinesis パーティションキー。
- "delimiter" (オプション) 読み取りに使用。classification が CSV の場合に使用される値の区切り文字。デフォルトは「,」です。
- "startingPosition": (オプション) 読み込みに使用。Kinesis データストリーム内の、データの読み取り開始位置。指定できる値は "latest"、"trim\_horizon"、"earliest"、または UTC 形式のタイムスタンプ文字列であり、この文字列のパターンは yyyy-mm-ddTHH:MM:SSZ です (Z は UTC タイムゾーンのオフセットを +/- で表します。例: 「2023-04-04T08:00:00-04:00」)。デフォルト値は、"latest" です。注: の UTC "startingPosition" 形式のタイムスタンプ文字列は、AWS Glue バージョン 4.0 以降でのみサポートされます。
- "failOnDataLoss": (オプション) アクティブなシャードがないか、有効期限が切れている場合、ジョブは失敗します。デフォルト値は、"false" です。
- "awsSTSRoleARN": (オプション) 読み取り/書き込みに使用。() を使用して引き受けるロールの Amazon リソースネーム AWS Security Token Service (ARN AWS STS)。このロールに

は、Kinesis データストリームのレコードの説明操作または読み取り操作の権限が必要です。このパラメーターは、別のアカウントのデータストリームにアクセスするときに使用する必要があります。"awsSTSSessionName" と組み合わせて使用します。

- "awsSTSSessionName": (オプション) 読み取り/書き込みに使用。AWS STSを使って、ロールを担うセッションの識別子。このパラメーターは、別のアカウントのデータストリームにアクセスするときに使用する必要があります。"awsSTSRoleARN" と組み合わせて使用します。
- "awsSTSEndpoint": (オプション) 引き受けたロールで Kinesis AWS STS に接続するときに使用するエンドポイント。これにより、AWS STS デフォルトのグローバルエンドポイントでは不可能な VPC 内のリージョナルエンドポイントを使用できます。
- "maxFetchTimeInMs": (オプション) 読み込みに使用。ジョブエグゼキューターが Kinesis データストリームから現在のバッチのレコードを読み取るために費やした最大時間は、ミリ秒 (ms) 単位で指定されます。この時間内に複数の GetRecords API コールを行うことができます。デフォルト値は、1000 です。
- "maxFetchRecordsPerShard": (オプション) 読み込みに使用。1 マイクロバッチあたりに Kinesis データストリームでシャードごとにフェッチするレコードの最大数。メモ: ストリーミングジョブが既に Kinesis (同じ get-records 呼び出しで) から余分なレコードを読み取っている場合、クライアントはこの制限を超えることができます。maxFetchRecordsPerShard が厳密である必要がある場合、maxRecordPerRead の倍数にする必要があります。デフォルト値は、100000 です。
- "maxRecordPerRead": (オプション) 読み込みに使用。getRecords オペレーションごとに、Kinesis データストリームからフェッチするレコードの最大数。デフォルト値は、10000 です。
- "addIdleTimeBetweenReads": (オプション) 読み込みに使用。2 つの連続する getRecords オペレーション間の遅延時間を追加します。デフォルト値は、"False" です。このオプションは、Glue バージョン 2.0 以降でのみ設定可能です。
- "idleTimeBetweenReadsInMs": (オプション) 読み込みに使用。2 つの連続する getRecords オペレーション間での、最短の遅延時間 (ミリ秒単位で指定)。デフォルト値は、1000 です。このオプションは、Glue バージョン 2.0 以降でのみ設定可能です。
- "describeShardInterval": (オプション) 読み込みに使用。スクリプトが呼び出す 2 つの ListShards API 間での、再シャーディングを考慮すべき最小時間。詳細については、Amazon Kinesis Data Streams デベロッパーガイドの「[リシャーディングのための戦略](#)」を参照してください。デフォルト値は、1s です。
- "numRetries": (オプション) 読み込みに使用。Kinesis Data Streams API リクエストを再試行する最大の回数。デフォルト値は、3 です。

- "retryIntervalMs": (オプション) 読み込みに使用。Kinesis Data Streams API 呼び出しを再試行するまでのクールオフ期間 (ミリ秒単位で指定)。デフォルト値は、1000です。
- "maxRetryIntervalMs": (オプション) 読み込みに使用。再試行で 2 つの Kinesis Data Streams API を呼び出す間の最大クールオフ期間 (ミリ秒単位で指定)。デフォルト値は、10000です。
- "avoidEmptyBatches": (オプション) 読み込みに使用。バッチ処理を開始する前に、Kinesis データストリームで未読のデータをチェックすることで、空のマイクロバッチジョブを作成しないようにします。デフォルト値は、"False"です。
- "schema": (inferSchema に false を設定した場合は必須) 読み取りに使用。ペイロードの処理に使用するスキーマ。分類が avro である場合、提供されるスキーマには Avro スキーマ形式を使用する必要があります。分類が avro 以外の場合、提供されるスキーマは DDL スキーマ形式である必要があります。

以下に、スキーマの例を示します。

Example in DDL schema format

```
`column1` INT, `column2` STRING , `column3` FLOAT
```

Example in Avro schema format

```
{
 "type": "array",
 "items":
 {
 "type": "record",
 "name": "test",
 "fields":
 [
 {
 "name": "_id",
 "type": "string"
 },
 {
 "name": "index",
 "type":
 [
 "int",
 "string",
 "float"
]
 }
]
 }
}
```

```
]
 }
}
```

- "inferSchema": (オプション) 読み込みに使用。デフォルト値は、「false」です。「true」に設定すると、実行時に、スキーマが foreachbatch 内のペイロードから検出されます。
- "avroSchema": (非推奨) 読み取りに使用。Avro 形式を使用する場合に、Avro データのスキーマを指定するために使用されるパラメータです。このパラメータは非推奨となりました。schema パラメータを使用します。
- "addRecordTimestamp": (オプション) 読み込みに使用。このオプションが「true」に設定されている場合、データ出力には、対応するレコードがストリームによって受信された時刻を表示する「\_\_src\_timestamp」という名前が付けられた追加の列が含まれます。デフォルト値は、「false」です。このオプションは AWS Glue バージョン 4.0 以降でサポートされています。
- "emitConsumerLagMetrics": (オプション) 読み込みに使用。このオプションを 'true' に設定すると、バッチごとに、ストリームが受信した最も古いレコードからそのレコードが到着するまでの間のメトリクスが出力されます。AWS Glue CloudWatchメトリクスの名前は「glue.driver.stream」です。maxConsumerLagInMs」。デフォルト値は、「false」です。このオプションは AWS Glue バージョン 4.0 以降でサポートされています。
- "fanoutConsumerARN": (オプション) 読み込みに使用。streamARN で指定されたストリームの Kinesis ストリームコンシューマの ARN。Kinesis 接続の拡張ファンアウトモードを有効にするために使用されます。拡張ファンアウトが使用された Kinesis ストリームの使用に関する詳細については、「[the section called “Kinesis ストリーミングジョブでの拡張ファンアウトの使用”](#)」を参照してください。
- "recordMaxBufferedTime" - (オプション) 書き込みに使用。デフォルト: 1000 (ミリ秒)。レコードが書き込まれるのを待っている間にバッファリングされる最大時間。
- "aggregationEnabled" - (オプション) 書き込みに使用。デフォルト: true。Kinesis に送信する前にレコードを集約するかどうかを指定します。
- "aggregationMaxSize" - (オプション) 書き込みに使用。デフォルト: 51200 (バイト) レコードがこの制限よりも大きい場合、そのレコードはアグリゲータをバイパスします。注: Kinesis では、レコードサイズに 50 KB の制限が適用されます。これを 50 KB を超えて設定すると、サイズ超過のレコードは Kinesis によって拒否されます。
- "aggregationMaxCount" - (オプション) 書き込みに使用。デフォルト: 4294967295。集計されたレコードにパックされる項目の最大数。

- "producerRateLimit" - (オプション) 書き込みに使用。デフォルト: 150 (%)。1つのプロデューサー (ジョブなど) からの送信されるシャード単位のスループットをバックエンド制限のパーセンテージとして制限できます。
- "collectionMaxCount" - (オプション) 書き込みに使用。デフォルト: 500。1 PutRecords 回のリクエストにまとめるアイテムの最大数。
- "collectionMaxSize" - (オプション) 書き込みに使用。デフォルト: 5242880 (バイト)。1 PutRecords 回のリクエストで送信するデータの最大量。

## Kinesis ストリーミングジョブでの拡張ファンアウトの使用

拡張ファンアウトコンシューマは、通常のコンシューマよりも高い専用スループットで Kinesis ストリームからレコードを受信できます。これは、ジョブなどの Kinesis コンシューマにデータを提供するために使用される転送プロトコルを最適化することによって行われます。Kinesis 拡張ファンアウトの詳細については、[Kinesis のドキュメント](#)を参照してください。

拡張ファンアウトモードでは、maxRecordPerRead および idleTimeBetweenReadsInMs 接続オプションは適用されなくなりました。拡張ファンアウトを使用する場合、これらのパラメータは設定できないためです。リトライの設定オプションは説明どおりに機能します。

以下の手順に従って、ストリーミングジョブの拡張ファンアウトを有効または無効にします。ストリームのデータを消費するジョブごとに、ストリームコンシューマを登録する必要があります。

ジョブの拡張ファンアウト消費を有効にするには:

1. Kinesis API を使用してジョブのストリームコンシューマを登録します。手順に従って、[Kinesis ドキュメント](#)の Kinesis Data Streams API を使用して、拡張ファンアウトでコンシューマを登録します。必要なのは、最初のステップである [RegisterStreamConsumer](#) を呼び出すことです。リクエストは ARN、*consumerARN* を返す必要があります。
2. 接続メソッドの引数で接続オプション fanoutConsumerARN を *consumerARN* に設定します。
3. ジョブを再開してください。

ジョブの拡張ファンアウト消費を無効にするには:

1. メソッド呼び出しから接続オプション fanoutConsumerARN を削除します。
2. ジョブを再開してください。
3. [Kinesis ドキュメント](#)の指示に従ってコンシューマを登録解除します。これらの手順はコンソールに適用されますが、Kinesis API を使用して実行することもできます。Kinesis

API によるストリームコンシューマの登録解除の詳細については、Kinesis ドキュメントの「[DeregisterStreamConsumer](#)」を参照してください。

## Amazon S3 接続

AWS Glue for Spark を使用して Amazon S3 内のファイルの読み込みと書き込みを行うことができます。AWS Glue for Spark は、CSV、Avro、JSON、Orc、Parquet など、Amazon S3 に保存されている多くの一般的なデータ形式をそのままサポートしています。サポートされるデータ形式の詳細については、「[the section called “データ形式に関するオプション”](#)」を参照してください。各データ形式が、異なる AWS Glue の機能のセットをサポートする場合があります。機能のサポートの詳細については、使用しているデータ形式のページを参照してください。さらに、Hudi、Iceberg、Delta Lake のデータレイクフレームワークに保存されているバージョン対応ファイルを読み書きできます。データレイクフレームワークの詳細については、「[the section called “データレイクフレームワーク”](#)」を参照してください。

AWS Glue を使用すると、書き込み中に Amazon S3 オブジェクトをフォルダ構造に分割し、簡単な設定でパーティションごとに取得してパフォーマンスを向上させることができます。また、データを変換する際に小さなファイルをまとめてグループ化するように設定してパフォーマンスを向上させることもできます。Amazon S3 では、bzip2 の読み込み、書き込み、gzip アーカイブを行うことができます。

## トピック

- [S3 接続の設定](#)
- [Amazon S3 接続のオプションのリファレンス](#)
- [廃止されたデータ形式の接続構文](#)
- [Amazon S3 ストレージクラスの除外](#)
- [AWS Glue での ETL 出力のパーティションの管理](#)
- [大きなグループの入力ファイルの読み取り](#)
- [Amazon S3 用の VPC エンドポイント](#)

## S3 接続の設定

AWS Glue with Spark ジョブで Amazon S3 に接続するには、いくつかの前提条件が必要です。

- AWS Glue ジョブには、関連する Amazon S3 バケットに対する IAM 権限が必要です。

場合によっては、追加の前提条件を設定する必要があります。

- クロスアカウントアクセスを設定するときは、Amazon S3 バケットで適切なアクセス制御を行います。
- セキュリティ上の理由から、Amazon S3 リクエストを Amazon VPC 経由でルーティングすることもできます。このアプローチでは、帯域幅や可用性の課題が生じる場合があります。詳細については、「[the section called “Amazon S3 用の VPC エンドポイント”](#)」を参照してください。

## Amazon S3 接続のオプションのリファレンス

Simple Storage Service (Amazon S3) への接続を指定します。

Amazon S3 はテーブルではなくファイルを管理するため、このドキュメントに記載されている接続プロパティを指定することに加えて、ファイルタイプに関する追加の設定を指定する必要があります。この情報はデータ形式オプションを使用して指定します。形式オプションの詳細については、「[the section called “データ形式に関するオプション”](#)」を参照してください。AWS Glue データカタログと統合してこの情報を指定することもできます。

接続オプションとフォーマットオプションの違いの例として、[the section called “create\\_dynamic\\_frame\\_from\\_options”](#) メソッドが `connection_type`、`connection_options`、`format` および `format_options` をどのように使用するかを考えてみます。このセクションでは、`connection_options` に提供されるパラメータについて具体的に説明します。

"`connectionType`": "s3" では、次の接続オプションを使用します。

- "`paths`": (必須) 読み取りのソースとなる Amazon S3 パスのリスト。
- "`exclusions`": (オプション) 除外する Unix スタイルの glob パターンの JSON リストを含む文字列。例えば、"`[\\\"**\\.pdf\\\"]`" はすべての PDF ファイルを除外します。AWS Glue がサポートする glob 構文の詳細については、「[包含パターンと除外パターンを使用する](#)」を参照してください。
- "`compressionType`": または「`compression`」: (オプション) データの圧縮方法を指定します。Simple Storage Service (Amazon S3) ソース用には "`compressionType`" を、Simple Storage Service (Amazon S3) ターゲット用には "`compression`" を使用します。データに標準のファイル拡張子が付いている場合、このオプションは一般的に不要です。指定できる値は "`gzip`" および "`bzip2`" です。特定の形式では、他の圧縮形式がサポートされている場合もあります。機能のサポートの詳細については、データ形式のページを参照してください。

- "groupFiles": (オプション) 入力ファイルが 50,000 個を超える場合、デフォルトでファイルのグループ化が有効化されます。入力ファイルが 50,000 個未満の場合にグループ化を有効化するには、このパラメータに "inPartition" を設定します。入力ファイルが 50,000 個を超える場合に、グループ化を無効にするには、このパラメータを "none" に設定します。
- "groupSize": (オプション) ターゲットグループのサイズ (バイト単位)。デフォルトは、入力データのサイズとクラスターのサイズに基づいて計算されます。入力ファイルが 50,000 個未満の場合、このオプションを有効にするには、"groupFiles" を "inPartition" に設定する必要があります。
- "recurse": (オプション) true に設定した場合は、指定したパスの下にあるすべてのサブディレクトリ内のファイルを再帰的に読み取ります。
- "maxBand": (オプション、詳細設定) このオプションでは、s3 リストの一貫性が認められるまでの期間をミリ秒単位で指定します。Amazon S3 の結果整合性を担保するために、直前の maxBand ミリ秒以内の変更タイムスタンプが付いたファイルが、特に JobBookmarks の使用時に追跡されます。ほとんどのユーザーはこのオプションを設定する必要はありません。デフォルトは 900000 ミリ秒 (15 分) です。
- "maxFilesInBand": (オプション、詳細設定) このオプションは、直前の maxBand 秒間に保存するファイルの最大数を指定します。この数を超えた場合、余分なファイルはスキップされ、次のジョブ実行時にのみ処理されます。ほとんどのユーザーはこのオプションを設定する必要はありません。
- "isFailFast": (オプション) このオプションにより、AWS Glue ETL ジョブがリーダー解析の例外をスローするかどうかを決定します。true に設定すると、Spark タスクが 4 回の再試行のうちにデータを正しく解析できなかった場合、ジョブは速やかに失敗します。
- "catalogPartitionPredicate": (オプション) 読み込みに使用。SQL WHERE 句の内容。非常に多数のパーティションを含むデータカタログテーブルから読み込むときに使用されます。データカタログインデックスから一致するパーティションを取得します。[the section called "create\\_dynamic\\_frame\\_from\\_catalog"](#) メソッド (および他の類似メソッド) のオプション `push_down_predicate` と共に使用します。詳細については、「[the section called "カタログパーティション述語"](#)」を参照してください。
- "partitionKeys": (オプション) 書き込みに使用。列ラベル文字列の配列。AWSGlue は、この設定で指定されたとおりにデータを分割します。詳細については、「[the section called "パーティションの書き込み"](#)」を参照してください。
- "excludeStorageClasses": (オプション) 読み込みに使用。Amazon S3 ストレージクラスを指定する文字列の配列。AWSGlue は、この設定に基づいて Amazon S3 オブジェクトを除外します。詳細については、「[the section called "Amazon S3 ストレージクラスの除外"](#)」を参照してください。

## 廃止されたデータ形式の接続構文

特定のデータ形式には、特定の接続タイプ構文を使用してアクセスできます。この構文は廃止されました。代わりに [the section called “データ形式に関するオプション”](#) で提供されている s3 接続タイプと形式オプションを使用して形式を指定することをお勧めします。

```
"connectionType": "orc"
```

[Apache Hive Optimized Row Columnar \(ORC\)](#) ファイル形式で、Amazon S3 に保存されるファイルへの接続を指定します。

"connectionType": "orc" では、次の接続オプションを使用します。

- paths: (必須) 読み取りのソースとなる Amazon S3 パスのリスト。
- (その他のオプション名/値ペア): 書式設定オプションなどのその他のオプションはすべて SparkSQL DataSource に直接渡されます。

```
"connectionType": "parquet"
```

[Apache Parquet](#) ファイル形式で、Amazon S3 に保存されるファイルへの接続を指定します。

"connectionType": "parquet" では、次の接続オプションを使用します。

- paths: (必須) 読み取りのソースとなる Amazon S3 パスのリスト。
- (その他のオプション名/値ペア): 書式設定オプションなどのその他のオプションはすべて SparkSQL DataSource に直接渡されます。

## Amazon S3 ストレージクラスの除外

Amazon Simple Storage Service (Amazon S3) からファイルまたはパーティションを読み取る AWS Glue ETL ジョブを実行している場合は、一部の Amazon S3 ストレージクラスタイプを除外できます。

Amazon S3 には、次のストレージクラスがあります。

- STANDARD – 頻繁にアクセスされるデータの汎用ストレージ向け。
- INTELLIGENT\_TIERING – アクセスパターンが不明または変化するデータ向け。
- STANDARD\_IA および ONEZONE\_IA – 保持期間が長くアクセス頻度の低いデータ向け。

- GLACIER、DEEP\_ARCHIVE、および REDUCED\_REDUNDANCY – 長期アーカイブおよびデジタル保存向け。

詳細については、Amazon S3 デベロッパーガイドの「[Amazon S3 ストレージクラス](#)」を参照してください。

このセクションの例では、GLACIER および DEEP\_ARCHIVE ストレージクラスを除外する方法を示します。これらのクラスではファイルをリストできますが、復元されない限り、ファイルを読み取ることにはできません (さらに詳細な情報については、Amazon S3 デベロッパーガイドの「[Restoring Archived Objects](#)」を参照してください。)

ストレージクラスの除外を使用することで、これらのストレージクラス層にまたがるパーティションを持つテーブルで AWS Glue ジョブを確実に動作させることができます。除外がない場合、これらの階層からデータを読み取るジョブは、次のエラーで失敗します。AmazonS3Exception: オペレーションはオブジェクトのストレージクラスに対して有効ではありません。

AWS Glue で Amazon S3 ストレージクラスをフィルタリングするには、さまざまな方法があります。

## トピック

- [動的フレームの作成時の Amazon S3 ストレージクラスの除外](#)
- [Data Catalog テーブルでの Amazon S3 ストレージクラスの除外](#)

## 動的フレームの作成時の Amazon S3 ストレージクラスの除外

ダイナミックフレームの作成時に Amazon S3 のストレージクラスを除外するには、`excludeStorageClasses` で `additionalOptions` を使います。AWS Glue では指定されたストレージクラスに対応するファイルを、独自の Amazon S3 Lister の実装を使用して自動的にリストアップし、除外します。

次に、動的フレームの作成時に GLACIER および DEEP\_ARCHIVE ストレージクラスを除外する、Python および Scala の例を示します。

## Python の例

```
glueContext.create_dynamic_frame.from_catalog(
 database = "my_database",
 tableName = "my_table_name",
 redshift_tmp_dir = "",
```

```
transformation_ctx = "my_transformation_context",
additional_options = {
 "excludeStorageClasses" : ["GLACIER", "DEEP_ARCHIVE"]
}
)
```

## Scala の例

```
val* *df = glueContext.getCatalogSource(
 nameSpace, tableName, "", "my_transformation_context",
 additionalOptions = JsonOptions(
 Map("excludeStorageClasses" -> List("GLACIER", "DEEP_ARCHIVE"))
)
).getDynamicFrame()
```

## Data Catalog テーブルでの Amazon S3 ストレージクラスの除外

AWS Glue ETL のジョブによるストレージクラスの除外は、AWS Glue Data Catalog のテーブルパラメータとして指定できます。このパラメータは、AWS Command Line Interface (AWS CLI) を使用するか、API を使用してプログラミングによって CreateTable オペレーションに含むことができます。詳細については、「[テーブル構造](#)」および「[CreateTable](#)」を参照してください。

AWS Glue コンソールで除外されるストレージクラスを指定することもできます。

### Amazon S3 ストレージクラスを除外するには (コンソール)

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. 左側のナビゲーションペインで、[Tables (テーブル)] を選択します。
3. リストでテーブル名を選択し、[Edit table (テーブルの編集)] を選択します。
4. [Table properties (テーブルプロパティ)] で、キーとして **excludeStorageClasses** を、値として **["GLACIER","DEEP\_ARCHIVE"]** を追加します。
5. [Apply] を選択します。

## AWS Glue での ETL 出力のパーティションの管理

パーティション分割は、データセットを整理して効率的にクエリを実行可能にする重要な手法です。1 つまたは複数の列の個別の値に基づいて、データを階層形式のディレクトリ構造に整理します。

たとえば、Amazon Simple Storage Service (Amazon S3) のアプリケーションログを、年、月、日で分類しながら、日付についてパーティション化できます。次に 1 日分のデータに対応するファイルを `s3://my_bucket/logs/year=2018/month=01/day=23/` などのプレフィックス別に配置します。Amazon Athena、Amazon Redshift Spectrum、そして AWS Glue などのシステムでは、基になるデータ全体を Amazon S3 から読み取ることなく、これらのパーティションを使用してパーティション値でデータをフィルタリングできます。

クローラは、ファイルタイプとスキーマを推定するだけでなく、AWS Glue Data Catalog を構成する際に、データセットのパーティション構造も自動的に特定します。これにより出力されるパーティション列に対しては、AWS Glue ETL ジョブやクエリエンジン (Amazon Athena など) からクエリを実行できます。

テーブルのクローラが完了すると、クローラが作成したパーティションを表示できます。AWS Glue コンソールの左のナビゲーションペインで、[Tables] (テーブル) をクリックします。クローラで作成されたテーブルを選択した後、[View Partitions] (パーティション) の表示をクリックします。

Apache Hive 形式のパーティション分割されたパス (`key=val` 形式) の場合、クローラはキー名を使用して自動的に列名を事前設定します。それ以外の場合は、`partition_0`、`partition_1` などのデフォルト名が使用されます。コンソールでデフォルトの名前を変更できます。そのためには、テーブルに移動します。[インデックス] タブにインデックスが存在するかどうかを確認します。ある場合は、削除してから続行する必要があります (後で新しい列名を使用して再作成できます)。次に、[スキーマの編集] を選択し、そこでパーティション列の名前を変更します。

次に、ETL スクリプトでパーティション列をフィルタリングできます。パーティション情報は Data Catalog に格納されるため、パーティション列を `DynamicFrame` に含めるには `from_catalog` API 呼び出しを使用します。例えば、`create_dynamic_frame.from_options` ではなく `create_dynamic_frame.from_catalog` を使用します。

パーティション化は、データスキャンを減らす最適化手法です。この手法が適切であることを特定するプロセスの詳細については、「AWS 規範的ガイダンス」の「Best practices for performance tuning AWS Glue for Apache Spark jobs」ガイドにある「[Reduce the amount of data scan](#)」を参照してください。

### プッシュダウン述語を使用した事前フィルタ処理

多くの場合、プッシュダウン述語を使用してパーティションをフィルタリングできます。データセットのすべてのファイルをリストアップして読み取る必要はありません。データセット全体を読み取って `DynamicFrame` でフィルタリングする代わりに、Data Catalog 内でパーティションのメタデータに直接フィルターを適用できます。次に、実際に必要なものだけをリストアップして `DynamicFrame` 内に読み取ることができます。

たとえば、Python では以下のように記述できます。

```
glue_context.create_dynamic_frame.from_catalog(
 database = "my_S3_data_set",
 table_name = "catalog_data_table",
 push_down_predicate = my_partition_predicate)
```

これによって作成される DynamicFrame では、Data Catalog のパーティションのうち、述語式を満たすものだけがロードされます。ロードするデータのサブセットを絞り込む度合いに応じて、処理時間を大幅に短縮できる場合があります。

述語式として、Spark SQL でサポートされている任意のブール式を使用できます。Spark SQL クエリで WHERE 句に指定できる条件は、すべて正常に動作します。例えば述語式 `pushDownPredicate = "(year=='2017' and month=='04')"` では、Data Catalog 内で `year` が「2017」に等しく、また `month` が「04」に等しいパーティションのみロードされます。詳細については、[Apache Spark SQL のドキュメント](#)を参照してください。特に [Scala SQL 関数リファレンス](#)が参考になります。

カタログのパーティション述語を使用したサーバー側のフィルタリング

`push_down_predicate` オプションは、カタログからすべてのパーティションを一覧表示した後、Amazon S3 にあるそれらのパーティションからファイルをリストする前に適用されます。テーブルに多数のパーティションがある場合、カタログパーティションのリストに、時間的なオーバーヘッドが余計に発生する可能性があります。このオーバーヘッドに対処するには、`catalogPartitionPredicate` オプションを指定して AWS Glue Data Catalog の [パーティションインデックス](#)を使用しながら、サーバー側のパーティションでプルーニングを行います。1つのテーブルに数百万のパーティションがある場合、これにより、パーティションのフィルタリングを大幅に高速化できます。カタログのパーティションインデックスではまだサポートされていない述語構文が、`catalogPartitionPredicate` で必要となる場合には、`additional_options` の中で `push_down_predicate` と `catalogPartitionPredicate` の両方を使用することもできます

Python:

```
dynamic_frame = glueContext.create_dynamic_frame.from_catalog(
 database=dbname,
 table_name=tablename,
 transformation_ctx="datasource0",
 push_down_predicate="day>=10 and customer_id like '10%'",
 additional_options={"catalogPartitionPredicate":"year='2021' and month='06'"}
```

```
)
```

Scala:

```
val dynamicFrame = glueContext.getCatalogSource(
 database = dbname,
 tableName = tablename,
 transformationContext = "datasource0",
 pushDownPredicate="day>=10 and customer_id like '10%'",
 additionalOptions = JsonOptions("""{
 "catalogPartitionPredicate": "year='2021' and month='06'"}""")
).getDynamicFrame()
```

#### Note

`push_down_predicate` と `catalogPartitionPredicate` では、使用される構文が異なります。前者では Spark SQL の標準構文を使用し、後者では JSQL パーサーを使用します。

## パーティションの書き込み

デフォルトでは、`DynamicFrame` は書き込むときにパーティション分割されません。すべての出力ファイルは、指定した出力パスの最上位レベルに書き込まれます。最近まで、`DynamicFrame` をパーティションに書き込む唯一の方法は、書き込む前に Spark SQL `DataFrame` に変換することでした。

ただし、`DynamicFrames` ではキーのシーケンスを使用したネイティブのパーティション分割がサポートされるようになりました。この場合、シンクの作成時に `partitionKeys` オプションを使用します。例えば、次の Python コードではデータセットを Parquet 形式で Amazon S3 のディレクトリに書き込みます。これらのディレクトリでは、型フィールドごとにパーティション分割されています。これらのパーティションに対しては、他のシステム (Amazon Athena など) を使用しての処理が行えます。

```
glue_context.write_dynamic_frame.from_options(
 frame = projectedEvents,
 connection_type = "s3",
 connection_options = {"path": "$outpath", "partitionKeys": ["type"]},
 format = "parquet")
```

## 大きなグループの入カファイルの読み取り

テーブルのプロパティを設定しておくことで、Amazon S3 データストアからファイルが読み取られる際に、AWS Glue の ETL ジョブでそれらのファイルをグループ化させることができます。これらのプロパティを使用すると、各 ETL タスクでは入カファイルのグループを単一のインメモリパーティションに読み取ることができます。これは、Amazon S3 データストアに多数の小さいファイルがある場合に便利です。特定のプロパティを定義する際には、Amazon S3 データパーティション内のファイルをグループ化し、そのグループの読み取り用のサイズを設定することを AWS Glue に指示します。また、`create_dynamic_frame.from_options` メソッドを使用した Amazon S3 データストアからの読み取りのために、これらのオプションを設定することもできます。

テーブルのファイルをグループ化するには、テーブル構造のパラメータフィールドにキーと値のペアを設定します。テーブルのパラメータフィールドに値を設定するには、JSON 表記を使用します。テーブルのプロパティを編集する詳しい方法については、「[テーブルの詳細の表示と編集](#)」を参照してください。

このメソッドを使用すると、Amazon S3 データストアにある Data Catalog でテーブルをグループ化できます。

### groupFiles

Amazon S3 データパーティションにあるファイルのグループ化を有効にするには、`groupFiles` に `inPartition` を設定します。入カファイル数が 50,000 を超える場合には、次の例のように AWS Glue がグループ化を自動的に有効にします。

```
'groupFiles': 'inPartition'
```

### groupSize

`groupSize` をグループのターゲットサイズ (バイト単位) に設定します。`groupSize` プロパティはオプションです。指定しない場合、AWS Glue はクラスター内のすべての CPU コアを使用すると同時に、ETL タスクとメモリ内パーティションの総数を減らすようにサイズを計算します。

たとえば、次の例では、グループサイズを 1 MB に設定します。

```
'groupSize': '1048576'
```

`groupsize` は、計算の結果を使用して設定する必要があることに注意してください。例:  $1024 * 1024 = 1048576$ 。

## recurse

`recurse` に `True` を設定し、`paths` によりパスの配列を指定すると、すべてのサブディレクトリのファイルを再帰的に読み取ることができます。次の例のように、`paths` が Amazon S3 のオブジェクトキーの配列である場合、または入力フォーマットが `parquet/orc` の場合、`recurse` を設定する必要はありません。

```
'recurse':True
```

`create_dynamic_frame.from_options` メソッドを使用して Amazon S3 から直接読み取る場合は、以下の接続オプションを追加します。たとえば、ファイルを 1 MB のグループにグループ化する方法は次のとおりです。

```
df = glueContext.create_dynamic_frame.from_options("s3", {'paths': ["s3://s3path/"],
'recurse':True, 'groupFiles': 'inPartition', 'groupSize': '1048576'}, format="json")
```

### Note

`groupFiles` は、`csv`、`lon`、`grokLog`、`json`、および `xml` のデータ形式から作成された `DynamicFrames` でサポートされています。このオプションは、`avro`、`parquet`、`orc` ではサポートされていません。

## Amazon S3 用の VPC エンドポイント

セキュリティ上の理由から、多くの AWS ユーザーがアプリケーションを Amazon Virtual Private Cloud 環境 (Amazon VPC) 内で実行しています。Amazon VPC を使用すると、Amazon EC2 インスタンスを仮想プライベートクラウドで作成できます。そのため、パブリックインターネットなどの他のネットワークから論理的に分離されます。Amazon VPC を使用すると、IP アドレスの範囲、サブネット、ルーティングテーブル、ネットワークゲートウェイ、セキュリティ設定を適切に管理できます。

**Note**

2013 年 12 月 4 日より後に AWS アカウントを作成した場合は、各 AWS リージョンにデフォルトで VPC が用意されています。追加設定なしにデフォルトの VPC をすぐに使用できます。

デフォルト VPC の詳細については、Amazon VPC ユーザーガイドの「[デフォルト VPC とデフォルトサブネット](#)」を参照してください。

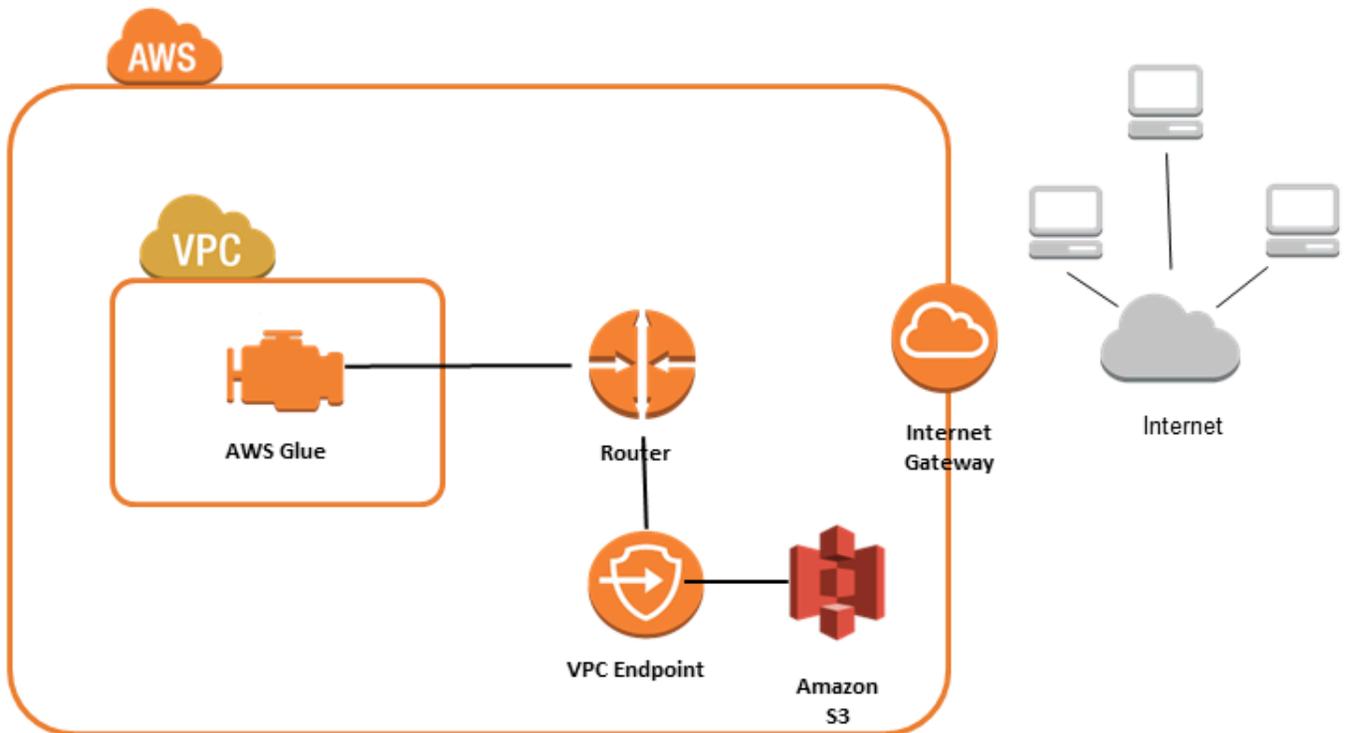
多くのお客様が、パブリックインターネット間のデータ送受信に関して、プライバシーとセキュリティに関する正当な懸念を抱いています。こういったお客様は、この懸念事項を解決するために、バーチャルプライベートネットワーク (VPN) を使用して、すべての Amazon S3 ネットワークトラフィックを、自社内の企業ネットワークのインフラストラクチャ経由でルーティングします。ただし、このアプローチでは、帯域幅や可用性の課題が生じる場合があります。

Amazon S3 用の VPC エンドポイントにより、これらの課題が軽減されます。Amazon S3 用 VPC エンドポイントを使用することで、AWS Glue はプライベート IP アドレスを使用して、パブリックインターネットに公開されることなく Amazon S3 にアクセスできるようになります。AWS Glue はパブリック IP アドレスを必要とせず、VPC のためのインターネットゲートウェイ、NAT デバイス、仮想プライベートゲートウェイは不要です。Amazon S3 へのアクセスを制御するには、エンドポイントのポリシーを使用します。VPC と AWS サービス間のトラフィックは、Amazon ネットワークを離れません。

Amazon S3 用に VPC エンドポイントを作成する際、リージョン内の Amazon S3 エンドポイント (例: s3.us-west-2.amazonaws.com) に対するリクエストはすべて、Amazon ネットワーク内のプライベートの Amazon S3 エンドポイントにルーティングされます。VPC の Amazon EC2 インスタンスで実行されているアプリケーションを変更する必要はありません。エンドポイント名は変わりませんが、Amazon S3 へのルーティングは完全に Amazon ネットワーク内で行われ、パブリックインターネットにアクセスすることはありません。

エンドポイントの詳細については、Amazon VPC ユーザーガイドの「[VPC Endpoints](#)」を参照してください。

AWS Glue が VPC エンドポイントを使用して Amazon S3 にアクセスする様子を、次の図に示します。



Amazon S3 のアクセスをセットアップするには

1. AWS Management Console にサインインして、Amazon VPC コンソール (<https://console.aws.amazon.com/vpc/>) を開きます。
2. 左のナビゲーションペインで [エンドポイント] を選択します。
3. [Create Endpoint] (エンドポイントの作成) をクリックし、ステップに従ってゲートウェイタイプの Amazon S3 VPC エンドポイントを作成します。

Amazon DocumentDB の接続

Spark 用の AWS Glue を使用して Amazon DocumentDB 内のテーブルに対する読み込みと書き込みを行うことができます。AWS Glue 接続を介して AWS Secrets Manager で保存されている認証情報を使用して Amazon DocumentDB に接続できます。

Amazon DocumentDB の詳細については、「[Amazon DocumentDB のドキュメント](#)」を参照してください。

**Note**

Amazon DocumentDB エラスティッククラスターは、現在 AWS Glue コネクタを使用する場合にはサポートされていません。エラスティッククラスターの詳細については、「[Amazon DocumentDB エラスティッククラスターの使用](#)」を参照してください。

## Amazon DocumentDB コレクションへの読み取りと書き込み

**Note**

Amazon DocumentDB に接続する ETL ジョブを作成する際、Connections ジョブのプロパティにおいて、Amazon DocumentDB が実行されている Virtual Private Cloud (VPC) を特定するための、接続オブジェクトを指定する必要があります。接続オブジェクトの場合、接続タイプは JDBC で、JDBC URL は `mongo://<DocumentDB_host>:27017` である必要があります。

**Note**

これらのコードサンプルは AWS Glue 3.0 用に開発されました。AWS Glue 4.0 への移行については、[the section called “MongoDB”](#) を参照してください。uri パラメータが変更されました。

**Note**

Amazon DocumentDB を使用する場合、書かれた文書に `_id` が指定されている場合など、特定の状況では、`retryWrites` は `false` に設定する必要があります。詳細については、Amazon DocumentDB ドキュメントの「[MongoDB との機能の違い](#)」を参照してください。

次の Python スクリプトでは、Amazon DocumentDB に対する読み取りと書き込みのための、接続タイプと接続オプションの使用方法を示しています。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext, SparkConf
from awsglue.context import GlueContext
from awsglue.job import Job
import time

@params: [JOB_NAME]
args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session

job = Job(glueContext)
job.init(args['JOB_NAME'], args)

output_path = "s3://some_bucket/output/" + str(time.time()) + "/"
documentdb_uri = "mongodb://<mongo-instanced-ip-address>:27017"
documentdb_write_uri = "mongodb://<mongo-instanced-ip-address>:27017"

read_docdb_options = {
 "uri": documentdb_uri,
 "database": "test",
 "collection": "coll",
 "username": "username",
 "password": "1234567890",
 "ssl": "true",
 "ssl.domain_match": "false",
 "partitioner": "MongoSamplePartitioner",
 "partitionerOptions.partitionSizeMB": "10",
 "partitionerOptions.partitionKey": "_id"
}

write_documentdb_options = {
 "retryWrites": "false",
 "uri": documentdb_write_uri,
 "database": "test",
 "collection": "coll",
 "username": "username",
 "password": "pwd"
}
```

```
Get DynamicFrame from DocumentDB
dynamic_frame2 =
 glueContext.create_dynamic_frame.from_options(connection_type="documentdb",
 connection_options=read_docdb_options)

Write DynamicFrame to MongoDB and DocumentDB
glueContext.write_dynamic_frame.from_options(dynamic_frame2,
 connection_type="documentdb",

 connection_options=write_documentdb_options)

job.commit()
```

次の Scala スクリプトでは、Amazon DocumentDB に対する読み取りと書き込みのための、接続タイプと接続オプションの使用法を示しています。

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.MappingSpec
import com.amazonaws.services.glue.errors.CallSite
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.DynamicFrame
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._

object GlueApp {
 val DOC_URI: String = "mongodb://<mongo-instanced-ip-address>:27017"
 val DOC_WRITE_URI: String = "mongodb://<mongo-instanced-ip-address>:27017"
 lazy val documentDBJsonOption = jsonOptions(DOC_URI)
 lazy val writeDocumentDBJsonOption = jsonOptions(DOC_WRITE_URI)
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)

 // Get DynamicFrame from DocumentDB
 val resultFrame2: DynamicFrame = glueContext.getSource("documentdb",
 documentDBJsonOption).getDynamicFrame()
```

```
// Write DynamicFrame to DocumentDB
glueContext.getSink("documentdb", writeJsonOption).writeDynamicFrame(resultFrame2)

Job.commit()
}

private def jsonOptions(uri: String): JsonOptions = {
 new JsonOptions(
 s""""{"uri": "${uri}",
 |"database":"test",
 |"collection":"coll",
 |"username": "username",
 |"password": "pwd",
 |"ssl":"true",
 |"ssl.domain_match":"false",
 |"partitioner": "MongoSamplePartitioner",
 |"partitionerOptions.partitionSizeMB": "10",
 |"partitionerOptions.partitionKey": "_id"}"""".stripMargin)
 }
}
```

## Amazon DocumentDB 接続のオプションのリファレンス

Amazon DocumentDB (MongoDB 互換) への接続を指定します。

接続オプションは、ソース接続とシンク接続とで異なります。

"connectionType": "Documentdb" ソースとする

"connectionType": "documentdb" をソースとして、次の接続オプションを使用します。

- "uri": (必須) 読み取り元の Amazon DocumentDB ホスト (mongodb://<host>:<port> 形式)。
- "database": (必須) 読み取り元の Amazon DocumentDB データベース。
- "collection": (必須) 読み取り元の Amazon DocumentDB コレクション。
- "username": (必須) Amazon DocumentDB のユーザー名。
- "password": (必須) Amazon DocumentDB のパスワード。
- "ssl": (SSL を使用する場合は必須) SSL を使用して接続する場合は、このオプションの値を "true" に設定する必要があります。
- "ssl.domain\_match": (SSL を使用する場合は必須) SSL を使用して接続する場合は、このオプションの値を "false" に設定する必要があります。

- "batchSize": (オプション): 内部バッチのカーソル内で使用される、バッチごとに返されるドキュメントの数。
- "partitioner": (オプション) Amazon DocumentDB から入力データを読み取るためのパーティショナーのクラス名。コネクタには、次のパーティショナーがあります。
  - MongoDefaultPartitioner (デフォルト) (AWS Glue 4.0 ではサポートされていません)
  - MongoSamplePartitioner (AWS Glue 4.0 ではサポートされていません)
  - MongoShardedPartitioner
  - MongoSplitVectorPartitioner
  - MongoPaginateByCountPartitioner
  - MongoPaginateBySizePartitioner (AWS Glue 4.0 ではサポートされていません)
- "partitionerOptions" (オプション): 指定されたパーティショナーのオプション。各パーティショナーでは、次のオプションがサポートされています。
  - MongoSamplePartitioner: partitionKey, partitionSizeMB, samplesPerPartition
  - MongoShardedPartitioner: shardkey
  - MongoSplitVectorPartitioner: partitionKey、partitionSizeMB
  - MongoPaginateByCountPartitioner: partitionKey, numberOfPartitions
  - MongoPaginateBySizePartitioner: partitionKey、partitionSizeMB

これらのオプションの詳細については、MongoDB のドキュメントの「[Partitioner Configuration](#)」を参照してください。

"connectionType": "documentdb" as Sink

"connectionType": "documentdb" をシンクとして、次の接続オプションを使用します。

- "uri": (必須) 書き込み先の Amazon DocumentDB ホスト (mongodb://<host>:<port> 形式)。
- "database": (必須) 書き込み先の Amazon DocumentDB データベース。
- "collection": (必須) 書き込み先の Amazon DocumentDB コレクション。
- "username": (必須) Amazon DocumentDB のユーザー名。
- "password": (必須) Amazon DocumentDB のパスワード。
- "extendedBsonTypes": (オプション) true が指定されている場合、Amazon DocumentDB へのデータ書き込み時に拡張 BSON 型を使用できます。デフォルト: true。

- "replaceDocument": (オプション) true の場合、\_id フィールドを含むデータセットを保存するときに、ドキュメント全体を置き換えます。false の場合、データセットのフィールドと一致するドキュメントのフィールドのみが更新されます。デフォルト: true。
- "maxBatchSize": (オプション): データを保存するときの一括オペレーションの最大バッチサイズ。デフォルトは 512 です。
- "retryWrites": (オプション): AWS Glue でネットワークエラーが発生した場合、特定の書き込みオペレーションを 1 回自動的に再試行します。

## OpenSearch Service 接続

AWS Glue for Spark を使用して、AWS Glue 4.0 以降のバージョンの OpenSearch Service のテーブルからの読み取りとテーブルへの書き込みを行うことができます。OpenSearch クエリを使用して、OpenSearch Service から何を読み取るかを定義できます。AWS Glue 接続を通じて AWS Secrets Manager で保存されている HTTP 基本認証情報を使用して OpenSearch Service に接続します。この機能は OpenSearch サービスのサーバーレスとは互換性がありません。

Amazon OpenSearch Service の詳細については、[Amazon OpenSearch Service ドキュメント](#)を参照してください。

## OpenSearch Service 接続の設定

AWS Glue から OpenSearch Service に接続するには、OpenSearch Service 認証情報を作成して AWS Secrets Manager シークレットに保存し、そのシークレットを OpenSearch Service AWS Glue 接続に関連付ける必要があります。

### 前提条件:

- 読み取り元とするドメインエンドポイント *aosEndpoint* とポート *aosPort* を特定するか、または Amazon OpenSearch Service ドキュメントの手順に従ってリソースを作成します。ドメインの作成の詳細については、Amazon OpenSearch Service ドキュメントの「[Amazon OpenSearch Service ドメインの作成と管理](#)」を参照してください。

Amazon OpenSearch Service ドメインエンドポイントのデフォルト形式は、`https://search-domainName-unstructuredIdContent.region.es.amazonaws.com` です。ドメインエンドポイントの識別の詳細については、Amazon OpenSearch Service ドキュメントの「[Amazon OpenSearch Service ドメインの作成と管理](#)」を参照してください。

ドメインの HTTP 基本認証情報、*aosUser*、および *aosPassword* を特定または生成します。

OpenSearch Service に対する接続を設定するには:

1. AWS Secrets Manager で、OpenSearch Service 認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
  - [key/value ペア] を選択する際に、*aosUser* という値を持つキー `opensearch.net.http.auth.user` のペアを作成します。
  - [key/value ペア] を選択する際に、*aosPassword* という値を持つキー `opensearch.net.http.auth.pass` のペアを作成します。
2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、将来的に AWS Glue で使用するために、接続名 *connectionName* を維持します。
  - [接続タイプ] を選択する場合は、[OpenSearch Service] を選択します。
  - ドメインエンドポイントを選択する場合は、*aosEndpoint* を入力します。
  - ポートを選択する場合は、*aosPort* を入力します。
  - [AWS Secret] をクリックして、*secretName* を入力します。

AWS Glue OpenSearch Service 接続を作成した後、AWS Glue ジョブを実行する前に次のステップを実行する必要があります。

- AWS Glue ジョブに関連付けられている IAM ロールに *secretName* を読み取るアクセス許可を付与します。
- AWS Glue ジョブ設定で、追加のネットワーク接続として *connectionName* を指定します。

OpenSearch Service のインデックスからの読み取り

前提条件:

- 読み取り元とする OpenSearch Service インデックスである *aosIndex*。
- 認証およびネットワーク位置情報を提供するように設定された AWS Glue OpenSearch Service 接続。これを取得するには、前の手順の「OpenSearch Service に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。

この例では、Amazon OpenSearch Service からインデックスを読み取ります。pushdown パラメータを入力する必要があります。

例:

```
opensearch_read = glueContext.create_dynamic_frame.from_options(
 connection_type="opensearch",
 connection_options={
 "connectionName": "connectionName",
 "opensearch.resource": "aosIndex",
 "pushdown": "true",
 }
)
```

クエリ文字列を指定して、DynamicFrame で返される結果をフィルタリングすることもできます。opensearch.query を設定する必要があります。

opensearch.query は、URL クエリパラメータ文字列である *queryString* またはクエリ DSL JSON オブジェクトである *queryObject* を受け取ることができます。クエリ DSL の詳細については、OpenSearch ドキュメントの「[クエリ DSL](#)」を参照してください。URL クエリパラメータ文字列を入力するには、完全修飾 URL の場合と同様に、クエリの前に ?q= を付加します。クエリ DSL オブジェクトを入力するには、入力する前に JSON オブジェクトを文字列エスケープします。

例:

```
 queryObject = "{ \"query\": { \"multi_match\": { \"query\": \"Sample\", \"fields\":
[\"sample\"] } } }"
 queryString = "?q=queryString"

 opensearch_read_query = glueContext.create_dynamic_frame.from_options(
 connection_type="opensearch",
 connection_options={
 "connectionName": "connectionName",
 "opensearch.resource": "aosIndex",
 "opensearch.query": queryString,
 "pushdown": "true",
 }
)
```

特定の構文以外でクエリを作成する方法の詳細については、OpenSearch ドキュメントの「[クエリ文字列構文](#)」を参照してください。

配列型データを含む OpenSearch コレクションから読み取る場合

は、`opensearch.read.field.as.array.include` パラメーターを使用してメソッド呼び出しでどのフィールドが配列型であるかを指定する必要があります。

例えば、次のドキュメントを読み取ると、`genre` と `actor` 配列フィールドが表示されます。

```
{
 "_index": "movies",
 "_id": "2",
 "_version": 1,
 "_seq_no": 0,
 "_primary_term": 1,
 "found": true,
 "_source": {
 "director": "Frankenheimer, John",
 "genre": [
 "Drama",
 "Mystery",
 "Thriller",
 "Crime"
],
 "year": 1962,
 "actor": [
 "Lansbury, Angela",
 "Sinatra, Frank",
 "Leigh, Janet",
 "Harvey, Laurence",
 "Silva, Henry",
 "Frees, Paul",
 "Gregory, James",
 "Bissell, Whit",
 "McGiver, John",
 "Parrish, Leslie",
 "Edwards, James",
 "Flowers, Bess",
 "Dhiegh, Khigh",
 "Payne, Julie",
 "Kleeb, Helen",
 "Gray, Joe",
 "Nalder, Reggie",
 "Stevens, Bert",
 "Masters, Michael",
 "Lowell, Tom"
]
 }
}
```

```
],
 "title": "The Manchurian Candidate"
 }
}
```

この場合、メソッド呼び出しにそれらのフィールド名を含めます。例:

```
"opensearch.read.field.as.array.include": "genre,actor"
```

配列フィールドがドキュメント構造内にネストされている場合は、ドット表記を使用して参照してください: "genre,actor,foo.bar.baz"。これにより、埋め込みドキュメント bar を含む埋め込みドキュメント foo を介してソースドキュメントに含まれる配列 baz が指定されます。

### OpenSearch Service テーブルへの書き込み

この例では、既存の DynamicFrame である *dynamicFrame* から OpenSearch Service に情報を書き込みます。インデックスに既に情報がある場合、AWS Glue は DynamicFrame からのデータを付加します。pushdown パラメータを入力する必要があります。

前提条件:

- 書き込み先とする OpenSearch Service テーブル。テーブルの識別情報が必要になります。これに *tableName* と名前を付けます。
- 認証およびネットワークの場所の情報を提供するように設定された AWS Glue OpenSearch Service 接続。これを取得するには、前の手順の「OpenSearch Service に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。

例:

```
glueContext.write_dynamic_frame.from_options(
 frame=dynamicFrame,
 connection_type="opensearch",
 connection_options={
 "connectionName": "connectionName",
 "opensearch.resource": "aosIndex",
 },
)
```

## OpenSearch Service 接続のオプションのリファレンス

- `connectionName` — 必須。読み込み/書き込みに使用されます。認証およびネットワークの場所の情報を接続方法に提供するように設定された AWS Glue OpenSearch Service 接続の名前。
- `opensearch.resource` — 必須。読み込み/書き込みに使用されます。有効な値: OpenSearch インデックス名。接続メソッドがインタラクションするインデックスの名前。
- `opensearch.query` — 読み取りに使用。有効な値: JSON をエスケープした文字列、またはこの文字列が ? で始まる場合は URL の検索部分。読み取り時に取得する内容をフィルタリングする OpenSearch クエリ。このパラメータの使用の詳細については、前のセクション「[the section called “OpenSearch Service から読み取る”](#)」を参照してください。
- `pushdown` — 次の場合は必須です。読み込みに使用されます。有効な値: ブール値。Spark に読み取りクエリを OpenSearch に渡し、データベースが関連ドキュメントのみを返すように指示します。
- `opensearch.read.field.as.array.include` — 配列タイプのデータを読み取る場合は必須。読み込みに使用されます。有効な値: カンマで区切られた、フィールド名のリスト。OpenSearch ドキュメントから配列として読み取るフィールドを指定します。このパラメータの使用の詳細については、前のセクション「[the section called “OpenSearch Service から読み取る”](#)」を参照してください。

## Redshift 接続

AWS Glue for Spark を使用して、Amazon Redshift データベース内のテーブルへの読み取りとテーブルへの書き込みを行うことができます。Amazon Redshift データベースに接続するとき、AWS Glue は Amazon Redshift SQL COPY UNLOAD とコマンドを使用してスループットを最大化するために Amazon S3 を介してデータを移動します。AWS Glue 4.0 以降では、[Apache Spark 用の Amazon Redshift インテグレーションを使用して](#)、Amazon Redshift 固有の最適化と機能を使用して読み取りと書き込みを行うことができます。また、以前のバージョンで接続する場合には利用できなかった機能もあります。

AWS Glue によって Amazon Redshift ユーザーがサーバーレスデータ統合と ETL のための AWS Glue への移行がこれまでになく簡単になった方法をご覧ください。

## Redshift 接続の設定

AWS Glue で Amazon Redshift クラスターを使用するには、いくつかの前提条件が必要です。

- データベースの読み込みおよび書き込み時に一時的なストレージとして使用する Amazon S3 デイレクトリがあること。

- Amazon Redshift クラスター、AWS Glue ジョブ、および Amazon S3 デイレトリ間の通信を可能にする Amazon VPC。
- AWS Glue ジョブと Amazon Redshift クラスターに対する適切な IAM アクセス許可。

## IAM ロールを設定する

### Amazon Redshift クラスターのロールを設定する

AWS Glue ジョブと統合するには、Amazon Redshift クラスターが Amazon S3 への読み取りと書き込みが可能である必要があります。これを可能にするには、接続する Amazon Redshift クラスターに IAM ロールを関連付けます。ロールには、Amazon S3 の一時ディレクトリからの読み取り、および一時ディレクトリへの書き込みを許可するポリシーが必要です。ロールには、`redshift.amazonaws.com` サービスで `AssumeRole` を許可する信頼関係が必要です。

### IAM ロールを Amazon Redshift に関連付けるには

1. 前提条件: ファイルの一時的なストレージとして使用する Amazon S3 バケットまたはディレクトリがあること。
2. Amazon Redshift クラスターに必要な Amazon S3 のアクセス許可を特定します。Amazon Redshift クラスターとの間でデータを移動するとき、AWS Glue ジョブは Amazon Redshift に対して COPY ステートメントと UNLOAD ステートメントを発行します。ジョブによって Amazon Redshift のテーブルが変更された場合、AWS Glue は CREATE LIBRARY ステートメントも発行します。Amazon Redshift がこれらのステートメントを実行するために必要な特定の Amazon S3 アクセス権限については、Amazon Redshift のドキュメント「[Amazon Redshift: 他のリソースにアクセスするためのアクセス許可](#)」を参照してください。AWS
3. IAM コンソールで、必要なアクセス許可を持つ IAM ポリシーを作成します。ポリシーの作成についての詳細は、「[IAM ポリシーの作成](#)」を参照してください。
4. IAM コンソールで、ロールと信頼関係を作成し、Amazon Redshift がロールを継承できるようにします。IAM ドキュメントの指示に従ってサービス ([コンソール](#)) のロールを作成するには [AWS](#)
  - AWS サービスのユースケースを選択するように求められたら、「Redshift-カスタマイズ可能」を選択します。
  - ポリシーをアタッチするように求められた場合、以前に定義したポリシーを選択します。

**Note**

Amazon Redshift のロールの設定の詳細については、Amazon Redshift ドキュメントの「[Amazon Redshift AWS がユーザーに代わって他のサービスにアクセスすることを許可する](#)」を参照してください。

5. Amazon Redshift コンソールで、ロールを Amazon Redshift クラスターに関連付けます。これについては、[Amazon Redshift のドキュメント](#)に記載されている手順に従ってください。

Amazon Redshift コンソールでハイライト表示されているオプションを選択して、この設定を行います。

The screenshot shows the Amazon Redshift console interface for a cluster named 'flight-2016'. The breadcrumb navigation is 'Amazon Redshift > Clusters > flight-2016'. The cluster name 'flight-2016' is prominently displayed. Below it, there are several tabs: 'Cluster performance', 'Query monitoring', and 'Properties'. The 'General information' section is visible, showing details like 'Cluster identifier: flight-2016', 'Status: Available', 'Cluster namespace', 'Date created', 'Storage used: 0.25% (0.41 of 160 GB)', and 'Multi-AZ: No'. An 'Actions' dropdown menu is open, listing various management options. The 'Permissions' section of the menu is expanded, and the 'Manage IAM roles' option is highlighted with a red rectangular box. Other options in the menu include 'Manage cluster', 'Backup and disaster recovery', and 'Permissions'.

**Note**

デフォルトでは、AWS Glue ジョブは、ジョブを実行するために指定したロールを使用して作成された Amazon Redshift の一時認証情報を渡します。これらの認証情報を使用すること

はお勧めしません。これらの認証情報は、セキュリティ上の理由により 1 時間後に失効します。

## AWS Glue ジョブのロールを設定する

AWS Glue ジョブには Amazon S3 バケットにアクセスするためのロールが必要です。Amazon Redshift クラスターの IAM アクセス許可は必要ありません。アクセスは Amazon VPC の接続とデータベースの認証情報によって制御されます。

### Amazon VPC を設定する

Amazon Redshift データストアへのアクセスをセットアップするには

1. AWS Management Console にサインインし、<https://console.aws.amazon.com/redshiftv2/> にある Amazon Redshift コンソールを開きます。
2. 左のナビゲーションペインで [クラスター] を選択します。
3. AWS Glue からアクセスするクラスターの名前を選択します。
4. [Cluster Properties (クラスターのプロパティ)] セクションで、[VPC security groups (VPC セキュリティグループ)] 内のセキュリティグループを選択し、AWS Glue が使用できるようにします。今後の参照用に選択したセキュリティグループの名前を記録します。Amazon EC2 コンソールでセキュリティグループを選択すると、[Security Groups] (セキュリティグループ) の一覧が開きます。
5. 変更するセキュリティグループを選択し、[Inbound (インバウンド)] タブに移動します。
6. 自己参照ルールを追加して、AWS Glue コンポーネントが通信できるようにします。具体的には、[Type (タイプ)] All TCP、[Protocol (プロトコル)] は TCP、[Port Range (ポート範囲)] にはすべてのポートが含まれ、[Source (ソース)] は [Group ID (グループ ID)] と同じセキュリティグループ名であるというルールを追加または確認します。

インバウンドルールは以下のようになります。

タイプ	プロトコル	ポート範囲	ソース
すべての TCP	TCP	0 ~ 65535	database-security-group

例:

7. アウトバウンドトラフィックのルールも追加します。すべてのポートへのアウトバウンドトラフィックを開きます。以下に例を示します。

タイプ	プロトコル	ポート範囲	デスティネーション
すべてのトラフィック	すべて	すべて	0.0.0.0/0

または、[Type (タイプ)] は All TCP、[Protocol (プロトコル)] は TCP、[Port Range (ポート範囲)] にすべてのポートが含まれ、[Destination (宛先)] は [Group ID (グループ ID)] と同じセキュリティグループ名の自己参照ルールを作成します。Amazon S3 VPC エンドポイントを使用している場合は、Amazon S3 にアクセスするための HTTPS ルールも追加します。VPC から Amazon *S3 VPC ##### s3-prefix-list-id* が必要です。

例:

タイプ	プロトコル	ポート範囲	デスティネーション
すべての TCP	TCP	0 ~ 65535	<i>security-group</i>
HTTPS	TCP	443	<i>s3-prefix-list-id</i>

## AWS Glue をセットアップ

Amazon VPC 接続情報を提供する AWS Glue データカタログ接続を作成する必要があります。

コンソールで AWS Glue への Amazon Redshift Amazon VPC 接続を設定するには

- 「[the section called “AWS Glue 接続の追加”](#)」のステップに従ってデータカタログ接続を作成します。接続を作成したら、次のステップのために接続名 *connectionName* を保存しておきます。
  - 接続タイプを選択するときは、Amazon Redshift を選択します。
  - Redshift クラスターを選択するときは、クラスターを名前で選択します。

- クラスターの Amazon Redshift ユーザーのデフォルト接続情報を指定します。
- Amazon VPC 設定は自動的に設定されます。

**Note**

AWS SDK を使用して Amazon Redshift 接続を作成する場合は、Amazon VPC の `PhysicalConnectionRequirements` を手動で指定する必要があります。

2. AWS Glue ジョブ設定で、追加のネットワーク接続として `ConnectionName #####`。

例: Amazon Redshift テーブルからの読み取り

Amazon Redshift クラスターと Amazon Redshift Serverless 環境から読み込みを行うことができます。

前提条件: 読み込む目的の Amazon Redshift テーブルがあること。 [the section called “Redshift の設定”](#) 前のセクションの手順に従うと、一時ディレクトリ用の Amazon S3 URI、 `temp-s3-dir`、および IAM ロール (アカウント内) が用意されているはずで、 `rs-role-namerole-account-id`

### Using the Data Catalog

その他の前提条件: 読み取る目的の Amazon Redshift テーブル用のテーブルとデータカタログのデータベースがあること。データカタログの詳細については、「[データ検出とカタログ化](#)」を参照してください Amazon Redshift テーブルのエントリを作成したら、 `redshift-dc-database-nameredshift-table-name` 接続をおよびで識別します。

設定: 関数オプションで、 `database` および `table_name` パラメータを使用してデータカタログテーブルを識別します。 `redshift_tmp_dir` を使用して、Amazon S3 の一時ディレクトリを識別します。また、 `rs-role-nameaws_iam_roleadditional_options` パラメータ内のキーを使用して指定することもできます。

```
glueContext.create_dynamic_frame.from_catalog(
 database = "redshift-dc-database-name",
 table_name = "redshift-table-name",
 redshift_tmp_dir = args["temp-s3-dir"],
 additional_options = {"aws_iam_role": "arn:aws:iam::role-account-id:role/rs-
role-name"})
```

## Connecting directly

その他の前提条件: Amazon Redshift テーブルの名前が必要になります (. *redshift-table-name* そのテーブルを格納する Amazon Redshift クラスターの JDBC 接続情報が必要です。####  
#####*redshift-database-name*

Amazon Redshift クラスターを操作するときは、Amazon Redshift コンソールから接続情報を取得できます。Amazon Redshift サーバーレスを使用する場合は、Amazon Redshift ドキュメントの「[Amazon Redshift Serverless への接続](#)」を参照してください。

設定: 関数オプションでは、url、dbtable、user および password を使用して接続パラメータを識別します。redshift\_tmp\_dir を使用して、Amazon S3 の一時ディレクトリを識別します。from\_options を使用する場合は IAM ロールを指定できます。この構文はデータカタログを介して接続するのと似ていますが、connection\_options マップにパラメータを配置します。

パスワードを AWS Glue スクリプトにハードコーディングするのは悪い習慣です。SDK for Python (Boto3) を使用して、AWS Secrets Manager パスワードをスクリプトに保存し、スクリプト内で取得することを検討してください。

```
my_conn_options = {
 "url": "jdbc:redshift://host:port/redshift-database-name",
 "dbtable": "redshift-table-name",
 "user": "username",
 "password": "password",
 "redshiftTmpDir": args["temp-s3-dir"],
 "aws_iam_role": "arn:aws:iam::account id:role/rs-role-name"
}

df = glueContext.create_dynamic_frame.from_options("redshift", my_conn_options)
```

### 例: Amazon Redshift テーブルへの書き込み

Amazon Redshift クラスターと Amazon Redshift Serverless 環境に書き込みを行うことができます。

前提条件: Amazon Redshift クラスターを作成し、[the section called “Redshift の設定”](#)前のセクションの手順に従うと、一時ディレクトリ用の Amazon S3 URI、*temp-s3-dir*、および IAM ロール (ア

ウント内) が必要です。 *rs-role-namerole-account-id* また、データベースに書き込むコンテンツがある DynamicFrame も必要です。

## Using the Data Catalog

その他の前提条件 Amazon Redshift クラスター用のデータカタログのデータベースと、読み取る目的のテーブルがあること。データカタログの詳細については、「[データ検出とカタログ化](#)」を参照してください 接続は *dc-connection-name* で、ターゲットテーブルは *redshift-database-name/redshift-table-name* で指定します。

設定: 関数オプションで、*database* パラメータを使用してデータカタログのデータベースを識別し、*table\_name* を使用してテーブルを指定します。 *redshift\_tmp\_dir* を使用して、Amazon S3 の一時ディレクトリを識別します。また、 *rs-role-name* *aws\_iam\_role* *additional\_options* パラメータ内のキーを使用して指定することもできます。

```
glueContext.write_dynamic_frame.from_catalog(
 frame = input dynamic frame,
 database = "redshift-dc-database-name",
 table_name = "redshift-table-name",
 redshift_tmp_dir = args["temp-s3-dir"],
 additional_options = {"aws_iam_role": "arn:aws:iam::account-id:role/rs-role-name"}
)
```

## Connecting through a AWS Glue connection

*write\_dynamic\_frame.from\_options* メソッドを使用して、Amazon Redshift に直接接続できます。ただし、接続の詳細をスクリプトに直接挿入するのではなく、*from\_jdbc\_conf* メソッドを使用してデータカタログの接続に保存されている接続の詳細を参照します。これは、データベースのデータカタログテーブルをクロールまたは作成することなく実行できます。データカタログの接続についての詳細は、「[データへの接続](#)」を参照してください。

その他の前提条件: データベース用のデータカタログの接続と、読み取る目的の Amazon Redshift テーブルがあること。

設定: データカタログとの接続を指定します *dc-connection-name*。 Amazon Redshift データベースとテーブルは、 *redshift-table-name/redshift-database-name* とで識別します。 *catalog\_connection* を使用してデータカタログの接続情報を指定し、*dbtable* および

database を使用して Amazon Redshift の情報を指定します。この構文はデータカタログを介して接続するのと似ていますが、connection\_options マップにパラメータを配置します。

```
my_conn_options = {
 "dbtable": "redshift-table-name",
 "database": "redshift-database-name",
 "aws_iam_role": "arn:aws:iam::role-account-id:role/rs-role-name"
}

glueContext.write_dynamic_frame.from_jdbc_conf(
 frame = input_dynamic_frame,
 catalog_connection = "dc-connection-name",
 connection_options = my_conn_options,
 redshift_tmp_dir = args["temp-s3-dir"])
```

## Amazon Redshift 接続のオプションに関するリファレンス

user やなどの情報を設定するためにすべての AWS Glue JDBC 接続に使用される基本的な接続オプションは url、すべての JDBC password タイプで一貫しています。JDBC のパラメータの詳細については、「[the section called “JDBC 接続パラメータ”](#)」を参照してください。

Amazon Redshift 接続のタイプには、その他の接続オプションがいくつかあります。

- "redshiftTmpDir": (必須) データベースからコピーする際に一時的なデータをステージングできる Amazon S3 パス。
- "aws\_iam\_role": (オプション) IAM ロールの ARN。AWS Glue ジョブはこのロールを Amazon Redshift クラスターに渡し、ジョブの指示を完了するために必要なアクセス権限をクラスターに付与します。

AWS Glue 4.0+ ではその他の接続オプションも利用可能

AWS Glue 接続オプションを使用して、新しい Amazon Redshift コネクタのオプションを渡すこともできます。サポートされているコネクタオプションの詳細なリストについては、「[Apache Spark 用の Amazon Redshift の統合](#)」の「Spark SQL パラメータ」セクションを参照してください。

参考までに、いくつかの新しいオプションをここで再度確認します。

名前	必須	デフォルト	[Description] (説明)
autopushdown	いいえ	TRUE	SQL オペレーションの Spark 論理プランをキャプチャして分析することにより、述語とクエリのプッシュダウンを適用します。オペレーションは SQL クエリに変換され、Amazon Redshift で実行されることでパフォーマンスが向上します。
autopushdown.s3_result_cache	いいえ	FALSE	SQL クエリをキャッシュして Amazon S3 パスマッピングのデータをメモリにアンロードします。これにより、同じクエリを同じ Spark セッションで再度実行する必要がなくなります。autopushdown が有効な場合のみサポートされます。
unload_s3_format	いいえ	[PARQUET]	PARQUET - クエリ結果をパーケット形式でアンロードします。  TEXT - クエリ結果をパイプ区切りのテキスト形式でアンロードします。

名前	必須	デフォルト	[Description] (説明)
sse_kms_key	いいえ	該当なし	AWS のデフォルトの暗号化の代わりに、UNLOADオペレーション中の暗号化に使用する SSE-KMS キー。AWS
extracopyoptions	いいえ	該当なし	<p>データロード時に Amazon Redshift COPY コマンドに付加する追加オプション (TRUNCATECOLUMNS、MAXERRORn など) のリスト (その他のオプションについては、「<a href="#">COPY: 任意指定のパラメータ</a>」を参照してください)。</p> <p>これらのオプションはCOPYコマンドの最後に付加されるため、使用できるのはコマンドの最後で意味のあるオプションのみであることに注意してください。これで、考えられるほとんどのユースケースがカバーされるはずです。</p>

名前	必須	デフォルト	[Description] (説明)
csvnullstring (実験的)	いいえ	NULL	CSV tempformat を使用するとき NULL として書き込む文字列値。これは実際のデータには出現しない値でなければなりません。

これらの新しいパラメータは、次の方法で使用できます。

#### パフォーマンス向上のための新しいオプション

新しいコネクタには、パフォーマンス向上のための新しいオプションがいくつか導入されています。

- autopushdown: デフォルトでは有効になっています。
- autopushdown.s3\_result\_cache: デフォルトでは無効になっています。
- unload\_s3\_format: デフォルトでは PARQUET になっています。

これらのオプションの使用方法については、「[Apache Spark 用の Amazon Redshift の統合](#)」を参照してください。キャッシュされた結果には古い情報が含まれている可能性があるため、読み取りと書き込みのオペレーションが混在している場合は `autopushdown.s3_result_cache` をオンにしないことをお勧めします。パフォーマンスを向上させ、ストレージコストを削減するために、UNLOAD コマンドの `unload_s3_format` オプションはデフォルトで PARQUET に設定されています。UNLOAD コマンドのデフォルト動作を使用するには、オプションを TEXT にリセットします。

#### 読み取り用の新しい暗号化オプション

デフォルトでは、Amazon Redshift テーブルからデータを読み取る際に AWS Glue が使用する一時フォルダ内のデータは、SSE-S3 暗号化を使用して暗号化されます。AWS Key Management Service (AWS KMS) の顧客管理キーを使用してデータを暗号化するには、バージョン 3.0 ("`sse_kms_key`" # `kmsKey`) の従来の設定オプションの代わりに AWS KMS、[KMSKey がキー ID の送信元となるように設定できます](#)。 ("`extraunloadoptions`" # `s"ENCRYPTED KMS_KEY_ID '$kmsKey'"`) AWS Glue

```
datasource0 = glueContext.create_dynamic_frame.from_catalog(
```

```

database = "database-name",
table_name = "table-name",
redshift_tmp_dir = args["TempDir"],
additional_options = {"sse_kms_key": "<KMS_KEY_ID>"},
transformation_ctx = "datasource0"
)

```

## IAM ベースの JDBC URL をサポート

新しいコネクタは IAM ベースの JDBC URL をサポートしているため、ユーザーやパスワード、またはシークレットを渡す必要はありません。IAM ベースの JDBC URL では、コネクタはジョブのランタイムロールを使用して Amazon Redshift データソースにアクセスします。

ステップ 1: 次の最低限必要なポリシーを AWS Glue ジョブランタイムロールにアタッチします。

```

{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "VisualEditor0",
 "Effect": "Allow",
 "Action": "redshift:GetClusterCredentials",
 "Resource": [
 "arn:aws:redshift:<region>:<account>:dbgroup:<cluster name>/*",
 "arn:aws:redshift:*:<account>:dbuser:/*/*",
 "arn:aws:redshift:<region>:<account>:dbname:<cluster name>/<database name>"
]
 },
 {
 "Sid": "VisualEditor1",
 "Effect": "Allow",
 "Action": "redshift:DescribeClusters",
 "Resource": "*"
 }
]
}

```

ステップ 2: IAM ベースの JDBC URL を次のように使用します。接続している Amazon Redshift ユーザー名で新しいオプション DbUser を指定します。

```

conn_options = {
 // IAM-based JDBC URL

```

```
"url": "jdbc:redshift:iam://<cluster name>:<region>/<database name>",
"dbtable": dbtable,
"redshiftTmpDir": redshiftTmpDir,
"aws_iam_role": aws_iam_role,
"DbUser": "<Redshift User name>" // required for IAM-based JDBC URL
}

redshift_write = glueContext.write_dynamic_frame.from_options(
 frame=dyf,
 connection_type="redshift",
 connection_options=conn_options
)

redshift_read = glueContext.create_dynamic_frame.from_options(
 connection_type="redshift",
 connection_options=conn_options
)
```

### Note

DynamicFrame は現在、GlueContext.create\_dynamic\_frame.from\_options ワークフローに DbUser が含まれている IAM ベースの JDBC URL のみをサポートしていません。

## AWS Glue バージョン 3.0 からバージョン 4.0 への移行

AWS Glue 4.0 では、ETL ジョブは新しい Amazon Redshift Spark コネクタと、さまざまなオプションと設定を備えた新しい JDBC ドライバーにアクセスできます。新しい Amazon Redshift コネクタとドライバーはパフォーマンスを念頭に置いて作成されており、データのトランザクションの一貫性が保たれます。これらの製品については、Amazon Redshift のドキュメントに記載されています。詳細については、以下を参照してください。

- [Apache Spark 用の Amazon Redshift インテグレーション](#)
- [Amazon Redshift JDBC ドライバー、バージョン 2.1](#)

## テーブル/列名と識別子の制限

新しい Amazon Redshift Spark コネクタとドライバーで、Redshift テーブル名の要件が厳しくなっています。Amazon Redshift のテーブル名定義の詳細については、「[名前と識別子](#)」を参照してください

い。ジョブブックマークのワークフローは、ルールに一致しないテーブル名やスペースなどの特定の文字があると機能しない場合があります。

[名前と識別子](#)のルールに適合しない名前のレガシーテーブルがあり、ブックマーク (古い Amazon Redshift テーブルデータを再処理するジョブ) に問題がある場合は、テーブル名を変更することをお勧めします。詳細については、「[ALTER TABLE の例](#)」を参照してください。

### データフレームのデフォルトの一時形式の変更

AWS Glue バージョン 3.0 の Spark コネクタでは、Amazon Redshift への書き込み時にデフォルトで `tempformat` が CSV に設定されます。一貫性を保つために、AWS Glue バージョン 3.0 でも、`DynamicFrame` ではデフォルトで `tempformat` に CSV を使用するようになっています。以前に Amazon Redshift Spark コネクタで Spark データフレーム API を直接使用していた場合は、`DataframeReader/Writer` オプションで `tempformat` を CSV に明示的に設定できます。それ以外の場合は、新しい Spark コネクタで `tempformat` はデフォルトで AVRO になります。

**動作の変更:** Amazon Redshift データ型 REAL を Spark データ型 DOUBLE ではなく FLOAT にマッピングする

AWS Glue バージョン 3.0 では、Amazon Redshift の REAL は Spark の DOUBLE 型に変換されます。新しい Amazon Redshift Spark コネクタでは、動作が更新され、Amazon Redshift の REAL 型が Spark の FLOAT 型に変換されたり、その逆の変換が行われたりするようになりました。従来のユースケースで引き続き Amazon Redshift の REAL 型を Spark の DOUBLE 型にマッピングする場合は、次の回避策を使用できます。

- `DynamicFrame` の場合、`DynamicFrame.ApplyMapping` で Float 型を Double 型にマッピングします。`Dataframe` の場合、`cast` を使用する必要があります。

コードサンプル:

```
dyf_cast = dyf.apply_mapping([('a', 'long', 'a', 'long'), ('b', 'float', 'b', 'double')])
```

### Kafka 接続

Kafka クラスターまたは Amazon Managed Streaming for Apache Kafka への接続を指定します。

Data Catalog テーブルに格納されている情報を使用するか、データストリームに直接アクセスするための情報を指定することにより、Kafka データストリームへ読み込むまたは書き込むことが

できます。Kafka から Spark に情報を読み取り DataFrame、それを Glue に変換できます。AWS DynamicFrameKafka には JSON DynamicFrames 形式で書き込むことができます。データストリームに直接アクセスする場合は、これらのオプションを使用して、データストリームへのアクセス方法に関する情報を提供します。

getCatalogSource または create\_data\_frame\_from\_catalog を使用して Kafka ストリーミングソースからレコードを消費するか、getCatalogSink または write\_dynamic\_frame\_from\_catalog を使用して Kafka にレコードを書き込み、ジョブに Data Catalog データベースおよびテーブル名の情報があり、それを使用して Kafka ストリーミングソースから読み込むためのいくつかの基本パラメータを取得できる場合。getSource、getCatalogSink、getSourceWithFormat、getSinkWithFormat、createDataFrame を使用する場合、ここで説明する接続オプションを使用し、これらの基本パラメータを指定する必要があります。

GlueContext クラスで指定されたメソッドの次の引数を使用し、Kafka の接続オプションを指定できます。

- Scala
  - connectionOptions: getSource、createDataFrameFromOptions、getSink で使用
  - additionalOptions: getCatalogSource、getCatalogSink で使用
  - options: getSourceWithFormat、getSinkWithFormat で使用
- Python
  - connection\_options:  
create\_data\_frame\_from\_options、write\_dynamic\_frame\_from\_options で使用
  - additional\_options:  
create\_data\_frame\_from\_catalog、write\_dynamic\_frame\_from\_catalog で使用
  - options: getSource、getSink で使用

ストリーミング ETL ジョブに関する注意事項と制限事項については、「[the section called “ストリーミング ETL に関する注意と制限”](#)」を参照してください。

## Kafka の設定

インターネット経由で利用可能な Kafka AWS ストリームに接続するための前提条件はありません。

AWS Glue Kafka 接続を作成して、接続認証情報を管理できます。詳細については、「[the section called “Kafka データストリームの接続の作成”](#)」を参照してください。AWS Glue ジョ

ブ設定で、`ConnectionName #####`、メソッド呼び出しでパラメーターに `ConnectionName #####`。connectionName

場合によっては、追加の前提条件を設定する必要があります。

- IAM 認証で Amazon Managed Streaming for Apache Kafka を使用する場合は、適切な IAM 設定が必要になります。
- Amazon VPC で Amazon Managed Streaming for Apache Kafka を使用する場合は、適切な Amazon VPC 設定が必要になります。Amazon VPC 接続情報を提供する AWS Glue 接続を作成する必要があります。AWS Glue 接続を追加ネットワーク接続として含めるには、ジョブ設定が必要です。

ストリーミング ETL ジョブの前提条件の詳細については、「[the section called “ストリーミング ETL ジョブ”](#)」を参照してください。

例: Kafka ストリームからの読み込み

[the section called “forEachBatch”](#) と組み合わせて使用します。

Kafka ストリーミングソースの例：

```
kafka_options =
 { "connectionName": "ConfluentKafka",
 "topicName": "kafka-auth-topic",
 "startingOffsets": "earliest",
 "inferSchema": "true",
 "classification": "json"
 }
data_frame_datasource0 =
 glueContext.create_data_frame.from_options(connection_type="kafka",
 connection_options=kafka_options)
```

例: Kafka ストリームへの書き込み

Kafka への書き込み例:

getSink メソッドを使った例:

```
data_frame_datasource0 =
 glueContext.getSink(
```

```

connectionType="kafka",
connectionOptions={
 JsonOptions("""{
 "connectionName": "ConfluentKafka",
 "classification": "json",
 "topic": "kafka-auth-topic",
 "typeOfData": "kafka"}
""")),
transformationContext="dataframe_ApacheKafka_node1711729173428")
.getDataFrame()

```

`write_dynamic_frame.from_options` メソッドを使った例:

```

kafka_options =
 { "connectionName": "ConfluentKafka",
 "topicName": "kafka-auth-topic",
 "classification": "json"
 }
data_frame_datasource0 =
 glueContext.write_dynamic_frame.from_options(connection_type="kafka",
 connection_options=kafka_options)

```

## Kafka 接続オプションのリファレンス

読み込むとき、`"connectionType": "kafka"` に次の接続オプションを使用します。

- `"bootstrap.servers"` (必須) ブートストラップサーバーの URL のリスト (例: `b-1.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094`)。このオプションは API 呼び出しで指定するか、データカタログ内のテーブルメタデータで定義する必要があります。
- `"security.protocol"` (必須) ブローカーと通信するために使用されるプロトコル。使用できる値は、`"SSL"` または `"PLAINTEXT"` です。
- `"topicName"` (必須) サブスクライブするトピックのカンマ区切りリスト。"`topicName`"、"`assign`"、または "`subscribePattern`" の中から、いずれか 1 つのみを指定する必要があります。
- `"assign"`: (必須) 消費する特定の `TopicPartitions` を指定する JSON 文字列。"`topicName`"、"`assign`"、または "`subscribePattern`" の中から、いずれか 1 つのみを指定する必要があります。

例: `'{"topicA":[0,1],"topicB":[2,4]}'`

- "subscribePattern": (必須) サブスクライブする先のトピックリストを識別する Java の正規表現文字列。"topicName"、"assign"、または "subscribePattern" の中から、いずれか 1 つのみを指定する必要があります。

例: 'topic.\*'

- "classification" (必須) レコード内のデータで使用されるファイル形式。データカタログを通じて提供されていない限り、必須です。
- "delimiter" (オプション) classification が CSV の場合に使用される値の区切り文字。デフォルトは「,」です。
- "startingOffsets": (オプション) Kafka トピック内で、データの読み取りを開始する位置 使用できる値は、"earliest" または "latest" です。デフォルト値は "latest" です。
- "startingTimestamp": (オプション、AWS Glue バージョン 4.0 以降でのみサポート) データを読み取る Kafka トピック内のレコードのタイムスタンプ。指定できる値は UTC 形式 (yyyy-mm-ddTHH:MM:SSZ のパターン) のタイムスタンプ文字列です (Z は UTC タイムゾーンのオフセットを +/- で表します。例: 「2023-04-04T08:00:00-04:00」)。

注: AWS Glue ストリーミングスクリプトの接続オプションリストには、「StartingOffsets」または「StartingTimestamp」のいずれかしか表示できません。これらのプロパティの両方を含めると、ジョブが失敗します。

- "endingOffsets": (オプション) バッチクエリの終了位置。設定が可能な値は、"latest" または、各 TopicPartition の終了オフセットを指定する JSON 文字列のいずれかです。

JSON 文字列の場合、{"topicA":{"0":23,"1":-1},"topicB":{"0":-1}} の形式を使用します。オフセットとして値 -1 を設定する場合、"latest" の意味になります。

- "pollTimeoutMs": (オプション) Spark ジョブエグゼキュータで、Kafka からデータをポーリングする際のタイムアウト値 (ミリ秒単位)。デフォルト値は、512 です。
- "numRetries": (オプション) Kafka オフセットのフェッチが失敗したと判断される前の再試行回数。デフォルト値は、3 です。
- "retryIntervalMs": (オプション) Kafka オフセットのフェッチを開始するまでの待機時間 (ミリ秒)。デフォルト値は、10 です。
- "maxOffsetsPerTrigger": (オプション) 処理されるオフセットの最大数を、トリガー間隔ごとのレート上限で指定する値。指定されたオフセットの合計数は、異なるボリュームの topicPartitions 間で均等に分割されます。デフォルト値は「null」です。この場合、コンシューマーは既知の最新のオフセットまで、すべてのオフセットを読み取ります。

- "minPartitions": (オプション) Kafka から読み取ることを想定する、最小のパーティション数。デフォルト値は「null」です。これは、Spark パーティションの数が Kafka パーティションの数に等しいことを意味します。
- "includeHeaders": (オプション) Kafka ヘッダーを含めるかどうかを決定します。このオプションが「true」に設定されている場合、データ出力には、「glue\_streaming\_kafka\_headers」という名前で Array[Struct(key: String, value: String)] 型の列が追加されます。デフォルト値は「false」です。このオプションは、AWS Glue バージョン 3.0 以降でのみ使用可能です。
- "schema": (inferSchema に false を設定した場合は必須) ペイロードの処理に使用するスキーマ。分類が avro である場合、提供されるスキーマには Avro スキーマ形式を使用する必要があります。分類が avro 以外の場合、提供されるスキーマは DDL スキーマ形式である必要があります。

以下に、スキーマの例を示します。

Example in DDL schema format

```
'column1' INT, 'column2' STRING , 'column3' FLOAT
```

Example in Avro schema format

```
{
 "type": "array",
 "items":
 {
 "type": "record",
 "name": "test",
 "fields":
 [
 {
 "name": "_id",
 "type": "string"
 },
 {
 "name": "index",
 "type":
 [
 "int",
 "string",
 "float"
]
 }
]
 }
}
```

```
}
}
```

- "inferSchema": (オプション) デフォルト値は「false」です。「true」に設定すると、実行時に、スキーマが foreachbatch 内のペイロードから検出されます。
- "avroSchema": (非推奨) Avro 形式を使用する場合に、Avro データのスキーマを指定するために使用されるパラメータです。このパラメータは非推奨となりました。schema パラメータを使用します。
- "addRecordTimestamp": (オプション) このオプションを「true」に設定すると、トピックが対応するレコードを受信した時刻を表示する「\_\_src\_timestamp」という列が、データ出力に追加で表示されます。デフォルト値は、「false」です。このオプションは AWS Glue バージョン 4.0 以降でサポートされています。
- "emitConsumerLagMetrics": (オプション) オプションを 'true' に設定すると、バッチごとに、トピックが受信した最も古いレコードからレコードが届くまでの期間のメトリクスが出力されます。AWS Glue CloudWatchメトリクスの名前は「glue.driver.stream」です。maxConsumerLagInMs」。デフォルト値は、「false」です。このオプションは AWS Glue バージョン 4.0 以降でサポートされています。

書き込むとき、"connectionType": "kafka" に次の接続オプションを使用します。

- "connectionName"(必須) Kafka クラスターへの接続に使用される AWS Glue 接続の名前 (Kafka ソースと同様)。
- "topic" (必須) トピック列が存在する場合、トピック設定オプションが設定されていない限り、指定された行を Kafka に書き込む際にトピック列の値がトピックとして使用されます。つまり、topic 設定オプションはトピック列をオーバーライドします。
- "partition" (オプション) 有効なパーティション番号が指定された場合、レコードの送信時にその partition が使用されます。

パーティションが指定されずに key が存在する場合、キーのハッシュを使用してパーティションが選択されます。

key と partition のどちらも存在しない場合、そのパーティションに少なくとも batch.size バイトが生成された際に変更内容のスティッキーパーティション分割に基づいてパーティションが選択されます。

- "key" (オプション) partition が null の場合、パーティション分割に使用されます。

- "classification" (オプション) レコードのデータに使用されるファイル形式。JSON、CSV、Avro のみをサポートしています。

Avro 形式を使用すると、シリアル化するためにカスタムの AvroSchema を提供できますが、シリアル化解除する場合にもソースに提供する必要があることに注意してください。それ以外の場合、デフォルトではシリアル化に Apache を使用します。AvroSchema

さらに、[Kafka プロデューサーの設定パラメーター](#)を更新することにより、必要に応じて Kafka シンクを微調整できます。接続オプションには許可リストがないことに注意してください。すべてのキー値のペアはそのままシンクに保持されます。

ただし、有効にならないオプションの小さな拒否リストがあります。詳細については、「[Apache 固有の設定](#)」を参照してください。

## Azure Cosmos DB 接続

AWS Glue for Spark を使用すると、AWS Glue 4.0 以降のバージョンの NoSQL API を使用して、Azure Cosmos DB の既存のコンテナとの間で読み取りと書き込みを行うことができます。SQL クエリを使用して、Azure Cosmos DB から何を読み取るかを定義できます。AWS Glue 接続を通じて AWS Secrets Manager に保存されている Azure Cosmos DB キーを使用して Azure Cosmos DB に接続します。

Azure Cosmos DB for NoSQL の詳細については、[Azure のドキュメント](#)を参照してください。

## Azure Cosmos DB 接続の設定

AWS Glue から Azure Cosmos DB に接続するには、Azure Cosmos DB キーを作成して AWS Secrets Manager シークレットに保存し、そのシークレットを Azure Cosmos DB AWS Glue 接続に関連付ける必要があります。

前提条件:

- Azure では、AWS Glue で使用する Azure Cosmos DB キー (cosmosKey) を特定または生成する必要があります。詳細については、Azure ドキュメントの「[Azure Cosmos DB のデータへのアクセスをセキュリティで保護する](#)」を参照してください。

Azure Cosmos DB に対する接続を設定するには:

1. AWS Secrets Manager で、Azure Cosmos DB キーを使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの

「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。

- [key/value ペア] を選択する際に、*cosmosKey* という値を持つキー `spark.cosmos.accountKey` のペアを作成します。
2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、将来的に AWS Glue で使用するために、接続名 *connectionName* を維持します。
- [接続タイプ] を選択する際に、[Azure Cosmos DB] を選択します。
  - [AWS Secret] をクリックして、*secretName* を入力します。

AWS Glue Azure Cosmos DB 接続を作成した後、AWS Glue ジョブを実行する前に次のステップを実行する必要があります。

- AWS Glue ジョブに関連付けられている IAM ロールに *secretName* を読み取るアクセス許可を付与します。
- AWS Glue ジョブ設定で、追加のネットワーク接続として *connectionName* を指定します。

## Azure Cosmos DB for NoSQL コンテナからの読み取り

### 前提条件:

- 読み取り元とする Azure Cosmos DB for NoSQL コンテナ。コンテナの識別情報が必要になります。

Azure Cosmos for NoSQL コンテナは、データベースとコンテナによって識別されます。Azure Cosmos for NoSQL API に接続する際には、データベース *cosmosDBName* とコンテナ *cosmosContainerName* の名前を指定する必要があります。

- 認証およびネットワークの場所の情報を提供するように設定された AWS Glue Azure Cosmos DB 接続。これを取得するには、前の手順の「Azure Cosmos DB に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。

### 例:

```
azurecosmos_read = glueContext.create_dynamic_frame.from_options(
```

```
connection_type="azurecosmos",
connection_options={
 "connectionName": connectionName,
 "spark.cosmos.database": cosmosDBName,
 "spark.cosmos.container": cosmosContainerName,
}
)
```

SELECT SQL クエリを指定して、DynamicFrame に返される結果をフィルタリングすることもできます。query を設定する必要があります。

例:

```
azurecosmos_read_query = glueContext.create_dynamic_frame.from_options(
 connection_type="azurecosmos",
 connection_options={
 "connectionName": "connectionName",
 "spark.cosmos.database": cosmosDBName,
 "spark.cosmos.container": cosmosContainerName,
 "spark.cosmos.read.customQuery": "query"
 }
)
```

## Azure Cosmos DB for NoSQL コンテナへの書き込み

この例では、既存の DynamicFrame である *dynamicFrame* から Azure Cosmos DB に情報を書き込みます。コンテナに既に情報がある場合、AWS Glue は DynamicFrame からのデータを付加します。コンテナ内の情報のスキーマが、書き込む情報と異なる場合、エラーが発生します。

前提条件:

- 書き込み先とする Azure Cosmos DB テーブル。コンテナの識別情報が必要になります。接続メソッドを呼び出す前にコンテナを作成する必要があります。

Azure Cosmos for NoSQL コンテナは、データベースとコンテナによって識別されます。Azure Cosmos for NoSQL API に接続する際には、データベース *cosmosDBName* とコンテナ *cosmosContainerName* の名前を指定する必要があります。

- 認証およびネットワークの場所の情報を提供するように設定された AWS Glue Azure Cosmos DB 接続。これを取得するには、前の手順の「Azure Cosmos DB に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。

例:

```
azurecosmos_write = glueContext.write_dynamic_frame.from_options(
 frame=dynamicFrame,
 connection_type="azurecosmos",
 connection_options={
 "connectionName": connectionName,
 "spark.cosmos.database": cosmosDBName,
 "spark.cosmos.container": cosmosContainerName
 }
)
```

## Azure Cosmos DB 接続オプションのリファレンス

- `connectionName` — 必須。読み込み/書き込みに使用されます。認証およびネットワークの場所の情報を接続方法に提供するように設定された AWS Glue Azure Cosmos DB 接続の名前。
- `spark.cosmos.database` — 必須。読み込み/書き込みに使用されます。有効な値: データベース名。Azure Cosmos DB for NoSQL のデータベース名。
- `spark.cosmos.container` — 必須。読み込み/書き込みに使用されます。有効な値: コンテナ名。Azure Cosmos DB for NoSQL のコンテナ名。
- `spark.cosmos.read.customQuery` — 読み取りに使用。有効な値: SELECT SQL クエリ。読み取るドキュメントを選択するためのカスタムクエリ。

## Azure SQL 接続

AWS Glue for Spark を使用して、AWS Glue 4.0 以降のバージョンの Azure SQL Managed Instances のテーブルからの読み取りとテーブルへの書き込みを行うことができます。SQL クエリを使用して、Azure SQL から何を読み取るかを定義できます。AWS Glue 接続を介して AWS Secrets Manager で保存されているユーザーおよびパスワードの認証情報を使用して Azure SQL に接続します。

Azure SQL の詳細については、[Azure SQL のドキュメント](#)を参照してください。

## Azure SQL 接続の設定

AWS Glue から Azure SQL に接続するには、Azure SQL 認証情報を作成して AWS Secrets Manager シークレットに保存し、そのシークレットを Azure SQL AWS Glue 接続に関連付ける必要があります。

Azure SQL に対する接続を設定するには:

1. AWS Secrets Manager で、Azure SQL 認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
  - [key/value ペア] を選択する際に、*azuresqlUsername* という値を持つキー user のペアを作成します。
  - [key/value ペア] を選択する際に、*azuresqlPassword* という値を持つキー password のペアを作成します。
2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、将来的に AWS Glue で使用するために、接続名 *connectionName* を維持します。

- [接続タイプ] を選択する際に、[Azure SQL] を選択します。
- Azure SQL URL を指定する場合は、JDBC エンドポイント URL を入力します。

URL は、次のような形式になります:

```
jdbc:sqlserver://databaseServerName:databasePort;databaseName=azuresqlDBName
```

AWS Glue には次の URL プロパティが必要です。

- *databaseName* – 接続先の Azure SQL のデフォルトデータベース。

Azure SQL Managed Instances の JDBC URL の詳細については、[Microsoft のドキュメント](#)を参照してください。

- [AWS Secret] をクリックして、*secretName* を入力します。

AWS Glue Azure SQL 接続を作成した後、AWS Glue ジョブを実行する前に次のステップを実行する必要があります。

- AWS Glue ジョブに関連付けられている IAM ロールに *secretName* を読み取るアクセス許可を付与します。
- AWS Glue ジョブ設定で、追加のネットワーク接続として *connectionName* を指定します。

## Azure SQL テーブルからの読み取り

### 前提条件:

- 読み取り元とする Azure SQL テーブル。テーブルの識別情報である *databaseName* および *tableIdentifier* が必要になります。

Azure SQL テーブルは、データベース、スキーマ、テーブル名によって識別されます。Azure SQL に接続する際には、データベース名とテーブル名を指定する必要があります。スキーマがデフォルトの「public」でない場合は、スキーマも指定する必要があります。データベースは、*connectionName* の URL プロパティ、dbtable を通じたスキーマおよびテーブル名を介して指定されます。

- 認証情報を提供するように設定された AWS Glue Azure SQL 接続。認証情報を設定するには、前の手順「Azure SQL に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。

### 例:

```
azuresql_read_table = glueContext.create_dynamic_frame.from_options(
 connection_type="azuresql",
 connection_options={
 "connectionName": "connectionName",
 "dbtable": "tableIdentifier"
 }
)
```

SELECT SQL クエリを指定して、DynamicFrame に返される結果をフィルタリングすることもできます。query を設定する必要があります。

### 例:

```
azuresql_read_query = glueContext.create_dynamic_frame.from_options(
 connection_type="azuresql",
 connection_options={
 "connectionName": "connectionName",
 "query": "query"
 }
)
```

## Azure SQL テーブルへの書き込み

この例では、既存の DynamicFrame である *dynamicFrame* から Azure SQL に情報を書き込みます。テーブルに既に情報がある場合、AWS Glue は DynamicFrame からのデータを付加します。

前提条件:

- 書き込み先とする Azure SQL テーブル。テーブルの識別情報である *databaseName* および *tableIdentifier* が必要になります。

Azure SQL テーブルは、データベース、スキーマ、テーブル名によって識別されます。Azure SQL に接続する際には、データベース名とテーブル名を指定する必要があります。スキーマがデフォルトの「public」でない場合は、スキーマも指定する必要があります。データベースは、*connectionName* の URL プロパティ、dbtable を通じたスキーマおよびテーブル名を介して指定されます。

- Azure SQL 認証情報。認証情報を設定するには、前の手順「Azure SQL に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。

例:

```
azuresql_write = glueContext.write_dynamic_frame.from_options(
 connection_type="azuresql",
 connection_options={
 "connectionName": "connectionName",
 "dbtable": "tableIdentifier"
 }
)
```

## Azure SQL 接続オプションのリファレンス

- *connectionName* — 必須。読み込み/書き込みに使用されます。認証情報を接続方法に提供するように設定された AWS Glue Azure SQL 接続の名前。
- *databaseName* — 読み込み/書き込みに使用されます。有効な値: Azure SQL データベース名。接続先の Azure SQL のデータベースの名前。
- *dbtable* — 書き込みの場合は必須。query が指定されていない限り、読み取りの場合は必須。読み込み/書き込みに使用されます。有効な値: Azure SQL テーブルの名前、またはピリオドで区切られたスキーマ/テーブル名の組み合わせ。接続先のテーブルを識別するテーブルとスキーマを指定するために使用されます。デフォルトのスキーマは「public」です。テーブルがデフォルト以外のスキーマにある場合は、この情報を *schemaName.tableName* の形式で入力します。

- query — 読み取りに使用。Azure SQL から読み取るときに何を取得するかを定義する Transact-SQL SELECT クエリ。詳細については、[Microsoft のドキュメント](#)を参照してください。

## BigQuery 接続

AWS Glue for Spark を使用して、AWS Glue 4.0 以降のバージョンの Google BigQuery のテーブルからの読み取りとテーブルへの書き込みを行うことができます。Google SQL クエリを使用して BigQuery からの読み取りを行うことができます。AWS Glue 接続を介して AWS Secrets Manager で保存されている認証情報を使用して BigQuery に接続します。

Google BigQuery の詳細については、[Google Cloud BigQuery のウェブサイト](#)を参照してください。

## BigQuery 接続を設定する

AWS Glue から Google BigQuery へ接続するには、Google Cloud Platform の認証情報を作成して AWS Secrets Manager シークレットに保存し、そのシークレットを Google BigQuery AWS Glue 接続に関連付ける必要があります。

BigQuery への接続を設定するには:

1. Google Cloud Platform で、関連するリソースを作成して特定します。
  - 接続したい BigQuery テーブルを含む GCP プロジェクトを作成または特定します。
  - BigQuery API を有効にします。詳細については、「[Use the BigQuery Storage Read API to read table data](#)」を参照してください。
2. Google Cloud Platform で、サービスアカウントの認証情報を作成してエクスポートします。

BigQuery 認証情報ウィザードを使用すると、「[認証情報の作成](#)」というステップを迅速に実行できます。

GCP でサービスアカウントを作成するには、「[サービスアカウントを作成する](#)」にあるチュートリアルに従ってください。

- プロジェクトを選択するときは、BigQuery テーブルを含むプロジェクトを選択します。
- サービスアカウントの GCP IAM ロールを選択するときは、BigQuery テーブルの読み取り、書き込み、作成を行う BigQuery ジョブを実行するための適切な権限を付与するロールを追加または作成します。

サービスアカウントの認証情報を作成するには、「[サービス アカウント キーを作成する](#)」にあるチュートリアルに従ってください。

- キータイプを選択するときは、[JSON] を選択します。

これで、サービスアカウントの認証情報が記載された JSON ファイルがダウンロードされたはずです。これは次のように表示されます。

```
{
 "type": "service_account",
 "project_id": "*****",
 "private_key_id": "*****",
 "private_key": "*****",
 "client_email": "*****",
 "client_id": "*****",
 "auth_uri": "https://accounts.google.com/o/oauth2/auth",
 "token_uri": "https://oauth2.googleapis.com/token",
 "auth_provider_x509_cert_url": "https://www.googleapis.com/oauth2/v1/certs",
 "client_x509_cert_url": "*****",
 "universe_domain": "googleapis.com"
}
```

3. ダウンロードした認証情報ファイルを base64 でエンコードします。AWS CloudShell セッションなどでは、コマンドラインから `cat credentialsFile.json | base64 -w 0` コマンドを実行してこれを実行できます。このコマンドの出力 *credentialString* を保持してください。
4. AWS Secrets Manager で、Google Cloud Platform の認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
  - キーと値のペアを選択するときは、*credentialString* という値を持つキー `credentials` のペアを作成します。
5. AWS Glue データカタログで、「[the section called “AWS Glue 接続の追加”](#)」にある手順に従って接続を作成します。接続を作成したら、次のステップのために接続名 *connectionName* を保存しておきます。

- [接続タイプ] を選択するときは、Google BigQuery を選択してください。
  - [AWS Secret] をクリックして、*secretName* を入力します。
6. AWS Glue ジョブに関連付けられている IAM ロールに *secretName* を読み取るアクセス許可を付与します。
  7. AWS Glue ジョブ設定で、追加のネットワーク接続として *connectionName* を指定します。

## BigQuery テーブルから読み取る

### 前提条件:

- 読み取る対象の BigQuery テーブル。フォーム [dataset].[table] には BigQuery のテーブル名とデータセット名が必要です。これに *tableName* と名前を付けます。
- BigQuery テーブルの請求プロジェクト。そのためには、プロジェクトの名前、*parentProject* が必要です。請求元の親プロジェクトがない場合は、テーブルを含むプロジェクトを使用してください。
- BigQuery の認証情報。「AWS Glue で接続認証情報を管理するには」の手順を実行して、認証情報を設定します。AWS Glue 接続、*connectionName* の名前が必要になります。

### 例:

```
bigquery_read = glueContext.create_dynamic_frame.from_options(
 connection_type="bigquery",
 connection_options={
 "connectionName": "connectionName",
 "parentProject": "parentProject",
 "sourceType": "table",
 "table": "tableName",
 }
)
```

クエリを指定して、DynamicFrame に返される結果をフィルタリングすることもできます。query、sourceType、viewsEnabled および materializationDataset を設定する必要があります。

### 例:

### 追加の前提条件:

BigQuery がクエリのマテリアライズドビューを書き込むことのできる BigQuery データセット、*materializationDataset* を作成または特定する必要があります。

*materializationDataset* にテーブルを作成するには、サービスアカウントに適切な GCP IAM アクセス許可を付与する必要があります。

```
glueContext.create_dynamic_frame.from_options(
 connection_type="bigquery",
 connection_options={
 "connectionName": "connectionName",
 "materializationDataset": materializationDataset,
 "parentProject": "parentProject",
 "viewsEnabled": "true",
 "sourceType": "query",
 "query": "select * from bqtest.test"
 }
)
```

## BigQuery テーブルへ書き込む

この例では BigQuery サービスに直接書き込みます。BigQuery は「間接」書き込み方法もサポートしています。間接書き込みの設定の詳細については、「[the section called “Google BigQuery での間接書き込みを使用する”](#)」を参照してください。

### 前提条件:

- 書き込む対象の BigQuery テーブル。フォーム [dataset].[table] には BigQuery のテーブル名とデータセット名が必要です。自動的に作成される新しいテーブル名を指定することもできます。これに *tableName* と名前を付けます。
- BigQuery テーブルの請求プロジェクト。そのためには、プロジェクトの名前、*parentProject* が必要です。請求元の親プロジェクトがない場合は、テーブルを含むプロジェクトを使用してください。
- BigQuery の認証情報。「AWS Glue で接続認証情報を管理するには」の手順を実行して、認証情報を設定します。AWS Glue 接続、*connectionName* の名前が必要になります。

### 例:

```
bigquery_write = glueContext.write_dynamic_frame.from_options(
 frame=frameToWrite,
```

```
connection_type="bigquery",
connection_options={
 "connectionName": "connectionName",
 "parentProject": "parentProject",
 "writeMethod": "direct",
 "table": "tableName",
}
)
```

## BigQuery 接続オプションのリファレンス

- **project** — デフォルト: Google Cloud サービスアカウントのデフォルト。読み込み/書き込みに使用されます。テーブルに関連付けられている Google Cloud プロジェクトの名前。
- **table** — (必須) 読み込み/書き込みに使用されます。[[project:]dataset.] フォーマット内の BigQuery テーブルの名前。
- **dataset** — **table** オプションで定義されていない場合は必須。読み込み/書き込みに使用されます。BigQuery テーブルを含むデータセットの名前。
- **parentProject** — デフォルト: Google Cloud サービスアカウントのデフォルト。読み込み/書き込みに使用されます。請求に使用される **project** に関連付けられた Google Cloud プロジェクトの名前。
- **sourceType** — 読み取りに使用。読み取り時には必須です。有効値: **table**、**query** AWS Glue にテーブルで読み取るかクエリで読み取るかを通知します。
- **materializationDataset** — 読み取りに使用。有効な値: 文字列。ビューのマテリアライゼーションの保存に使用される BigQuery データセットの名前。
- **viewsEnabled** — 読み取りに使用。デフォルト: **false**。有効な値: **true**、**false**。BigQuery がビューを使用するかどうかを設定します。
- **query** — 読み取りに使用。 **viewsEnabled** が **true** のときに使用されます。GoogleSQL DQL クエリ。
- **temporaryGcsBucket** — 書き込みに使用。 **writeMethod** がデフォルト (**indirect**) に設定されている場合は必須です。BigQuery への書き込み中に中間形式のデータを保存するために使用される Google Cloud Storage バケツの名前。
- **writeMethod** - デフォルト: **indirect**。有効な値: **direct**、**indirect**。書き込みに使用。データの書き込みに使用する方法を指定します。
  - **direct** に設定すると、コネクタは BigQuery Storage Write API を使用して書き込みを行います。

- `indirect` に設定すると、コネクタは Google Cloud Storage に書き込み、読み取り操作を使用して BigQuery に転送します。Google Cloud サービスアカウントには適切な GCS アクセス許可が必要です。

## Google BigQuery での間接書き込みを使用する

この例では、Google Cloud ストレージにデータを書き込み、それを Google BigQuery にコピーする間接書き込みを使用しています。

前提条件:

Google Cloud Storage の一時的なバケット、*temporaryBucket* が必要になります。

AWS Glue の GCP サービスアカウントの GCP IAM ロールには、*temporaryBucket* にアクセスするための適切な GCS アクセス許可が必要です。

追加設定:

BigQuery による間接書き込みを設定するには:

1. [the section called “BigQuery を設定する”](#) にアクセスして検索するか、GCP 認証情報の JSON ファイルを再ダウンロードします。ジョブで使用する Google BigQuery AWS Glue 接続の AWS Secrets Manager シークレットである *secretName* を特定します。
2. 認証情報 JSON ファイルを、適切に安全な Amazon S3 の場所にアップロードします。ファイルへのパス、*s3secretpath* は、今後の手順に備えて保持してください。
3. *secretName* を編集し、`spark.hadoop.google.cloud.auth.service.account.json.keyfile` キーを追加します。*s3secretpath* に値を設定します。
4. AWS Glue ジョブに、*s3secretpath* にアクセスするための Amazon S3 IAM アクセス許可を付与します。

一時的な GCS バケットの場所を書き込みメソッドに指定できるようになりました。`indirect` がこれまでデフォルトだったように、`writeMethod` を指定する必要はありません。

```
bigquery_write = glueContext.write_dynamic_frame.from_options(
 frame=frameToWrite,
 connection_type="bigquery",
 connection_options={
 "connectionName": "connectionName",
```

```

 "parentProject": "parentProject",
 "temporaryGcsBucket": "temporaryBucket",
 "table": "tableName",
 }
)

```

## JDBC 接続

特定の (通常はリレーショナル) データベースタイプは、JDBC 標準による接続をサポートします。JDBC の詳細については、[Java JDBC API](#) に関するドキュメントを参照してください。AWS Glue は JDBC コネクタを介した特定のデータベースへの接続をネイティブにサポートしています。JDBC ライブラリは AWS Glue Spark ジョブで提供されます。AWS Glue ライブラリを使用してこれらのデータベースタイプに接続すると、標準のオプションセットにアクセスできます。

JDBC connectionType の値には次のようなものがあります。

- "connectionType": "sqlserver": Microsoft SQL Server データベースへの接続を指定します。
- "connectionType": "mysql": MySQL データベースへの接続を指定します。
- "connectionType": "oracle": Oracle データベースへの接続を指定します。
- "connectionType": "postgresql": PostgreSQL データベースへの接続を指定します。
- "connectionType": "redshift": Amazon Redshift データベースへの接続を指定します。詳細については、「[the section called “Redshift 接続”](#)」を参照してください。

次の表に、AWS Glue がサポートする JDBC ドライバーのバージョンを示します。

製品	Glue 4.0 の JDBC ドライバーのバージョン	Glue 3.0 の JDBC ドライバのバージョン	Glue 0.9、1.0、2.0 の JDBC ドライバのバージョン
Microsoft SQL Server	9.4.0	7.x	6.x
MySQL	8.0.23	8.0.23	5.1
Oracle Database	21.7	21.1	11.2
PostgreSQL	42.3.6	42.2.18	42.1.x

製品	Glue 4.0 の JDBC ドライバーのバージョン	Glue 3.0 の JDBC ドライバのバージョン	Glue 0.9、1.0、2.0 の JDBC ドライバのバージョン
MongoDB	4.7.2	4.0.0	2.0.0
Amazon Redshift *	redshift-jdbc42-2.1.0.16	redshift-jdbc41-1.2.12.1017	redshift-jdbc41-1.2.12.1017

\* Amazon Redshift 接続タイプの場合は、書式設定オプションを含め、JDBC 接続用に接続オプションに含まれるその他のオプション名/値ペアはすべて、基になる SparkSQL DataSource に直接渡されます。AWS Glue 4.0 以降のバージョンの AWS Glue with Spark ジョブでは、Amazon Redshift の AWS Glue ネイティブコネクタは Apache Spark の Amazon Redshift インテグレーションを使用します。詳細については、「[Apache Spark 用の Amazon Redshift の統合](#)」を参照してください。以前のバージョンでは、「[Spark 用の Amazon Redshift データソース](#)」を参照してください。

JDBC を使用して Amazon RDS データストアに接続するように Amazon VPC を設定するには、「[the section called “Amazon RDS データストアに接続するための VPC の設定”](#)」を参照してください。

#### Note

AWS Glue ジョブは、1 回の実行で 1 つのサブネットにのみ関連付けられます。これは、同じジョブで複数のデータソースに接続する能力に影響する可能性があります。この動作は JDBC ソースに限定されません。

#### トピック

- [JDBC 接続オプションのリファレンス](#)
- [sampleQuery を使用](#)
- [カスタム JDBC ドライバを使用](#)
- [JDBC テーブルからの並列読み取り](#)
- [AWS Glue から Amazon RDS データストアに JDBC 接続するための Amazon VPC の設定](#)

## JDBC 接続オプションのリファレンス

既に JDBC AWS Glue 接続を定義している場合、URL、ユーザー、パスワードなど、その中で定義されている設定プロパティを再利用できるため、コードで接続オプションとして指定する必要はありません。この機能は、AWS Glue バージョン 3.0 以降でのみ有効です。これを行うには、次の接続プロパティを使用します。

- "useConnectionProperties": この設定を接続から使用するときは、「true」に設定します。
- "connectionName": この設定を取得する接続名を入力します。接続は、ジョブと同じリージョンで定義されている必要があります。

JDBC 接続では、次の接続オプションを使用します。

- "url": (必須) データベースの JDBC URL。
- "dbtable": (必須) 読み込み元のデータベーステーブル。データベース内でスキーマをサポートする JDBC データストアの場合、`schema.table-name` を指定します。スキーマを指定しない場合、デフォルトの「パブリック」スキーマが使用されます。
- "user": (必須) 接続時に使用するユーザー名。
- "password": (必須) 接続時に使用するパスワード。
- (オプション) 次のオプションを使用すると、カスタム JDBC ドライバーを指定できます。AWS Glue がネイティブでサポートしていないドライバーを使用する必要がある場合は、これらのオプションを使用します。

ソースおよびターゲットが同じデータベース製品の場合でも、ETL ジョブは、データソースおよびターゲットに対して、異なるバージョンの JDBC ドライバーを使用することができます。これにより、異なるバージョンのソースデータベースとターゲットデータベース間でデータを移行できます。これらのオプションを使用するには、最初に JDBC ドライバーの JAR ファイルを Amazon S3 にアップロードする必要があります。

- "customJdbcDriverS3Path": カスタム JDBC ドライバーの Amazon S3 パス。
- "customJdbcDriverClassName": JDBC ドライバーのクラス名。
- "bulkSize": (オプション) JDBC ターゲットへのバルクロードを高速化するためのパラレル挿入を構成するために使用します。データの書き込みまたは挿入時に使用する並列度の整数値を指定します。このオプションは、Arch User Repository (AUR) などのデータベースへの書き込みのパフォーマンスを向上させるのに役立ちます。
- "hashfield" (オプション) JDBC テーブルから並列で読み込むときにデータをパーティションに分割するために使用される JDBC テーブル内の列の名前を指定するために使用される文字列。

「hashfield」または「hashexpression」を指定します。詳細については、「[the section called “JDBC からの並列読み取り”](#)」を参照してください。

- "hashexpression" (オプション) 整数を返す SQL SELECT 句。JDBC テーブルから並列で読み込むときに、JDBC テーブル内のデータをパーティションに分割するために使用されます。「hashfield」または「hashexpression」を指定します。詳細については、「[the section called “JDBC からの並列読み取り”](#)」を参照してください。
- "hashpartitions" (オプション) 正の整数。JDBC テーブルから並列で読み込むときに、JDBC テーブルの並列読み込み数を指定するために使用されます。デフォルト: 7。詳細については、「[the section called “JDBC からの並列読み取り”](#)」を参照してください。
- "sampleQuery": (オプション) カスタム SQL クエリステートメント。テーブル内の情報のサブセットを指定して、テーブルの内容のサンプルを取得するために使用されます。データを考慮せずに設定すると、DynamicFrame メソッドよりも効率が悪くなり、タイムアウトやメモリ不足エラーが発生する可能性があります。詳細については、「[the section called “sampleQuery を使用”](#)」を参照してください。
- "enablePartitioningForSampleQuery": (オプション) ブール値。デフォルト: false。sampleQuery の指定時に、JDBC テーブルから並列で読み込むことを可能にするために使用されます。true に設定すると、AWS Glue でパーティショニング条件を追加するためには、sampleQuery が「where」または「and」で終わる必要があります。詳細については、「[the section called “sampleQuery を使用”](#)」を参照してください。
- "sampleSize": (オプション) 正の整数。サンプルクエリによって返される行数を制限します。enablePartitioningForSampleQuery が true の場合にのみ動作します。パーティショニングが有効になっていない場合は、代わりに sampleQuery で "limit x" を追加してサイズを制限します。詳細については、「[the section called “sampleQuery を使用”](#)」を参照してください。

## sampleQuery を使用

このセクションでは、sampleQuery、sampleSize、enablePartitioningForSampleQuery を使用する方法について説明します。

sampleQuery ではデータセットのいくつかの行を効率的にサンプリングできます。デフォルトでは、クエリは単一のエグゼキューターによって実行されます。データを考慮せずに設定すると、DynamicFrame メソッドよりも効率が悪くなり、タイムアウトやメモリ不足エラーが発生する可能性があります。ETL パイプラインの一部として基盤となるデータベースで SQL を実行することは、通常、パフォーマンス上の目的でのみ必要になります。データセットの数行をプレビューする場合は、[the section called “show”](#) を使用することを検討してください。SQL を使用してデータセッ

トを変換する場合は、[the section called “toDF”](#) を使用して DataFrame フォーム内のデータに対して SparkSQL 変換を定義することを検討してください。

クエリではさまざまなテーブルを操作できますが、dbtable はやはり必須です。

sampleQuery を使用してテーブルのサンプルを取得する

デフォルトの sampleQuery の動作を使用してデータのサンプルを取得する場合、AWS Glue は大きなスループットを期待しないため、クエリは 1 つのエグゼキューターで実行されます。提供するデータを制限し、パフォーマンス上の問題を引き起こさないようにするには、SQL に LIMIT 句を指定することをお勧めします。

Example パーティション化せずに sampleQuery を使用する

次のコード例は、パーティション化せずに sampleQuery を使用する方法を示しています。

```
//A full sql query statement.
val query = "select name from $tableName where age > 0 limit 1"
val connectionOptions = JsonOptions(Map(
 "url" -> url,
 "dbtable" -> tableName,
 "user" -> user,
 "password" -> password,
 "sampleQuery" -> query))
val dyf = glueContext.getSource("mysql", connectionOptions)
 .getDynamicFrame()
```

sampleQuery を大きなデータセットに対して使用する

大きなデータセットを読み取る場合は、JDBC パーティション化を有効にして、テーブルを並列にクエリする必要がある場合があります。詳細については、「[the section called “JDBC からの並列読み取り”](#)」を参照してください。sampleQuery を JDBC パーティショニングで使用する場合は、enablePartitioningForSampleQuery を true に設定します。この機能を有効にするには、sampleQuery にいくつかの変更を加える必要があります。

sampleQuery で JDBC パーティショニングを使用するとき、パーティショニング条件を追加するためには、AWS Glue でクエリが「where」または「and」で終わる必要があります。

JDBC テーブルから並列で読み込むときに sampleQuery の結果を制限したい場合は、LIMIT 句を指定するのではなく、「sampleSize」パラメータを設定してください。

## Example JDBC パーティション化で sampleQuery を使用する

次のコード例は、JDBC パーティション化で sampleQuery を使用する方法を示しています。

```
//note that the query should end with "where" or "and" if use with JDBC partitioning.
val query = "select name from $tableName where age > 0 and"

//Enable JDBC partitioning by setting hashfield.
//to use sampleQuery with partitioning, set enablePartitioningForSampleQuery.
//use sampleSize to limit the size of returned data.
val connectionOptions = JsonOptions(Map(
 "url" -> url,
 "dbtable" -> tableName,
 "user" -> user,
 "password" -> password,
 "hashfield" -> primaryKey,
 "sampleQuery" -> query,
 "enablePartitioningForSampleQuery" -> true,
 "sampleSize" -> "1"))
val dyf = glueContext.getSource("mysql", connectionOptions)
 .getDynamicFrame()
```

### 注意と制限

サンプルクエリはジョブブックマークと一緒に使用することはできません。両方の設定を指定しても、ブックマークの状態は無視されます。

### カスタム JDBC ドライバーを使用

次のコード例は、カスタム JDBC ドライバーを使用して JDBC データベースから読み書きする方法を示しています。データベース製品の 1 つのバージョンから読み取り、同じ製品の新しいバージョンに書き込んでいます。

### Python

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext, SparkConf
from awsglue.context import GlueContext
from awsglue.job import Job
import time
from pyspark.sql.types import StructType, StructField, IntegerType, StringType
```

```
sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session

Construct JDBC connection options
connection_mysql5_options = {
 "url": "jdbc:mysql://<jdbc-host-name>:3306/db",
 "dbtable": "test",
 "user": "admin",
 "password": "pwd"}

connection_mysql8_options = {
 "url": "jdbc:mysql://<jdbc-host-name>:3306/db",
 "dbtable": "test",
 "user": "admin",
 "password": "pwd",
 "customJdbcDriverS3Path": "s3://DOC-EXAMPLE-BUCKET/mysql-connector-
java-8.0.17.jar",
 "customJdbcDriverClassName": "com.mysql.cj.jdbc.Driver"}

connection_oracle11_options = {
 "url": "jdbc:oracle:thin:@//<jdbc-host-name>:1521/ORCL",
 "dbtable": "test",
 "user": "admin",
 "password": "pwd"}

connection_oracle18_options = {
 "url": "jdbc:oracle:thin:@//<jdbc-host-name>:1521/ORCL",
 "dbtable": "test",
 "user": "admin",
 "password": "pwd",
 "customJdbcDriverS3Path": "s3://DOC-EXAMPLE-BUCKET/ojdbc10.jar",
 "customJdbcDriverClassName": "oracle.jdbc.OracleDriver"}

Read from JDBC databases with custom driver
df_mysql8 = glueContext.create_dynamic_frame.from_options(connection_type="mysql",
 connection_options=connection_mysql8_options)

Read DynamicFrame from MySQL 5 and write to MySQL 8
df_mysql5 = glueContext.create_dynamic_frame.from_options(connection_type="mysql",
 connection_options=connection_mysql5_options)
```

```
glueContext.write_from_options(frame_or_dfc=df_mysql5, connection_type="mysql",
 connection_options=connection_mysql8_options)

Read DynamicFrame from Oracle 11 and write to Oracle 18
df_oracle11 =
 glueContext.create_dynamic_frame.from_options(connection_type="oracle",

 connection_options=connection_oracle11_options)
glueContext.write_from_options(frame_or_dfc=df_oracle11, connection_type="oracle",
 connection_options=connection_oracle18_options)
```

## Scala

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.MappingSpec
import com.amazonaws.services.glue.errors.CallSite
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.DynamicFrame
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._

object GlueApp {
 val MYSQL_5_URI: String = "jdbc:mysql://<jdbc-host-name>:3306/db"
 val MYSQL_8_URI: String = "jdbc:mysql://<jdbc-host-name>:3306/db"
 val ORACLE_11_URI: String = "jdbc:oracle:thin:@//<jdbc-host-name>:1521/ORCL"
 val ORACLE_18_URI: String = "jdbc:oracle:thin:@//<jdbc-host-name>:1521/ORCL"

 // Construct JDBC connection options
 lazy val mysql5JsonOption = jsonOptions(MYSQL_5_URI)
 lazy val mysql8JsonOption = customJDBCdriverJsonOptions(MYSQL_8_URI, "s3://DOC-EXAMPLE-BUCKET/mysql-connector-java-8.0.17.jar", "com.mysql.cj.jdbc.Driver")
 lazy val oracle11JsonOption = jsonOptions(ORACLE_11_URI)
 lazy val oracle18JsonOption = customJDBCdriverJsonOptions(ORACLE_18_URI, "s3://DOC-EXAMPLE-BUCKET/ojdbc10.jar", "oracle.jdbc.OracleDriver")

 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
```

```
Job.init(args("JOB_NAME"), glueContext, args.asJava)

// Read from JDBC database with custom driver
val df_mysql8: DynamicFrame = glueContext.getSource("mysql",
mysql8JsonOption).getDynamicFrame()

// Read DynamicFrame from MySQL 5 and write to MySQL 8
val df_mysql5: DynamicFrame = glueContext.getSource("mysql",
mysql5JsonOption).getDynamicFrame()
glueContext.getSink("mysql", mysql8JsonOption).writeDynamicFrame(df_mysql5)

// Read DynamicFrame from Oracle 11 and write to Oracle 18
val df_oracle11: DynamicFrame = glueContext.getSource("oracle",
oracle11JsonOption).getDynamicFrame()
glueContext.getSink("oracle", oracle18JsonOption).writeDynamicFrame(df_oracle11)

Job.commit()
}

private def jsonOptions(url: String): JsonOptions = {
 new JsonOptions(
 s""""{"url": "${url}",
 |"dbtable": "test",
 |"user": "admin",
 |"password": "pwd"}"""".stripMargin)
}

private def customJDBCdriverJsonOptions(url: String, customJdbcDriverS3Path:
String, customJdbcDriverClassName: String): JsonOptions = {
 new JsonOptions(
 s""""{"url": "${url}",
 |"dbtable": "test",
 |"user": "admin",
 |"password": "pwd",
 |"customJdbcDriverS3Path": "${customJdbcDriverS3Path}",
 |"customJdbcDriverClassName" :
"${customJdbcDriverClassName}"""".stripMargin)
}
}
```

## JDBC テーブルからの並列読み取り

AWS Glue がデータを並列に読み込むように、JDBC テーブルのプロパティを設定できます。特定のプロパティを設定するときに、データの論理パーティションに対して、並列 SQL クエリを実行するように AWS Glue に指示します。ハッシュフィールド、またはハッシュ式を設定して、パーティション分割を制御できます。データにアクセスするために使用される、並列読み込みの数を制御することもできます。

JDBC テーブルからの並列読み取りは、パフォーマンスを向上させる最適化手法です。この手法が適切であることを特定するプロセスの詳細については、「AWS 規範的ガイダンス」の「Best practices for performance tuning AWS Glue for Apache Spark jobs」ガイドにある「[Reduce the amount of data scan](#)」を参照してください。

並列読み込みするに、テーブル構造のパラメータフィールドにキーと値のペアを設定できます。テーブルのパラメータフィールドに値を設定するには、JSON 表記を使用します。テーブルのプロパティを編集する詳しい方法については、「[テーブルの詳細の表示と編集](#)」を参照してください。ETL (抽出、変換、ロード) メソッド `create_dynamic_frame_from_options` および `create_dynamic_frame_from_catalog` を呼び出すときに、並列読み込みを有効にすることもできます。これらのメソッドのオプションの指定の詳細については、「[from\\_options](#)」および「[from\\_catalog](#)」を参照してください。

JDBC テーブル、つまり、ほとんどの基本データが JDBC データストアのテーブルで、このメソッドを使用できます。Amazon Redshift および Amazon S3 テーブルの読み取り時は、これらのプロパティは無視されます。

### hashfield

データをパーティションに分割するために使用する JDBC テーブルの列の名前を `hashfield` に設定します。最良の結果を得るために、この列の値は、パーティション間にデータを分散する、値の均等分散である必要があります。この列には任意のデータ型を使用できます。AWS Glue は、この列によってパーティション分割されたデータの並列読み込みを実行する、重複しないクエリを生成します。たとえば、データが毎月均等に分散される場合、`month` 列を使用して、各月のデータを並行して読み込むことができます。

```
'hashfield': 'month'
```

AWS Glue は、フィールド値をパーティション番号にハッシュするクエリを作成して、すべてのパーティションにクエリを並行して実行します。パーティションのテーブル読み込みに独自のクエリを使用するには、`hashfield` の代わりに `hashexpression` を指定します。

### hashexpression

`hashexpression` を、すべての番号を返す SQL 式 (JDBC データベースエンジンの文法に準拠する) に設定します。単純式は、テーブル内の任意の数値列の名前です。AWS Glue は、パーティションデータの WHERE 句で、`hashexpression` を使用して並行して JDBC データを読み込む SQL クエリを生成します。

たとえば、数値列 `customerID` を使用して、パーティション分割されたデータを顧客番号を使用して読み込みます。

```
'hashexpression': 'customerID'
```

AWS Glue がパーティション分割を制御するようにするには、`hashexpression` の代わりに `hashfield` を指定します。

### hashpartitions

JDBC テーブルの並列読み込みの数を `hashpartitions` に設定します。このプロパティが設定されていない場合、デフォルト値は、7 です。

たとえば、並列読み込みの数に 5 を設定すると、AWS Glue は 5 (以下) のクエリを使用して、データを読み込みます。

```
'hashpartitions': '5'
```

## AWS Glue から Amazon RDS データストアに JDBC 接続するための Amazon VPC の設定

JDBC を使用して Amazon RDS のデータベースに接続すると、追加の設定を行う必要があります。AWS Glue コンポーネントが Amazon RDS を通信できるようにするには、Amazon VPC で Amazon RDS データストアへのアクセスを設定する必要があります。AWS Glue がコンポーネント間で通信できるようにするには、すべての TCP ポートに対して自己参照のインバウンドルールを持つセキュリティグループを指定します。自己参照ルールを作成することにより、ソースを VPC の

同じセキュリティグループに制限することができます。自己参照ルールは、VPC をすべてのネットワークに開放しません。VPC のデフォルトのセキュリティグループには、すでに ALL Traffic (すべてのトラフィック) の自己参照インバウンドルールがある場合があります。

AWS Glue と Amazon RDS データストア間のアクセスをセットアップするには

1. AWS Management Console にサインインし、<https://console.aws.amazon.com/rds/> にある Amazon RDS コンソールを開きます。
2. Amazon RDS コンソールで、Amazon RDS データベースへのアクセスを管理するために使用するセキュリティグループを特定します。

左側のナビゲーションペインで [データベース] を選択し、メインペインのリストから接続するインスタンスを選択します。

データベースの詳細ページで、[接続とセキュリティ] タブで [VPC セキュリティグループ] を見つけます。

3. ネットワークアーキテクチャに基づいて、AWS Glue サービスへのアクセスを許可するために変更するのが最適な関連セキュリティグループを特定します。future 参照できるように、*database-security-group* 名前を保存してください。適切なセキュリティグループがない場合、Amazon RDS ドキュメントの「[セキュリティグループを作成して VPC 内の DB インスタンスへのアクセスを提供する](#)」の指示に従ってください。
4. AWS Management Console にサインインし、<https://console.aws.amazon.com/vpc/> にある Amazon VPC コンソールを開きます。
5. Amazon VPC コンソールで、*database-security-group* 更新する方法を特定します。

左側のナビゲーションペインで [セキュリティグループ] を選択し、*database-security-group* メインペインのリストから選択します。

6. のセキュリティグループ ID *database-sg-id* を特定します。*database-security-group* 今後の参照のために保存します。

セキュリティグループの詳細ページで、[セキュリティグループ ID] を探します。

7. *database-security-group* のインバウンドルールを変更し、AWS Glue コンポーネントが通信できるように自己参照ルールを追加します。具体的には、「タイプ」、「プロトコル」**All TCP**、「ポート範囲」にはすべてのポート **TCP**、「ソース」というルールを追加または確認します。*database-sg-id* [ソース] に入力したセキュリティグループが、編集中のセキュリティグループと同じであることを確認します。

セキュリティグループの詳細ページで、[インバウンドルールの編集] を選択します。

インバウンドルールは次のようになります。

タイプ	プロトコル	ポート範囲	ソース
すべての TCP	TCP	0 ~ 65535	<i>database-sg-id</i>

8. アウトバウンドトラフィックのルールを追加します。

セキュリティグループの詳細ページで、[アウトバウンドルールの編集] を選択します。

セキュリティグループがすべてのアウトバウンドトラフィックを許可する場合、個別のルールは必要ありません。例:

タイプ	プロトコル	ポート範囲	デスティネーション
すべてのトラフィック	すべて	すべて	0.0.0.0/0

ネットワークアーキテクチャがアウトバウンドトラフィックを制限するために設計されている場合、次のアウトバウンドルールを作成してください。

[タイプ] All TCP、[プロトコル]、[ポート範囲] がすべてのポート TCP、[宛先] の自己参照ルールを作成します。*database-sg-id*[送信先] に入力したセキュリティグループが、編集中のセキュリティグループと同じであることを確認します。

Amazon S3 VPC エンドポイントを使用している場合、VPC から Amazon S3 へのトラフィックを許可する HTTPS ルールを追加します。[タイプ]、[プロトコル] HTTPS、[ポート範囲] TCP、[宛先] が Amazon S3 ゲートウェイエンドポイント *s3 -* の管理対象プレフィックスリストの ID であるルールを作成します *prefix-list-id*。443プレフィックスのリストと Amazon S3 ゲートウェイエンドポイントの詳細については、Amazon VPC ドキュメントの「[Gateway endpoints for Amazon S3](#)」をご参照ください。

例:

タイプ	プロトコル	ポート範囲	デスティネーション
すべての TCP	TCP	0 ~ 65535	<i>database-sg-id</i>

タイプ	プロトコル	ポート範囲	デスティネーション
HTTPS	TCP	443	<i>s3-prefix-list-id</i>

## MongoDB 接続

AWS Glue for Spark を使用して、AWS Glue 4.0 以降のバージョンの MongoDB および MongoDB Atlas のテーブルからの読み取りとテーブルへの書き込みを行うことができます。AWS Glue 接続を介して AWS Secrets Manager で保存されているユーザー名およびパスワードの認証情報を使用して MongoDB に接続できます。

MongoDB の詳細については、[MongoDB のドキュメント](#)を参照してください。

## MongoDB 接続の設定

AWS Glue から MongoDB に接続するには、MongoDB 認証情報、*mongodbUser* および *mongodbPass* が必要です。

AWS Glue から MongoDB に接続するには、いくつかの前提条件を満たす必要がある場合があります。

- MongoDB インスタンスが Amazon VPC 内にある場合は、トラフィックがパブリックインターネットを経由することなく、AWS Glue ジョブが MongoDB インスタンスと通信できるように Amazon VPC を設定します。

Amazon VPC で、AWS Glue がジョブの実行中に使用する [VPC]、[サブネット]、および [セキュリティグループ] を特定または作成します。さらに、MongoDB インスタンスとこの場所の間のネットワークトラフィックを許可するように Amazon VPC が設定されているようにする必要があります。ネットワークレイアウトに基づいて、セキュリティグループルール、ネットワーク ACL、NAT ゲートウェイ、およびピアリング接続の変更が必要になる場合があります。

その後、MongoDB で使用できるように AWS Glue を設定する作業に進むことができます。

MongoDB に対する接続を設定するには:

1. 必要に応じて、AWS Secrets Manager で、MongoDB 認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください

い。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。

- [key/value ペア] を選択する際に、*mongodbUser* という値を持つキー `username` のペアを作成します。

[key/value ペア] を選択する際に、*mongodbPass* という値を持つキー `password` のペアを作成します。

2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、将来的に AWS Glue で使用するために、接続名 *connectionName* を維持します。

- [接続タイプ] を選択する際には、[MongoDB] または [MongoDB Atlas] を選択します。
- [MongoDB URL] または [MongoDB Atlas URL] を選択する場合は、MongoDB インスタンスのホスト名を入力します。

MongoDB URL は、`mongodb://mongoHost:mongoPort/mongoDBname` の形式で指定されます。

MongoDB Atlas URL は、`mongodb+srv://mongoHost:mongoPort/mongoDBname` の形式で指定されます。

接続にデフォルトのデータベースを指定する場合、*mongoDBname* はオプションです。

- Secrets Manager シークレットを作成することを選択した場合は、AWS Secrets Manager の [認証情報タイプ] を選択します。

その後、[AWS シークレット] で *secretName* を入力します。

- [ユーザー名とパスワード] を入力することを選択した場合は、*mongodbUser* および *mongodbPass* を入力します。

3. 次の状況では、追加の設定が必要になる場合があります。

- Amazon VPC の AWS でホストされている MongoDB インスタンスの場合
  - MongoDB セキュリティ認証情報を定義する AWS Glue 接続に、Amazon VPC 接続に関する情報を提供する必要があります。接続を作成または更新する際に、[ネットワークオプション] で [VPC]、[サブネット]、および [セキュリティグループ] を設定します。

AWS Glue MongoDB 接続を作成した後、接続メソッドを呼び出す前に次のアクションを実行する必要があります。

- Secrets Manager シークレットを作成することを選択した場合は、AWS Glue ジョブに関連付けられた IAM ロールに *secretName* を読み取るための許可を付与します。
- AWS Glue ジョブ設定で、追加のネットワーク接続として *connectionName* を指定します。

AWS Glue for Spark で AWS Glue MongoDB 接続を使用するには、接続メソッド呼び出しで *connectionName* オプションを指定します。あるいは、[the section called “MongoDB との統合”](#) のステップに従って、AWS Glue データカタログと組み合わせて接続を使用することもできます。

## AWS Glue 接続を使用した MongoDB からの読み取り

### 前提条件:

- 読み取り元とする MongoDB コレクション。コレクションの識別情報が必要になります。

MongoDB コレクションは、データベース名とコレクション名である *mongodbName* および *mongodbCollection* によって識別されます。

- 認証情報を提供するように設定された AWS Glue MongoDB 接続。認証情報を設定するには、前の手順「MongoDB に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。

### 例:

```
mongodb_read = glueContext.create_dynamic_frame.from_options(
 connection_type="mongodb",
 connection_options={
 "connectionName": "connectionName",
 "database": "mongodbName",
 "collection": "mongodbCollection",
 "partitioner":
 "com.mongodb.spark.sql.connector.read.partitionner.SinglePartitionPartitioner",
 "partitionerOptions.partitionSizeMB": "10",
 "partitionerOptions.partitionKey": "_id",
 "disableUpdateUri": "false",
 }
)
```

## MongoDB テーブルへの書き込み

この例では、既存の DynamicFrame である *dynamicFrame* から MongoDB に情報を書き込みます。

### 前提条件:

- 書き込み先とする MongoDB コレクション。コレクションの識別情報が必要になります。

MongoDB コレクションは、データベース名とコレクション名である *mongodbName* および *mongodbCollection* によって識別されます。

- 認証情報を提供するように設定された AWS Glue MongoDB 接続。認証情報を設定するには、前の手順「MongoDB に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。

### 例:

```
glueContext.write_dynamic_frame.from_options(
 frame=dynamicFrame,
 connection_type="mongodb",
 connection_options={
 "connectionName": "connectionName",
 "database": "mongodbName",
 "collection": "mongodbCollection",
 "disableUpdateUri": "false",
 "retryWrites": "false",
 },
)
```

### MongoDB テーブルに対する読み書き

この例では、既存の DynamicFrame である *dynamicFrame* から MongoDB に情報を書き込みます。

### 前提条件:

- 読み取り元とする MongoDB コレクション。コレクションの識別情報が必要になります。

書き込み先とする MongoDB コレクション。コレクションの識別情報が必要になります。

MongoDB コレクションは、データベース名とコレクション名である *mongodbName* および *mongodbCollection* によって識別されます。

- MongoDB 認証情報である *mongodbUser* および *mongodbPassword*。

### 例:

## Python

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext, SparkConf
from awsglue.context import GlueContext
from awsglue.job import Job
import time

@params: [JOB_NAME]
args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session

job = Job(glueContext)
job.init(args['JOB_NAME'], args)

output_path = "s3://some_bucket/output/" + str(time.time()) + "/"
mongo_uri = "mongodb://<mongo-instanced-ip-address>:27017"
mongo_ssl_uri = "mongodb://<mongo-instanced-ip-address>:27017"
write_uri = "mongodb://<mongo-instanced-ip-address>:27017"

read_mongo_options = {
 "uri": mongo_uri,
 "database": "mongodbName",
 "collection": "mongodbCollection",
 "username": "mongodbUsername",
 "password": "mongodbPassword",
 "partitioner": "MongoSamplePartitioner",
 "partitionerOptions.partitionSizeMB": "10",
 "partitionerOptions.partitionKey": "_id"}

ssl_mongo_options = {
 "uri": mongo_ssl_uri,
 "database": "mongodbName",
 "collection": "mongodbCollection",
 "ssl": "true",
 "ssl.domain_match": "false"
}

write_mongo_options = {
```

```
"uri": write_uri,
"database": "mongodbName",
"collection": "mongodbCollection",
"username": "mongodbUsername",
"password": "mongodbPassword",
}

Get DynamicFrame from MongoDB
dynamic_frame =
glueContext.create_dynamic_frame.from_options(connection_type="mongodb",

connection_options=read_mongo_options)

Write DynamicFrame to MongoDB
glueContext.write_dynamic_frame.from_options(dynamicFrame,
connection_type="mongodb", connection_options=write_mongo_options)

job.commit()
```

## Scala

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.MappingSpec
import com.amazonaws.services.glue.errors.CallSite
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.DynamicFrame
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._

object GlueApp {
 val DEFAULT_URI: String = "mongodb://<mongo-instanced-ip-address>:27017"
 val WRITE_URI: String = "mongodb://<mongo-instanced-ip-address>:27017"
 lazy val defaultJsonOption = jsonOptions(DEFAULT_URI)
 lazy val writeJsonOption = jsonOptions(WRITE_URI)
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)

 // Get DynamicFrame from MongoDB
```

```
val dynamicFrame: DynamicFrame = glueContext.getSource("mongodb",
defaultJsonOption).getDynamicFrame()

// Write DynamicFrame to MongoDB
glueContext.getSink("mongodb", writeJsonOption).writeDynamicFrame(dynamicFrame)

Job.commit()
}

private def jsonOptions(uri: String): JsonOptions = {
 new JsonOptions(
 s""""{"uri": "${uri}",
 |"database": "mongodbName",
 |"collection": "mongodbCollection",
 |"username": "mongodbUsername",
 |"password": "mongodbPassword",
 |"ssl": "true",
 |"ssl.domain_match": "false",
 |"partitioner": "MongoSamplePartitioner",
 |"partitionerOptions.partitionSizeMB": "10",
 |"partitionerOptions.partitionKey": "_id"}"""".stripMargin)
 }
}
```

## MongoDB 接続オプションのリファレンス

MongoDB への接続を指定します。接続オプションは、ソース接続とシンク接続とで異なります。

これらの接続プロパティは、ソース接続とシンク接続の間で共有されます。

- `connectionName` — 読み込み/書き込みに使用されます。認証およびネットワークの情報を接続方法に提供するように設定された AWS Glue MongoDB 接続の名前。前のセクション「[the section called “MongoDB の設定”](#)」で説明したように AWS Glue 接続が設定されている場合、`connectionName` を入力することで、`"uri"`、`"username"`、および `"password"` 接続オプションを入力する必要がなくなります。
- `"uri"`: (必須) 読み込み元の MongoDB ホスト (形式: `mongodb://<host>:<port>`)。AWS Glue 4.0 より前のバージョンの AWS Glue で使用されます。
- `"connection.uri"`: (必須) 読み込み元の MongoDB ホスト (形式: `mongodb://<host>:<port>`)。AWS Glue 4.0 以降のバージョンで使用されます。
- `"username"`: (必須) MongoDB のユーザー名。

- "password": (必須) MongoDB のパスワード。
- "database": (必須) 読み込み元の MongoDB データベース。このオプションは、ジョブスクリプトで `glue_context.create_dynamic_frame_from_catalog` を呼び出す際に、`additional_options` を介して渡すことも可能です。
- "collection": (必須) 読み込み元の MongoDB コレクション。このオプションは、ジョブスクリプトで `glue_context.create_dynamic_frame_from_catalog` を呼び出す際に、`additional_options` を介して渡すことも可能です。

"connectionType": "mongodb" ソースとする

"connectionType": "mongodb" をソースとして、次の接続オプションを使用します。

- "ssl": (オプション) true の場合、SSL 接続を開始します。デフォルト: false。
- "ssl.domain\_match": (任意) true と ssl が true の場合、ドメイン一致チェックが実行されます。デフォルト: true。
- "batchSize": (オプション): 内部バッチのカーソル内で使用される、バッチごとに返されるドキュメントの数。
- "partitioner": (オプション): MongoDB から入力データを読み取るためのパーティショナーのクラス名。コネクタには、次のパーティショナーがあります。
  - MongoDefaultPartitioner (デフォルト) (AWS Glue 4.0 ではサポートされていません)
  - MongoSamplePartitioner (MongoDB 3.2 以降が必要) (AWS Glue 4.0 ではサポートされていません)
  - MongoShardedPartitioner (AWS Glue 4.0 ではサポートされていません)
  - MongoSplitVectorPartitioner (AWS Glue 4.0 ではサポートされていません)
  - MongoPaginateByCountPartitioner (AWS Glue 4.0 ではサポートされていません)
  - MongoPaginateBySizePartitioner (AWS Glue 4.0 ではサポートされていません)
  - `com.mongodb.spark.sql.connector.read.partitionner.SinglePartitionPartitioner`
  - `com.mongodb.spark.sql.connector.read.partitionner.ShardedPartitioner`
  - `com.mongodb.spark.sql.connector.read.partitionner.PaginateIntoPartitionsPartitioner`
- "partitionerOptions" (オプション): 指定されたパーティショナーのオプション。各パーティショナーでは、次のオプションがサポートされています。
  - MongoSamplePartitioner: `partitionKey`, `partitionSizeMB`, `samplesPerPartition`
  - MongoShardedPartitioner: `shardkey`

- `MongoSplitVectorPartitioner`: `partitionKey`, `partitionSizeMB`
- `MongoPaginateByCountPartitioner`: `partitionKey`, `numberOfPartitions`
- `MongoPaginateBySizePartitioner`: `partitionKey`, `partitionSizeMB`

これらのオプションの詳細については、MongoDB のドキュメントの「[Partitioner Configuration](#)」を参照してください。

"connectionType": "mongodb" as Sink

"connectionType": "mongodb" をシンクとして、次の接続オプションを使用します。

- "ssl": (オプション) true の場合、SSL 接続を開始します。デフォルト: false。
- "ssl.domain\_match": (任意) true と ssl が true の場合、ドメイン一致チェックが実行されます。デフォルト: true。
- "extendedBsonTypes": (オプション) true が設定されている場合、MongoDB にデータを書き込む際に拡張 BSON 型を使用することを許可します。デフォルト: true。
- "replaceDocument": (オプション) true の場合、\_id フィールドを含むデータセットを保存するときに、ドキュメント全体を置き換えます。false の場合、データセットのフィールドと一致するドキュメントのフィールドのみが更新されます。デフォルト: true。
- "maxBatchSize": (オプション): データを保存するときの一括オペレーションの最大バッチサイズ。デフォルトは 512 です。
- "retryWrites": (オプション): AWS Glue でネットワークエラーが発生した場合、特定の書き込みオペレーションを 1 回自動的に再試行します。

## SAP HANA 接続

AWS Glue for Spark を使用して、AWS Glue 4.0 以降のバージョンの SAP HANA のテーブルからの読み取りとテーブルへの書き込みを行うことができます。SQL クエリを使用して、SAP HANA から何を読み取るかを定義できます。AWS Glue SAP HANA 接続を通じて AWS Secrets Manager に保存されている JDBC 認証情報を使用して SAP HANA に接続します。

SAP HANA JDBC の詳細については、[SAP HANA のドキュメント](#)を参照してください。

## SAP HANA 接続の設定

AWS Glue から SAP HANA に接続するには、SAP HANA 認証情報を作成して AWS Secrets Manager シークレットに保存し、そのシークレットを SAP HANA AWS Glue 接続に関連付ける必

必要があります。SAP HANA サービスと AWS Glue の間のネットワーク接続を設定する必要があります。

SAP HANA に接続するには、いくつかの前提条件を満たす必要がある場合があります。

- SAP HANA サービスが Amazon VPC 内にある場合は、トラフィックがパブリックインターネットを経由することなく、AWS Glue ジョブが SAP HANA サービスと通信できるように Amazon VPC を設定します。

Amazon VPC で、AWS Glue がジョブの実行中に使用する [VPC]、[サブネット]、および [セキュリティグループ] を特定または作成します。さらに、SAP HANA エンドポイントとこの場所の間のネットワークトラフィックを許可するように Amazon VPC が設定されているようにする必要があります。ジョブでは、SAP HANA JDBC ポートとの TCP 接続を確立する必要があります。SAP HANA ポートの詳細については、[SAP HANA のドキュメント](#)を参照してください。ネットワークレイアウトに基づいて、セキュリティグループルール、ネットワーク ACL、NAT ゲートウェイ、およびピアリング接続の変更が必要になる場合があります。

- SAP HANA エンドポイントがインターネットにアクセスできる場合、追加の前提条件はありません。

SAP HANA に対する接続を設定するには:

1. AWS Secrets Manager で、SAP HANA 認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
  - [key/value ペア] を選択する際に、*saphanaUsername* という値を持つキー *user* のペアを作成します。
  - [key/value ペア] を選択する際に、*saphanaPassword* という値を持つキー *password* のペアを作成します。
2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、将来的に AWS Glue で使用するために、接続名 *connectionName* を維持します。
  - [接続タイプ] を選択する際に、[SAP HANA] を選択します。
  - [SAP HANA URL] を入力する場合は、インスタンスの URL を入力します。

SAP HANA JDBC URL の形式は

```
jdbc:sap://saphanaHostname:saphanaPort/?databaseName=saphanaDBname,Parameter
です
```

AWS Glue には次の JDBC URL パラメータが必要です。

- `databaseName` – 接続先の SAP HANA のデフォルトデータベース。
- [AWS Secret] をクリックして、`secretName` を入力します。

AWS Glue SAP HANA 接続を作成した後、AWS Glue ジョブを実行する前に次のステップを実行する必要があります。

- AWS Glue ジョブに関連付けられている IAM ロールに `secretName` を読み取るアクセス許可を付与します。
- AWS Glue ジョブ設定で、追加のネットワーク接続として `connectionName` を指定します。

## SAP HANA テーブルからの読み取り

前提条件:

- 読み取り元とする SAP HANA テーブル。テーブルの識別情報が必要になります。

テーブルは、SAP HANA テーブル名とスキーマ名 (形式: `schemaName.tableName`) を使用して指定できます。テーブルがデフォルトのスキーマ「public」にある場合、スキーマ名と「.」区切り文字は必要ありません。この `tableIdentifier` を呼び出します。データベースは `connectionName` の JDBC URL パラメータとして指定されることに注意してください。

- 認証情報を提供するように設定された AWS Glue SAP HANA 接続。認証情報を設定するには、前の手順「SAP HANA に対する接続を設定するには」のステップを実行します。AWS Glue 接続、`connectionName` の名前が必要になります。

例:

```
saphana_read_table = glueContext.create_dynamic_frame.from_options(
 connection_type="saphana",
 connection_options={
 "connectionName": "connectionName",
 "dbtable": "tableIdentifier",
 }
)
```

```
)
```

SELECT SQL クエリを指定して、DynamicFrame に返される結果をフィルタリングすることもできます。query を設定する必要があります。

例:

```
saphana_read_query = glueContext.create_dynamic_frame.from_options(
 connection_type="saphana",
 connection_options={
 "connectionName": "connectionName",
 "query": "query"
 }
)
```

## SAP HANA テーブルへの書き込み

この例では、既存の DynamicFrame である *dynamicFrame* から SAP HANA に情報を書き込みます。テーブルに既に情報がある場合、AWS Glue はエラーになります。

前提条件:

- 書き込み先とする SAP HANA テーブル。

テーブルは、SAP HANA テーブル名とスキーマ名 (形式: *schemaName.tableName*) を使用して指定できます。テーブルがデフォルトのスキーマ「public」にある場合、スキーマ名と「.」区切り文字は必要ありません。この *tableIdentifier* を呼び出します。データベースは *connectionName* の JDBC URL パラメータとして指定されることに注意してください。

- SAP HANA 認証情報。認証情報を設定するには、前の手順「SAP HANA に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。

例:

```
options = {
 "connectionName": "connectionName",
 "dbtable": 'tableIdentifier'
}

saphana_write = glueContext.write_dynamic_frame.from_options(
 frame=dynamicFrame,
```

```
connection_type="saphana",
connection_options=options
)
```

## SAP HANA 接続オプションのリファレンス

- `connectionName` — 必須。読み込み/書き込みに使用されます。認証およびネットワークの情報を接続方法に提供するように設定された AWS Glue SAP HANA 接続の名前。
- `databaseName` — 読み込み/書き込みに使用されます。有効な値: SAP HANA のデータベースの名前。接続先のデータベースの名前。
- `dbtable` — 書き込みの場合は必須。query が指定されていない限り、読み取りの場合は必須。読み込み/書き込みに使用されます。有効な値: SAP HANA SQL FROM 句の内容。接続先の SAP HANA 内のテーブルを識別します。サブクエリなど、テーブル名以外の他の SQL を入力することもできます。詳細については、SAP HANA ドキュメントの「[From 句](#)」を参照してください。
- `query` — 読み取りに使用。SAP HANA から読み取るときに何を取得するかを定義する SAP HANA SQL SELECT クエリ。

## Snowflake 接続

AWS Glue for Spark を使用して、AWS Glue 4.0 以降のバージョンの Snowflake のテーブルからの読み込みとテーブルへの書き込みを行うことができます。Snowflake からの読み込みは SQL クエリで行うことができます。ユーザーとパスワードを使用して Snowflake に接続できます。AWS Glue データカタログを使用して AWS Secrets Manager に保存されている Snowflake 認証情報を参照できます。AWS Glue for Spark のデータカタログ Snowflake 認証情報は、クローラー用のデータカタログ Snowflake 認証情報とは別に保存されます。Snowflake に接続するように設定された JDBC タイプの接続ではなく、SNOWFLAKE タイプの接続を選択する必要があります。

Snowflake の詳細については、「[Snowflake ウェブサイト](#)」を参照してください。AWS 上の Snowflake の詳細については、「[Snowflake Data Warehouse on Amazon Web Services](#)」を参照してください。

## Snowflake 接続の設定

インターネット経由で利用できる Snowflake データベースに接続するための AWS の前提条件はありません。

オプションで、次の設定を実行して AWS Glue で接続認証情報を管理できます。

## AWS Glue で接続認証情報を管理するには

1. Snowflake でユーザー *snowflakeUser* とパスワード *snowflakePassword* を作成します。
2. AWS Secrets Manager で、Snowflake の認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
  - [キー/値のペア] をクリックして、sfUser キーを使用して *snowflakeUser* のペアを作成します。
  - [キー/値のペア] をクリックして、sfPassword キーを使用して *snowflakePassword* のペアを作成します。
  - キーと値のペアを選択するときに、Snowflake ウェアハウスにキー sfWarehouse を提供できます。
3. AWS Glue データカタログで、「[the section called “AWS Glue 接続の追加”](#)」にある手順に従って接続を作成します。接続を作成したら、次のステップのために接続名 *connectionName* を保存しておきます。
  - [接続タイプ] をクリックして、[Snowflake] を選択します。
  - [Snowflake URL] をクリックして、Snowflake インスタンスの URL を指定します。URL では、*account\_identifier*.snowflakecomputing.com という形式のホスト名を使用します。
  - [AWS Secret] をクリックして、*secretName* を入力します。
4. AWS Glue ジョブ設定で、追加のネットワーク接続として *connectionName* を指定します。

次のような状況では、以下が必要になることがあります。

- Amazon VPC の AWS でホストされている Snowflake の場合
  - Snowflake には適切な Amazon VPC 設定が必要です。Amazon VPC の設定方法の詳細については、Snowflake ドキュメントの「[AWS PrivateLink & Snowflake](#)」を参照してください。
  - AWS Glue には適切な Amazon VPC 設定が必要です。[the section called “VPC エンドポイント \(AWS PrivateLink\)”](#)。
  - (Snowflake のセキュリティ認証情報を定義する AWS Secrets Manager シークレットの ID に加えて) Amazon VPC 接続情報を提供する AWS Glue データカタログ接続を作成する必要がある

ります。前の項目でリンクされている Snowflake ドキュメントで説明されているように、AWS PrivateLink の使用時に URL が変更されます。

- ジョブ設定には、データカタログ接続を追加ネットワーク接続として含める必要があります。

## Snowflake テーブルからの読み込み

前提条件: 読み込む目的の Snowflake テーブルがあること。Snowflake のテーブル名 *tableName* が必要になります。Snowflake の URL *snowflakeUrl*、ユーザー名 *snowflakeUser*、パスワード *snowflakePassword* が必要です。Snowflake ユーザーにデフォルトの名前空間が設定されていない場合は、Snowflake データベース名 *databaseName*、スキーマ名 *schemaName* が必要になります。さらに、Snowflake ユーザーにデフォルトのウェアハウスセットがない場合は、ウェアハウス名 *warehouseName* が必要になります。

例:

その他の前提条件: 「AWS Glue で接続認証情報を管理するには」の手順を実行して、*snowflakeUrl*、*snowflakeUsername*、および *snowflakePassword* を設定します。これらのステップを確認するには、前のセクション「[the section called “Snowflake の設定”](#)」を参照してください。接続するその他のネットワーク接続を選択するには、*connectionName* パラメータを使用します。

```
snowflake_read = glueContext.create_dynamic_frame.from_options(
 connection_type="snowflake",
 connection_options={
 "connectionName": "connectionName",
 "dbtable": "tableName",
 "sfDatabase": "databaseName",
 "sfSchema": "schemaName",
 "sfWarehouse": "warehouseName",
 }
)
```

さらに、*autopushdown* および *query* パラメータを使用して Snowflake テーブルの一部を読み込むことができます。これは、Spark に読み込まれた後に結果をフィルタリングするよりもはるかに効率的です。すべての売上が同じテーブルに格納されているが、分析する必要があるのは休日の特定の店舗の売上だけだという例を考えてみましょう。その情報がテーブルに格納されている場合は、述語プッシュダウンを使用して次のように結果を取得できます。

```
snowflake_node = glueContext.create_dynamic_frame.from_options(

```

```
connection_type="snowflake",
connection_options={
 "autopushdown": "on",
 "query": "select * from sales where store='1' and IsHoliday='TRUE'",
 "connectionName": "snowflake-glue-conn",
 "sfDatabase": "databaseName",
 "sfSchema": "schemaName",
 "sfWarehouse": "warehouseName",
}
)
```

## Snowflake テーブルへの書き込み

前提条件: 書き込み先の Snowflake データベース。現在の、または希望するテーブル名 *tableName* が必要です。Snowflake の URL *snowflakeUrl*、ユーザー名 *snowflakeUser*、パスワード *snowflakePassword* が必要です。Snowflake ユーザーにデフォルトの名前空間が設定されていない場合は、Snowflake データベース名 *databaseName*、スキーマ名 *schemaName* が必要になります。さらに、Snowflake ユーザーにデフォルトのウェアハウスセットがない場合は、ウェアハウス名 *warehouseName* が必要になります。

例:

その他の前提条件: 「AWS Glue で接続認証情報を管理するには」の手順を実行して、*snowflakeUrl*、*snowflakeUsername*、および *snowflakePassword* を設定します。これらのステップを確認するには、前のセクション「[the section called “Snowflake の設定”](#)」を参照してください。接続するその他のネットワーク接続を選択するには、*connectionName* パラメータを使用します。

```
glueContext.write_dynamic_frame.from_options(
 connection_type="snowflake",
 connection_options={
 "connectionName": "connectionName",
 "dbtable": "tableName",
 "sfDatabase": "databaseName",
 "sfSchema": "schemaName",
 "sfWarehouse": "warehouseName",
 },
)
```

## Snowflake 接続オプションのリファレンス

Snowflake 接続タイプには次の接続オプションがあります。

このセクションのパラメータの一部はデータカタログ接続 (sfUrl、sfUser、sfPassword) から取得できます。その場合は指定する必要はありません。これを行うには、connectionName パラメータを含めます。

このセクションのパラメータの一部は、AWS Secrets Manager シークレット (sfUser、sfPassword) から取得できます。その場合は指定する必要はありません。シークレットは、sfUser および sfPassword キーの下にコンテンツを提供する必要があります。これを行うには、secretId パラメータを含めます。

以下のパラメータは、Snowflake に接続する際に一般的に使用されます。

- sfDatabase — Snowflake でユーザーデフォルトが設定されていない場合は必須です。読み込み/書き込みに使用されます。接続後にセッションに使用するデータベース。
- sfSchema — Snowflake でユーザーデフォルトが設定されていない場合は必須です。読み込み/書き込みに使用されます。接続後にセッションに使用するスキーマ。
- sfWarehouse — Snowflake でユーザーデフォルトが設定されていない場合は必須です。読み込み/書き込みに使用されます。接続後にセッションに使用するデフォルトの仮想ウェアハウス。
- sfRole — Snowflake でユーザーデフォルトが設定されていない場合は必須です。読み込み/書き込みに使用されます。接続後にセッションに使用するデフォルトのセキュリティロール。
- sfUrl — (必須) 読み込み/書き込みに使用されます。アカウントのホスト名を `account_identifier.snowflakecomputing.com` の形式で指定します。アカウント識別子の詳細については、Snowflake ドキュメントの「[アカウント識別子](#)」を参照してください。
- sfUser — (必須) 読み込み/書き込みに使用されます。Snowflake ユーザーのログイン名。
- sfPassword - (pem\_private\_key 指定がない場合は必須) 読み取り/書き込みに使用します。Snowflake ユーザーのパスワード。
- dbtable — フルテーブルで作業する場合は必須です。読み込み/書き込みに使用されます。読み込まれるテーブルまたはデータが書き込まれるテーブルの名前。読み込み時には、すべての列とレコードが取得されます。
- pem\_private\_key — 読み込み/書き込みに使用されます。暗号化されていない b64 でエンコードされたプライベートキーの文字列。Snowflake ユーザーのプライベートキー。これを PEM ファイルからコピーするのが一般的です。詳細については、Snowflake ドキュメントの「[キーペア認証とキーペアローテーション](#)」を参照してください。
- query — クエリで読み込みを行う場合は必須です。読み込みに使用されます。実行する正確なクエリ (SELECT ステートメント)

以下のオプションを使用して、Snowflake への接続プロセス中の特定の動作を設定します。

- `preactions` — 読み込み/書き込みに使用されます。有効値: セミコロンで区切られた文字列形式の SQL ステートメントのリスト。SQL ステートメントは、AWS Glue と Snowflake の間でデータが転送される前に実行されます。ステートメントに `%s` が含まれる場合、`%s` は操作で参照されるテーブル名に置き換えられます。
- `postactions` — 読み込み/書き込みに使用されます。SQL ステートメントは、AWS Glue と Snowflake 間でデータが転送された後に実行されます。ステートメントに `%s` が含まれる場合、`%s` は操作で参照されるテーブル名に置き換えられます。
- `autopushdown` - デフォルト: "on"。有効な値: "on"、"off"。このパラメータは、自動クエリプッシュダウンを有効にするかどうかを制御します。プッシュダウンが有効になっている場合、Spark でクエリを実行すると、クエリの一部が Snowflake サーバーに「プッシュダウン」できる場合にクエリがプッシュダウンされます。これにより、一部のクエリのパフォーマンスが向上します。クエリをプッシュダウンできるかどうかについては、Snowflake ドキュメントの「[プッシュダウン](#)」を参照してください。

さらに、Snowflake Spark コネクタで使用できるオプションの一部が AWS Glue でサポートされている場合があります。Snowflake Spark コネクタで使用できるオプションの詳細については、Snowflake ドキュメントの「[コネクタの構成オプションの設定](#)」を参照してください。

## Snowflake コネクタの制限事項

AWS Glue for Spark を使用してスノーflakeに接続することには、以下の制限があります。

- このコネクタはジョブブックマークをサポートしていません。ブックマークの詳細については、「[the section called “ジョブのブックマークを使用した処理済みデータの追跡”](#)」を参照してください。
- このコネクタは、`create_dynamic_frame.from_catalog` および `write_dynamic_frame.from_catalog` メソッドを使用した AWS Glue データカタログのテーブルを介した Snowflake の読み込みと書き込みをサポートしていません。
- このコネクタは、ユーザーとパスワード以外の認証情報による Snowflake への接続をサポートしていません。
- このコネクタはストリーミングジョブではサポートされていません。
- このコネクタは、情報を取得する (query パラメータを使用してなど) 際の SELECT ステートメントベースのクエリをサポートします。他の種類のクエリ (SHOW、DESC、DML ステートメントなど) はサポートされていません。

- Snowflake は、Snowflake クライアントを介して送信されるクエリテキスト (つまり、SQL ステートメント) のサイズをステートメントあたり 1 MB に制限しています。詳細については、「[クエリテキストサイズの制限](#)」を参照してください。

## Teradata Vantage 接続

AWS Glue for Spark を使用して、AWS Glue 4.0 以降のバージョンの Teradata Vantage の既存のテーブルからの読み取りとテーブルへの書き込みを行うことができます。SQL クエリを使用して、Teradata から何を読み取るかを定義できます。AWS Glue 接続を介して AWS Secrets Manager で保存されているユーザー名およびパスワードの認証情報を使用して Teradata に接続できます。

Teradata の詳細については、[Teradata ドキュメント](#)を参照してください。

## Teradata 接続の設定

AWS Glue から Teradata に接続するには、Teradata 認証情報を作成して AWS Secrets Manager シークレットに保存し、そのシークレットを Teradata AWS Glue 接続に関連付ける必要があります。Teradata インスタンスが Amazon VPC 内にある場合は、AWS Glue Teradata 接続にネットワークオプションを提供する必要もあります。

AWS Glue から Teradata に接続するには、いくつかの前提条件を満たす必要がある場合があります。

- Amazon VPC を通じて Teradata 環境にアクセスしている場合は、AWS Glue ジョブが Teradata 環境と通信できるように Amazon VPC を設定します。パブリックインターネット経由で Teradata 環境にアクセスすることは推奨されていません。

Amazon VPC で、AWS Glue がジョブの実行中に使用する [VPC]、[サブネット]、および [セキュリティグループ] を特定または作成します。さらに、Teradata インスタンスとこの場所の間のネットワークトラフィックを許可するように Amazon VPC が設定されているようにする必要があります。ジョブでは、Teradata クライアントポートとの TCP 接続を確立する必要があります。Teradata ポートの詳細については、[Teradata のドキュメント](#)を参照してください。

ネットワークレイアウトに応じて、安全な VPC 接続を実現するには、Amazon VPC および他のネットワークサービスの変更が必要な場合があります。AWS 接続の詳細については、Teradata ドキュメントの「[AWS 接続オプション](#)」を参照してください。

AWS Glue Teradata 接続を設定するには:

1. Teradata 設定で、AWS Glue が接続するユーザーとパスワード (*teradataUser* および *teradataPassword*) を識別または作成します。詳細については、Teradata ドキュメントの「[Vantage セキュリティの概要](#)」を参照してください。
2. AWS Secrets Manager で、Teradata 認証情報を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。
  - [key/value ペア] を選択する際に、*teradataUsername* という値を持つキー *user* のペアを作成します。
  - [key/value ペア] を選択する際に、*teradataPassword* という値を持つキー *password* のペアを作成します。
3. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、次のステップのために接続名 *connectionName* を保存しておきます。
  - [接続タイプ] を選択する際に、[Teradata] を選択します。
  - [JDBC URL] を入力する場合は、インスタンスの URL を入力します。JDBC URL に特定のカンマ区切りの接続パラメータをハードコーディングすることもできます。URL は次の形式に準拠する必要があります:  
`jdbc:teradata://teradataHostname/ParameterName=ParameterValue,ParameterName`

サポートされる URL パラメータには次が含まれます。

  - DATABASE – デフォルトでアクセスするホスト上のデータベースの名前。
  - DBS\_PORT – データベースポート。非標準ポートで実行する場合に使用されます。
  - [認証情報タイプ] を選択する場合は、[AWS Secrets Manager] を選択し、[AWS シークレット] を *secretName* に設定します。
4. 次の状況では、追加の設定が必要になる場合があります。
  - Amazon VPC の AWS でホストされている Teradata インスタンスの場合
    - Teradata セキュリティ認証情報を定義する AWS Glue 接続に、Amazon VPC 接続に関する情報を提供する必要があります。接続を作成または更新する際に、[ネットワークオプション] で [VPC]、[サブネット]、および [セキュリティグループ] を設定します。

AWS Glue Teradata 接続を作成した後、接続メソッドを呼び出す前に次のステップを実行する必要があります。

- AWS Glue ジョブに関連付けられている IAM ロールに *secretName* を読み取るアクセス許可を付与します。
- AWS Glue ジョブ設定で、追加のネットワーク接続として *connectionName* を指定します。

## Teradata からの読み取り

前提条件:

- 読み取り元とする Teradata テーブル。テーブル名 *tableName* が必要になります。
- 認証情報を提供するように設定された AWS Glue Teradata 接続。認証情報を設定するには、「Teradata に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。

例:

```
teradata_read_table = glueContext.create_dynamic_frame.from_options(
 connection_type="teradata",
 connection_options={
 "connectionName": "connectionName",
 "dbtable": "tableName"
 }
)
```

SELECT SQL クエリを指定して、DynamicFrame に返される結果をフィルタリングすることもできます。query を設定する必要があります。

例:

```
teradata_read_query = glueContext.create_dynamic_frame.from_options(
 connection_type="teradata",
 connection_options={
 "connectionName": "connectionName",
 "query": "query"
 }
)
```

## Teradata テーブルへの書き込み

前提条件: 書き込み先とする Teradata テーブルである *tableName*。接続メソッドを呼び出す前にテーブルを作成する必要があります。

例:

```
teradata_write = glueContext.write_dynamic_frame.from_options(
 connection_type="teradata",
 connection_options={
 "connectionName": "connectionName",
 "dbtable": "tableName"
 }
)
```

## Teradata 接続オプションのリファレンス

- *connectionName* — 必須。読み込み/書き込みに使用されます。認証およびネットワークの情報を接続方法に提供するように設定された AWS Glue Teradata 接続の名前。
- *dbtable* — 書き込みの場合は必須。query が指定されていない限り、読み取りの場合は必須。読み込み/書き込みに使用されます。接続メソッドがインタラクシオンするテーブルの名前。
- *query* — 読み取りに使用。Teradata から読み取るときに何を取得するかを定義する SELECT SQL クエリ。

## Vertica 接続

AWS Glue for Spark を使用して、AWS Glue 4.0 以降のバージョンの Vertica のテーブルからの読み取りとテーブルへの書き込みを行うことができます。SQL クエリを使用して、Vertica から何を読み取るかを定義できます。AWS Glue 接続を介して AWS Secrets Manager で保存されているユーザー名およびパスワードの認証情報を使用して Vertica に接続します。

Vertica の詳細については、[Vertica ドキュメント](#)を参照してください。

## Vertica 接続の設定

AWS Glue から Vertica に接続するには、Vertica 認証情報を作成して AWS Secrets Manager シークレットに保存し、そのシークレットを Vertica AWS Glue 接続に関連付ける必要があります。Vertica インスタンスが Amazon VPC 内にある場合は、AWS Glue Vertica 接続にネットワークオプションを提供する必要があります。データベースの読み取りおよび書き込み時に一時ストレージとして使用する Amazon S3 バケットまたはフォルダが必要です。

AWS Glue から Vertica に接続するには、いくつかの前提条件を満たす必要があります。

- データベースの読み取りおよび書き込み時に一時ストレージとして使用する Amazon S3 バケットまたはフォルダ。 *tempS3Path* によって参照されます。

#### Note

AWS Glue ジョブデータプレビューで Vertica を使用する場合、一時ファイルは *tempS3Path* から自動的に削除されない場合があります。一時ファイルを確実に削除するには、[データプレビュー] ペインで [セッションの終了] を選択して、データプレビューセッションを直接終了します。

データプレビューセッションが直接終了することを保証できない場合は、Amazon S3 ライフサイクル構成を設定して古いデータを削除することを検討してください。ジョブの最大実行時間にマージンを加えた値に基づいて、49 時間を超える時間が経過しているデータを削除することをお勧めします。Amazon S3 ライフサイクルの設定の詳細については、Amazon S3 ドキュメントの「[ストレージのライフサイクルの管理](#)」を参照してください。

- AWS Glue ジョブロールに関連付けることができる、Amazon S3 パスに対する適切な許可を持つ IAM ポリシー。
- Vertica インスタンスが Amazon VPC 内にある場合は、トラフィックがパブリックインターネットを経由することなく、AWS Glue ジョブが Vertica インスタンスと通信できるように Amazon VPC を設定します。

Amazon VPC で、AWS Glue がジョブの実行中に使用する [VPC]、[サブネット]、および [セキュリティグループ] を特定または作成します。さらに、Vertica インスタンスとこの場所の間のネットワークトラフィックを許可するように Amazon VPC が設定されているようにする必要があります。ジョブは、Vertica クライアントポート (デフォルトは 5433) との TCP 接続を確立する必要があります。ネットワークレイアウトに基づいて、セキュリティグループルール、ネットワーク ACL、NAT ゲートウェイ、およびピアリング接続の変更が必要になる場合があります。

その後、Vertica で使用できるように AWS Glue を設定する作業に進むことができます。

Vertica に対する接続を設定するには:

1. AWS Secrets Manager で、Vertica 認証情報、*verticaUsername* および *verticaPassword* を使用してシークレットを作成します。Secrets Manager でシークレットを作成するには、AWS Secrets Manager ドキュメントの「[AWS Secrets Manager シークレットを作成する](#)」

にあるチュートリアルに従ってください。シークレットを作成したら、次のステップのためにシークレット名 *secretName* を保存しておきます。

- [key/value ペア] を選択する際に、*verticaUsername* という値を持つキー user のペアを作成します。
  - [key/value ペア] を選択する際に、*verticaPassword* という値を持つキー password のペアを作成します。
2. AWS Glue コンソールで、「[the section called “AWS Glue 接続の追加”](#)」にあるステップに従って接続を作成します。接続を作成したら、次のステップのために接続名 *connectionName* を保存しておきます。
- [接続タイプ] を選択する際に、[Vertica] を選択します。
  - [Vertica ホスト] を選択する場合は、Vertica インストールのホスト名を入力します。
  - [Vertica ポート] を選択すると、Vertica インストールを使用できるポートが選択されます。
  - [AWS Secret] をクリックして、*secretName* を入力します。
3. 次の状況では、追加の設定が必要になる場合があります。
- Amazon VPC の AWS でホストされている Vertica インスタンスの場合
    - Vertica セキュリティ認証情報を定義する AWS Glue 接続に、Amazon VPC 接続に関する情報を提供します。接続を作成または更新する際に、[ネットワークオプション] で [VPC]、[サブネット]、および [セキュリティグループ] を設定します。

AWS Glue Vertica 接続を作成した後、接続メソッドを呼び出す前に次のステップを実行する必要があります。

- AWS Glue ジョブの許可に関連付けられた IAM ロールを *tempS3Path* に付与します。
- AWS Glue ジョブに関連付けられている IAM ロールに *secretName* を読み取るアクセス許可を付与します。
- AWS Glue ジョブ設定で、追加のネットワーク接続として *connectionName* を指定します。

## Vertica からの読み取り

### 前提条件:

- 読み取り元とする Vertica テーブル。Vertica データベース名である *dbName* とテーブル名である *tableName* が必要になります。

- 認証情報を提供するように設定された AWS Glue Vertica 接続。認証情報を設定するには、前の手順「Vertica に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。
- 前述した一時ストレージに使用する Amazon S3 バケットまたはフォルダ。*tempS3Path* という名前が必要になります。s3a プロトコルを使用してこの場所に接続する必要があります。

例:

```
dynamicFrame = glueContext.create_dynamic_frame.from_options(
 connection_type="vertica",
 connection_options={
 "connectionName": "connectionName",
 "staging_fs_url": "s3a://tempS3Path",
 "db": "dbName",
 "table": "tableName",
 }
)
```

SELECT SQL クエリを入力して、DynamicFrame に返された結果をフィルタリングしたり、複数のテーブルからデータセットにアクセスしたりすることもできます。

例:

```
dynamicFrame = glueContext.create_dynamic_frame.from_options(
 connection_type="vertica",
 connection_options={
 "connectionName": "connectionName",
 "staging_fs_url": "s3a://tempS3Path",
 "db": "dbName",
 "query": "select * FROM tableName",
 },
)
```

## Vertica テーブルへの書き込み

この例では、既存の DynamicFrame である *dynamicFrame* から Vertica に情報を書き込みます。テーブルに既に情報がある場合、AWS Glue は DynamicFrame からのデータを付加します。

前提条件:

- 書き込み先とする現在または希望のテーブル名である *tableName*。対応する Vertica データベース名である *dbName* も必要になります。
- 認証情報を提供するように設定された AWS Glue Vertica 接続。認証情報を設定するには、前の手順「Vertica に対する接続を設定するには」のステップを実行します。AWS Glue 接続、*connectionName* の名前が必要になります。
- 前述した一時ストレージに使用する Amazon S3 バケットまたはフォルダ。*tempS3Path* という名前が必要になります。s3a プロトコルを使用してこの場所に接続する必要があります。

例:

```
glueContext.write_dynamic_frame.from_options(
 frame=dynamicFrame,
 connection_type="vertica",
 connection_options={
 "connectionName": "connectionName",
 "staging_fs_url": "s3a://tempS3Path",
 "db": "dbName",
 "table": "tableName",
 }
)
```

## Vertica 接続オプションのリファレンス

- *connectionName* — 必須。読み込み/書き込みに使用されます。認証およびネットワークの情報を接続方法に提供するように設定された AWS Glue Vertica 接続の名前。
- *db* — 必須。読み込み/書き込みに使用されます。接続方法がインタラクションする Vertica 内のデータベースの名前。
- *dbSchema* — テーブルを識別する必要がある場合は必須です。読み込み/書き込みに使用されます。デフォルト: *public*。接続メソッドがインタラクションするスキーマの名前。
- *table* — 書き込みの場合は必須。query が指定されていない限り、読み取りの場合は必須。読み込み/書き込みに使用されます。接続メソッドがインタラクションするテーブルの名前。
- *query* — 読み取りに使用。Teradata から読み取るときに何を取得するかを定義する SELECT SQL クエリ。
- *staging\_fs\_url* — 必須。読み込み/書き込みに使用されます。有効な値: s3a URL。一時ストレージに使用する Amazon S3 バケットまたはフォルダの URL。

## カスタムと AWS Marketplace での、connectionType の値

これには以下が含まれます。

- "connectionType": "marketplace.athena": Amazon Athena データストアへの接続を指定します。この接続では、AWS Marketplace が提供するコネクタを使用します。
- "connectionType": "marketplace.spark": Apache Spark データストアへの接続を指定します。この接続では、AWS Marketplace が提供するコネクタを使用します。
- "connectionType": "marketplace.jdbc": JDBC データストアへの接続を指定します。この接続では、AWS Marketplace が提供するコネクタを使用します。
- "connectionType": "custom.athena": Amazon Athena データストアへの接続を指定します。接続には、AWS Glue Studio にアップロードしたカスタムコネクタを使用します。
- "connectionType": "custom.spark": Apache Spark データストアへの接続を指定します。接続には、AWS Glue Studio にアップロードしたカスタムコネクタを使用します。
- "connectionType": "custom.jdbc": JDBC データストアへの接続を指定します。接続には、AWS Glue Studio にアップロードしたカスタムコネクタを使用します。

### custom.jdbc または marketplace.jdbc 型での接続オプション

- className – (必須) ドライバクラス名を示す文字列。
- connectionName – (必須) コネクタに関連付けられている接続の名前を示す文字列。
- url – (必須) データソースへの接続を構築するために使用される、プレースホルダを含む JDBC URL (`{}`) を示す文字列。プレースホルダー `{secretKey}` は、AWS Secrets Manager 内にある同じ名前のシークレットにより置き換えられます。URL の構築の詳細については、データストアのドキュメントを参照してください。
- secretId または user/password – (必須) URL の認証情報を取得するために使用される文字列。
- dbTable または query – (必須) データを取得するテーブルまたは SQL クエリを示す文字列。dbTable または query を指定できます。両方を指定することはできません。
- partitionColumn – (オプション) パーティション化に使用される整数カラムの名前を示す文字列。このオプションは、lowerBound、upperBound、および numPartitions に含まれている場合にのみ機能します。このオプションの機能は、Spark SQL JDBC リーダーのものと同様です。詳細については、Apache Spark SQL, DataFrames and Datasets Guide の「[JDBC To Other Databases](#)」を参照してください。

lowerBound および upperBound 値は、パーティションのストライドを決定するために使用されます (テーブル内の行のフィルタリングには使用しません)。返されるテーブル内のすべての行は、パーティション化されています。

#### Note

テーブル名の代わりにクエリを使用する場合は、指定されたパーティショニング条件でクエリが動作することを確認する必要があります。例:

- "SELECT col1 FROM table1" の形式のクエリでパーティション列を使用する場合、末尾に WHERE 句を追加してそのクエリをテストします。
- クエリ形式が SELECT col1 FROM table1 WHERE col2=val" の場合は、WHERE 句を AND とパーティション列を使用する式で拡張することで、そのクエリをテストします。

- lowerBound – パーティションストライドを決定するために使用される partitionColumn の最小値を示す整数 (オプション)。
- upperBound – パーティションストライドを決定するために使用される partitionColumn の最大値を示す整数 (オプション)。
- numPartitions – パーティション数を示す整数 (オプション)。この値は、(範囲に含まれる) lowerBound と (範囲に含まれない) upperBound とともに使用され、partitionColumn の分割で使用するために生成された WHERE 句の式のための、パーティションストライドを形成します。

#### Important

パーティションが多すぎると、外部データベースシステムで問題が発生する可能性があるため、パーティションの数には注意を払ってください。

- filterPredicate – ソースからのデータをフィルタリングする、追加の条件句を示す文字列 (オプション)。例:

```
BillingCity='Mountain View'
```

table 名の代わりに query を使用した場合は、指定された filterPredicate でクエリが動作することを確認します。例:

- クエリの形式が "SELECT col1 FROM table1" の場合は、フィルタ述語を使用するクエリの末尾に WHERE 句を追加して、そのクエリをテストします。
- クエリ形式が "SELECT col1 FROM table1 WHERE col2=val" の場合は、WHERE 句を AND およびフィルター述語を使用する式で拡張して、そのクエリをテストします。
- dataTypeMapping – (デイクシヨナリ、オプション) JDBC データ型 から Glue データ型 に対するマッピングを構築する、カスタムのデータ型マッピング。例えば、オプション "dataTypeMapping":{"FLOAT":"STRING"} では、ドライバーの ResultSet.getString() メソッドを呼び出すことで、JDBC タイプ FLOAT のデータフィールドが Java String 型にマッピングされ、それを使用して AWS Glue レコードが構築されます。ResultSet オブジェクトは各ドライバによって実装されるため、その動作は使用するドライバにより決定されます。ドライバによる変換の実行方法については、JDBC ドライバのドキュメントを参照してください。
- AWS Glue で現在サポートされているデータ型は次のとおりです。
  - DATE
  - STRING
  - TIMESTAMP
  - INT
  - FLOAT
  - LONG
  - BIGDECIMAL
  - BYTE
  - SHORT
  - DOUBLE

JDBC データ型としては、[Java8 java.sql.types](#) がサポートされています。

デフォルトの (JDBC から AWS Glue への) データ型マッピングは以下のとおりです。

- DATE -> DATE
- VARCHAR -> STRING
- CHAR -> STRING
- LONGNVARCHAR -> STRING
- TIMESTAMP -> TIMESTAMP
- INTEGER -> INT
- FLOAT -> FLOAT

- REAL -> FLOAT
- BIT -> BOOLEAN
- BOOLEAN -> BOOLEAN
- BIGINT -> LONG
- DECIMAL -> BIGDECIMAL
- NUMERIC -> BIGDECIMAL
- TINYINT -> SHORT
- SMALLINT -> SHORT
- DOUBLE -> DOUBLE

カスタムのデータタイプマッピングで `dataTypeMapping` オプションを使用すると、デフォルトのデータ型マッピングをオーバーライドできます。この影響を受けるのは、`dataTypeMapping` オプションでリストされた JDBC データ型のみです。他のすべての JDBC データ型に対しては、デフォルトのマッピングが使用されます。必要に応じて、別の JDBC データ型のマッピングを追加することも可能です。デフォルトまたはカスタムのマッピングのいずれにも JDBC データ型が含まれていない場合、データ型はデフォルトで AWS GlueSTRING データ型に変換されます。

次の Python コード例は、AWS Marketplace の JDBC ドライバーを使用して、JDBC データベースからの読み取りを実行する方法を示しています。ここでは、データベースからの読み取りと、S3 口ケーションへの書き込みの方法を知ることができます。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job

@params: [JOB_NAME]
args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'], args)
@type: DataSource
@args: [connection_type = "marketplace.jdbc", connection_options =
```

```

 {"dataTypeMapping":{"INTEGER":"STRING"},"upperBound":"200","query":"select id,
 name, department from department where id < 200","numPartitions":"4",
 "partitionColumn":"id","lowerBound":"0","connectionName":"test-connection-
jdbc"},
 transformation_ctx = "DataSource0"]
@return: DataSource0
@inputs: []
DataSource0 = glueContext.create_dynamic_frame.from_options(connection_type =
 "marketplace.jdbc", connection_options = {"dataTypeMapping":{"INTEGER":"STRING"},
 "upperBound":"200","query":"select id, name, department from department where
 id < 200","numPartitions":"4","partitionColumn":"id","lowerBound":"0",
 "connectionName":"test-connection-jdbc"}, transformation_ctx = "DataSource0")
@type: ApplyMapping
@args: [mappings = [("department", "string", "department", "string"), ("name",
"string",
 "name", "string"), ("id", "int", "id", "int")], transformation_ctx =
"Transform0"]
@return: Transform0
@inputs: [frame = DataSource0]
Transform0 = ApplyMapping.apply(frame = DataSource0, mappings = [("department",
"string",
 "department", "string"), ("name", "string", "name", "string"), ("id", "int",
"id", "int")],
 transformation_ctx = "Transform0")
@type: DataSink
@args: [connection_type = "s3", format = "json", connection_options = {"path":
"s3://<S3 path>/", "partitionKeys": []}, transformation_ctx = "DataSink0"]
@return: DataSink0
@inputs: [frame = Transform0]
DataSink0 = glueContext.write_dynamic_frame.from_options(frame = Transform0,
 connection_type = "s3", format = "json", connection_options = {"path":
 "s3://<S3 path>/", "partitionKeys": []}, transformation_ctx = "DataSink0")
job.commit()

```

## custom.athena または marketplace.athena タイプ用の接続オプション

- `className` – (必須) ドライバクラス名を示す文字列。Athena-CloudWatch コネクタを使用している場合、このパラメータ値はクラス名にプレフィックスされます (例: "com.amazonaws.athena.connectors")。Athena-CloudWatch コネクタは、メタデータハンドラーとレコードハンドラーの 2 つのクラスで構成されています。共通のプレフィックスを指定することで、API がそのプレフィックスに基づいた適切なクラスをロードします。

- `tableName` – (必須) 読み込む CloudWatch ログストリームの名前を示す文字列。このコードスニペットでは、ビューに特別な名前 `all_log_streams` を使用しています。この場合、返された動的データフレームには、ロググループ内のすべてのログストリームからのデータが含まれます。
- `schemaName` – (必須) 読み取りのソースとなる CloudWatch ロググループの名前を示す文字列。例えば、`/aws-glue/jobs/output` と指定します。
- `connectionName` – (必須) コネクタに関連付けられている接続の名前を示す文字列。

このコネクタの追加オプションについては、GitHub の [Amazon Athena CloudWatch Connector README](#) ファイルを参照してください。

次の Python コード例は、AWS Marketplace コネクタを使用しながら、Athena データストアからの読み取りを実行する方法を示しています。ここでは、Athena からの読み取りと、S3 ロケーションへの書き込みを行う方法を知ることができます。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job

@params: [JOB_NAME]
args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'], args)
@type: DataSource
@args: [connection_type = "marketplace.athena", connection_options =
 {"tableName": "all_log_streams", "schemaName": "/aws-glue/jobs/output",
 "connectionName": "test-connection-athena"}, transformation_ctx = "DataSource0"]
@return: DataSource0
@inputs: []
DataSource0 = glueContext.create_dynamic_frame.from_options(connection_type =
 "marketplace.athena", connection_options = {"tableName": "all_log_streams",,
 "schemaName": "/aws-glue/jobs/output", "connectionName":
 "test-connection-athena"}, transformation_ctx = "DataSource0")
@type: ApplyMapping
```

```

 ## @args: [mappings = [("department", "string", "department", "string"), ("name",
"string",
 "name", "string"), ("id", "int", "id", "int")], transformation_ctx =
"Transform0"]
 ## @return: Transform0
 ## @inputs: [frame = DataSource0]
 Transform0 = ApplyMapping.apply(frame = DataSource0, mappings = [("department",
"string",
 "department", "string"), ("name", "string", "name", "string"), ("id", "int",
"id", "int")],
 transformation_ctx = "Transform0")
 ## @type: DataSink
 ## @args: [connection_type = "s3", format = "json", connection_options = {"path":
"s3://<S3 path>/", "partitionKeys": []}, transformation_ctx = "DataSink0"]
 ## @return: DataSink0
 ## @inputs: [frame = Transform0]
 DataSink0 = glueContext.write_dynamic_frame.from_options(frame = Transform0,
 connection_type = "s3", format = "json", connection_options = {"path":
"s3://<S3 path>/", "partitionKeys": []}, transformation_ctx = "DataSink0")
 job.commit()

```

### custom.spark または marketplace.spark タイプ用の接続オプション

- `className` – (必須) コネクタのクラス名を支援する文字列。
- `secretId` – (オプション) コネクタ接続の認証情報を取得するために使用される文字列。
- `connectionName` – (必須) コネクタに関連付けられている接続の名前を示す文字列。
- その他のオプションは、データストアによって異なります。例えば、[Elasticsearch for Apache Hadoop](#) ドキュメントの説明にあるように、OpenSearch の設定オプションは「es」でプレフィックスされます。Spark から Snowflake への接続では、[Connecting to Snowflake](#) ガイドの「[Using the Spark Connector](#)」で説明されているように、`sfUser` および `sfPassword` のオプションを使用します。

次の Python コード例に、`marketplace.spark` 接続を使用して、OpenSearch のデータストアからの読み取りを実行する方法を示します。

```

import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext

```

```
from awsglue.job import Job

@params: [JOB_NAME]
args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'], args)
@type: DataSource
@args: [connection_type = "marketplace.spark", connection_options =
{"path":"test",
 "es.nodes.wan.only":"true","es.nodes":"https://<AWS endpoint>",
 "connectionName":"test-spark-es","es.port":"443"}, transformation_ctx =
"DataSource0"]
@return: DataSource0
@inputs: []
DataSource0 = glueContext.create_dynamic_frame.from_options(connection_type =
 "marketplace.spark", connection_options = {"path":"test","es.nodes.wan.only":
 "true","es.nodes":"https://<AWS endpoint>","connectionName":
 "test-spark-es","es.port":"443"}, transformation_ctx = "DataSource0")
@type: DataSink
@args: [connection_type = "s3", format = "json", connection_options = {"path":
 "s3://<S3 path>/", "partitionKeys": []}, transformation_ctx = "DataSink0"]
@return: DataSink0
@inputs: [frame = DataSource0]
DataSink0 = glueContext.write_dynamic_frame.from_options(frame = DataSource0,
 connection_type = "s3", format = "json", connection_options = {"path":
 "s3://<S3 path>/", "partitionKeys": []}, transformation_ctx = "DataSink0")
job.commit()
```

## 汎用オプション

このセクションのオプションは `connection_options` として提供されていますが、特定のコネクタには当てはまりません。

以下のパラメータは、ブックマークを設定する際に一般的に使用されます。Amazon S3 または JDBC ワークフローに適用される場合があります。詳細については、「[the section called “ジョブのブックマークを使用する”](#)」を参照してください。

- `jobBookmarkKeys` - 列名の配列。

- `jobBookmarkKeysSortOrder` — ソート順序に基づいて値を比較する方法を定義する文字列。有効な値: "asc"、"desc"。
- `useS3ListImplementation` — Amazon S3 バケットの内容を一覧表示する際のメモリパフォーマンスの管理に使用されます。詳しくは、「[Optimize memory management in AWS Glue](#)」を参照してください。

## AWS Glue for Spark での入出力のデータ形式に関するオプション

これらのページでは、AWS Glue for Spark でサポートされているデータ形式の機能のサポートと設定パラメータについて説明します。この情報の使用と適用可能性の説明については、以下を参照してください。

### AWS Glue でのデータ形式間の機能のサポート

各データ形式が、異なる AWS Glue の機能をサポートする場合があります。次の一般的な機能は、形式のタイプによってサポートされる場合とサポートされない場合があります。データ形式についてのドキュメントを参照し、当社の機能を活用して要件を満たす方法を理解してください。

読み取り	AWS Glue では、コネクタなどのリソースを追加することなく、データ形式を認識し解釈できます。
書き込み	AWS Glue では、リソースを追加することなく、この形式でデータを書き込むことができます。他の Spark 環境と同様に、サードパーティライブラリをジョブに含め、標準の Apache Spark 関数を使用してデータを書き込むことができます。ライブラリを含める方法の詳細については、「 <a href="#">the section called “Python ライブラリ”</a> 」を参照してください。
ストリーミングの読み取り	AWS Glue は、Apache Kafka、Amazon Managed Streaming for Apache Kafka、または Amazon Kinesis メッセージストリームから、このデータ形式を認識して解釈できます。ストリームは一貫した形式でデータを表示すること

を想定しているため、DataFrames として読み込まれます。

小さなファイルのグループ化

AWS Glue は、AWS Glue 変換を実行する際にファイルをグループ化し、各ノードに送信されるバッチを処理できます。これにより、大量の小さなファイルを扱うワークロードのパフォーマンスを大幅に向上できます。詳細については、「[the section called “入力ファイルのグループ化”](#)」を参照してください。

ジョブのブックマーク

AWS Glue は、ジョブブックマークを使用して、ジョブの実行間で同じデータセットに対して同じ作業を実行する変換の進行状況を追跡できます。これにより、最後のジョブ実行以降の新しいデータに対してのみ作業を行う必要があるデータセットを含むワークロードのパフォーマンスを向上させることができます。詳細については、「[the section called “ジョブのブックマークを使用した処理済みデータの追跡”](#)」を参照してください。

## AWS Glue でデータ形式を操作するために使用されるパラメータ

特定の AWS Glue 接続タイプは複数の format タイプをサポートしており、`GlueContext.write_dynamic_frame.from_options` のようなメソッドを使用する場合は、`format_options` オブジェクトを使用してデータ形式に関する情報を指定する必要があります。

- s3 - 詳細については、AWS Glue: [S3 接続パラメータ](#) の「Connection types and options for ETL」(ETL の接続タイプとオプション) を参照してください。また、Python でのこの接続タイプを容易にするメソッド [the section called “create\\_dynamic\\_frame\\_from\\_options”](#) および [the section called “write\\_dynamic\\_frame\\_from\\_options”](#) と、対応する Scala メソッド [the section called “getSourceWithフォーマット”](#) および [the section called “getSinkWithフォーマット”](#) についてのドキュメントもご覧ください。

- kinesis - 詳細については、AWS Glue: [Kinesis 接続パラメータ](#) の「Connection types and options for ETL」(ETL の接続タイプとオプション) を参照してください。また、この接続タイプを容易にするメソッドについてのドキュメント: [the section called “create\\_data\\_frame\\_from\\_options”](#) と、対応する Scala メソッド [the section called “createDataFrameFromOptions”](#) についてのドキュメントもご覧ください。
- kafka - 詳細については、AWS Glue: [Kafka 接続パラメータ](#) の「Connection types and options for ETL」(ETL の接続タイプとオプション) を参照してください。また、この接続タイプを容易にするメソッドについてのドキュメント: [the section called “create\\_data\\_frame\\_from\\_options”](#) と、対応する Scala メソッド [the section called “createDataFrameFromOptions”](#) についてのドキュメントもご覧ください。

一部の接続タイプでは、`format_options` は必要ありません。例えば、通常の使用においては、リレーショナルデータベースへの JDBC 接続では整合性のある表形式のデータ形式でデータを取得します。したがって、JDBC 接続からの読み取りでは、`format_options` は必要ありません。

Glue でデータを読み書きする一部のメソッドでは、`format_options` が必要です。例えば、AWS Glue クローラーで `GlueContext.create_dynamic_frame.from_catalog` を使用します。クローラーによってデータの形式が決定されます。クローラーを使用する場合、AWS Glue 分類子がデータを調べ、データ形式を表す方法について賢く判断します。次に、データの表現を AWS Glue Data Catalog に保存します。これを AWS Glue ETL スクリプト内で使用し、`GlueContext.create_dynamic_frame.from_catalog` メソッドによってデータを取得できます。クローラーを使用すると、データ形式に関する情報を手動で指定する必要がなくなります。

AWS Lake Formation government table にアクセスするジョブの場合、AWS Glue は Lake Formation の government table でサポートされているすべてのフォーマットの読み書きをサポートします。現在サポートされている AWS Lake Formation government table のフォーマット一覧は、AWS Lake Formation デベロッパーガイドの「[ガバメントテーブルの注意点と制限事項](#)」を参照してください。

#### Note

Apache 寄木細工を書くために、AWS Glue ETL では、ダイナミックフレーム用に最適化されたカスタム Parquet ライタータイプのオプションを指定することで、管理されたテーブルへの書き込みのみをサポートします。parquet フォーマットでガバメントテーブルに書き込む場合は、テーブルパラメータの `useGlueParquetWriter` の値でキー `true` に追加する必要があります。

## トピック

- [AWS Glue で CSV 形式を使用する](#)
- [AWS Glue で Parquet 形式を使用する](#)
- [AWS Glue で XML 形式を使用する](#)
- [AWS Glue で Avro 形式を使用する](#)
- [AWS Glue で grokLog 形式を使用する](#)
- [AWS Glue で Ion 形式を使用する](#)
- [AWS Glue での JSON フォーマットの使用](#)
- [AWS Glue で ORC 形式を使用する](#)
- [AWS Glue ETL ジョブでのデータレイクフレームワークの使用](#)
- [共有設定リファレンス](#)

### AWS Glue で CSV 形式を使用する

AWS Glue はソースからデータを取得し、さまざまなデータ形式で保存および転送されたターゲットにデータを書き込みます。このドキュメントでは、データが CSV データ形式で保存または転送される場合に、AWS Glue でデータを使用する際に利用できる機能について説明します。

AWS Glue は、カンマ区切り値 (CSV) ファイル形式をサポートしています。この形式は、最小限の行をベースにしたデータ形式です。CSV は厳密に基準に準拠していないことがよくありますが、詳細については [RFC 4180](#) および [RFC 7111](#) を参照してください。

AWS Glue を使用して、Amazon S3 およびストリーミングソースから CSV を読み取ることができ、Amazon S3 に CSV を書き込むこともできます。S3 から、CSV ファイルが含まれる bzip および gzip アーカイブを読み書きすることができます。このページで説明する設定ではなく、[S3 接続パラメータ](#) 上で圧縮動作を設定します。

次の表は、CSV 形式オプションをサポートする一般的な AWS Glue 機能を示しています。

読み込み	書き込み	ストリーミングの読み取り	小さなファイルのグループ化	ジョブのブックマーク
サポート	サポート対象	サポート対象	サポート対象	サポート

例: S3 から CSV ファイルまたはフォルダを読み取る

前提条件: 読み取りたい CSV ファイルまたはフォルダへの S3 パス (s3path) が必要です。

設定: 関数オプションで format="csv" を指定します。connection\_options で、paths キーを使用して s3path を指定します。リーダーが S3 とやり取りする方法は、connection\_options で設定できます。詳細については、AWS Glue: [S3 接続パラメータ](#) の「ETL の接続タイプとオプション」を参照してください。リーダーが CSV ファイルを解釈する方法は、format\_options で設定できます。詳細については、「[CSV 設定リファレンス](#)」を参照してください。

次の AWS Glue ETL スクリプトは、S3 から CSV ファイルまたはフォルダを読み取るプロセスを示しています。

optimizePerformance 設定キーを使用して、一般的なワークフローのパフォーマンスを最適化するカスタム CSV リーダーを提供します。このリーダーがワークロードに適しているかどうかを判断するには、「[the section called “最適化された CSV リーダーを使用する”](#)」を参照してください。

Python

この例では [create\\_dynamic\\_frame.from\\_options](#) メソッドを使用します。

```
Example: Read CSV from S3
For show, we handle a CSV with a header row. Set the withHeader option.
Consider whether optimizePerformance is right for your workflow.

from pyspark.context import SparkContext
from awsglue.context import GlueContext

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)
spark = glueContext.spark_session

dynamicFrame = glueContext.create_dynamic_frame.from_options(
 connection_type="s3",
 connection_options={"paths": ["s3://s3path"]},
 format="csv",
 format_options={
 "withHeader": True,
 # "optimizePerformance": True,
 },
)
```

スクリプト (pyspark.sql.DataFrame) では DataFrame を使用することもできます。

```
dataFrame = spark.read\
 .format("csv")\
 .option("header", "true")\
 .load("s3://s3path")
```

## Scala

この例では [getSourceWithFormat](#) 操作を使用します。

```
// Example: Read CSV from S3
// For show, we handle a CSV with a header row. Set the withHeader option.
// Consider whether optimizePerformance is right for your workflow.

import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.{DynamicFrame, GlueContext}
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 val dynamicFrame = glueContext.getSourceWithFormat(
 formatOptions=JsonOptions("""{"withHeader": true}"""),
 connectionType="s3",
 format="csv",
 options=JsonOptions("""{"paths": ["s3://s3path"], "recurse": true}""")
).getDynamicFrame()
 }
}
```

スクリプト (`org.apache.spark.sql.DataFrame`) では `DataFrame` を使用することもできます。

```
val dataFrame = spark.read
 .option("header", "true")
 .format("csv")
 .load("s3://s3path")
```

例: CSV ファイルおよびフォルダを S3 に書き込む

前提条件: 初期化された DataFrame (dataFrame) または DynamicFrame (dynamicFrame) が必要です。また、予想される S3 出力パスである s3path も必要になります。

設定: 関数オプションで format="csv" を指定します。connection\_options で、paths キーを使用して s3path を指定します。ライターが S3 とやり取りする方法は、connection\_options で設定できます。詳細については、AWS Glue: [S3 接続パラメータ](#) の「ETL の接続タイプとオプション」を参照してください。操作によってファイルの内容が format\_options にどのように書き込まれるかを設定できます。詳細については、「[CSV 設定リファレンス](#)」を参照してください。次の AWS Glue ETL スクリプトは、S3 に CSV ファイルとフォルダを書き込むプロセスを示しています。

Python

この例では [write\\_dynamic\\_frame.from\\_options](#) メソッドを使用します。

```
Example: Write CSV to S3
For show, customize how we write string type values. Set quoteChar to -1 so our
values are not quoted.

from pyspark.context import SparkContext
from awsglue.context import GlueContext

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

glueContext.write_dynamic_frame.from_options(
 frame=dynamicFrame,
 connection_type="s3",
 connection_options={"path": "s3://s3path"},
 format="csv",
 format_options={
 "quoteChar": -1,
 },
)
```

スクリプト (pyspark.sql.DataFrame) では DataFrame を使用することもできます。

```
dataFrame.write\
 .format("csv")\
 .option("quote", None)\
```

```
.mode("append")\
.save("s3://s3path")
```

## Scala

この例では [getSinkWithFormat](#) メソッドを使用します。

```
// Example: Write CSV to S3
// For show, customize how we write string type values. Set quoteChar to -1 so our
// values are not quoted.

import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.{DynamicFrame, GlueContext}
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 glueContext.getSinkWithFormat(
 connectionType="s3",
 options=JsonOptions("""{"path": "s3://s3path"}"""),
 format="csv"
).writeDynamicFrame(dynamicFrame)
 }
}
```

スクリプト (`org.apache.spark.sql.DataFrame`) では `DataFrame` を使用することもできます。

```
dataFrame.write
 .format("csv")
 .option("quote", null)
 .mode("Append")
 .save("s3://s3path")
```

## CSV 設定リファレンス

AWS Glue ライブラリが `format="csv"` を指定している場合には、以下の `format_options` を使用することができます。

- `separator` - 区切り文字を指定します。デフォルトはカンマですが、他の任意の文字を指定できます。
  - タイプ: テキスト、デフォルト: `","`
- `escaper` - エスケープに使用する文字を指定します。このオプションは、CSV ファイルを書き込む場合ではなく、読み取る場合にのみ使用します。有効にすると、直後の文字はそのまま使用されます。ただし、よく知られている小さいエスケープセット (`\n`、`\r`、`\t`、`\0`) を除きます。
  - タイプ: テキスト、デフォルト: なし
- `quoteChar` - 引用に使用する文字を指定します。デフォルト値は二重引用符です。これに `-1` を設定すると、全体的に引用が無効になります。
  - タイプ: テキスト、デフォルト: `'"'`
- `multiLine` - 単一のレコードが複数行にまたがるかどうかを指定します。これが発生するのは、フィールドに引用符で囲まれた改行文字がある場合などです。複数行にまたがるレコードがある場合は、このオプションを `True` に設定する必要があります。 `multiLine` を有効にすると、解析の際にファイル分割をより慎重に行う必要があるため、パフォーマンスが低下する可能性があります。
  - タイプ: ブール値、デフォルト: `false`
- `withHeader` - 最初の行をヘッダーとして扱うかどうかを指定します。このオプションは `DynamicFrameReader` クラスで使用できます。
  - タイプ: ブール値、デフォルト: `false`
- `writeHeader` - 出力にヘッダーを書き込むかどうかを指定します。このオプションは `DynamicFrameWriter` クラスで使用できます。
  - タイプ: ブール値、デフォルト: `true`
- `skipFirst` - 最初のデータ行をスキップするかどうかを指定します。
  - タイプ: ブール値、デフォルト: `false`
- `optimizePerformance` - 高度な SIMD CSV リーダーで、Apache Arrow ベースの列指向メモリ形式を使用するかどうかを指定します。これは、AWS Glue 3.0 以降でのみ使用できます。
  - タイプ: ブール値、デフォルト: `false`
- `strictCheckForQuoting` - CSV を作成する際、Glue は文字列と解釈した値に引用符を追加することがあります。これは、書き出される内容があいまいにならないようにするためです。Glue は、書き出す内容を定める時間を節約するために、特定の状況では引用符が不要でも引用符を追加することがあります。厳密なチェックを有効にすると、より強力な計算が行われ、厳密に必要な場合にのみ引用符が追加されます。これは、AWS Glue 3.0 以降でのみ使用できます。
  - タイプ: ブール値、デフォルト: `false`

## ベクトル化された SIMD CSV リーダーで読み取りパフォーマンスを最適化する

AWS Glue バージョン 3.0 では、行ベースの CSV リーダーと比較して、全体的なジョブパフォーマンスを大幅に高速化できる最適化された CSV リーダーが追加されています。

最適化されたリーダーでは:

- CPU SIMD 命令を使用してディスクから読み取ります
- レコードを列形式 (Apache Arrow) で即座にメモリに書き込みます
- レコードをバッチに分割します

これにより、レコードがバッチ処理されるときや、後で列形式に変換されるときの処理時間が節約されます。一例としては、スキーマを変更したり、列ごとにデータを取得したりする場合などが挙げられます。

最適化されたリーダーを使用するには、`format_options` またはテーブルプロパティで `"optimizePerformance"` を `true` に設定します。

```
glueContext.create_dynamic_frame.from_options(
 frame = datasource1,
 connection_type = "s3",
 connection_options = {"paths": ["s3://s3path"]},
 format = "csv",
 format_options={
 "optimizePerformance": True,
 "separator": ",",
 },
 transformation_ctx = "datasink2")
```

## ベクトル化された CSV リーダーでの制限事項

ベクトル化された CSV リーダーには、次の制限があるので注意が必要です。

- `multiLine` および `escaper` 形式オプションはサポートされません。デフォルトの `escaper` として二重引用符文字 `'\"'` を使用します。これらのオプションを設定すると、AWS Glue は自動的にフォールバックし、行ベースの CSV リーダーを使用するようになります。
- [ChoiceType](#) の `DynamicFrame` の作成はサポートされません。
- [error records](#) の `DynamicFrame` の作成はサポートされません。

- 日本語や中国語など、マルチバイト文字を含む CSV ファイルからの読み取りはサポートされません。

## AWS Glue で Parquet 形式を使用する

AWS Glue はソースからデータを取得し、さまざまなデータ形式で保存および転送されたターゲットにデータを書き込みます。このドキュメントでは、データが Parquet データ形式で保存または転送される場合に、AWS Glue でデータを使用する際に利用できる機能について説明します。

AWS Glue は Parquet 形式の使用をサポートしています。この形式は、パフォーマンス指向の列ベースのデータ形式です。標準局から発行されている形式の概要については、「[Apache Parquet ドキュメントの概要](#)」を参照してください。

AWS Glue を使用して、Amazon S3 およびストリーミングソースから Parquet ファイルを読み取ることができ、Amazon S3 に Parquet ファイルを書き込むこともできます。S3 から、Parquet ファイルが含まれる bzip および gzip アーカイブを読み書きすることができます。このページで説明する設定ではなく、[S3 接続パラメータ](#) 上で圧縮動作を設定します。

次の表は、Parquet 形式オプションをサポートする一般的な AWS Glue 機能を示しています。

読み込み	書き込み	ストリーミングの読み取り	小さなファイルのグループ化	ジョブのブックマーク
サポート	サポート対象	サポート	サポートされていません	サポート対象*

\* AWS Glue バージョン 1.0+ でサポート

例: S3 から Parquet ファイルまたはフォルダを読み取る

前提条件: 読み取りたい Parquet ファイルまたはフォルダへの S3 パス (s3path) が必要です。

設定: 関数オプションで format="parquet" を指定します。connection\_options で、paths キーを使用して s3path を指定します。

リーダーが S3 とやり取りする方法は、connection\_options で設定できます。詳細については、AWS Glue: [S3 接続パラメータ](#) の「ETL の接続タイプとオプション」を参照してください。

リーダーが Parquet ファイルを解釈する方法は、format\_options で設定できます。詳細については、「[Parquet 設定リファレンス](#)」を参照してください。

次の AWS Glue ETL スクリプトは、S3 から Parquet ファイルまたはフォルダを読み取るプロセスを示しています。

## Python

この例では [create\\_dynamic\\_frame.from\\_options](#) メソッドを使用します。

```
Example: Read Parquet from S3

from pyspark.context import SparkContext
from awsglue.context import GlueContext

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)
spark = glueContext.spark_session

dynamicFrame = glueContext.create_dynamic_frame.from_options(
 connection_type = "s3",
 connection_options = {"paths": ["s3://s3path/"]},
 format = "parquet"
)
```

スクリプト (`pyspark.sql.DataFrame`) では `DataFrame` を使用することもできます。

```
dataFrame = spark.read.parquet("s3://s3path/")
```

## Scala

この例では [getSourceWithFormat](#) メソッドを使用します。

```
// Example: Read Parquet from S3

import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.{DynamicFrame, GlueContext}
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 val dynamicFrame = glueContext.getSourceWithFormat(
```

```
 connectionType="s3",
 format="parquet",
 options=JsonOptions("""{"paths": ["s3://s3path"]}""")
).getDynamicFrame()
}
}
```

スクリプト (`org.apache.spark.sql.DataFrame`) では `DataFrame` を使用することもできます。

```
spark.read.parquet("s3://s3path/")
```

例: Parquet ファイルおよびフォルダを S3 に書き込む

前提条件: 初期化された `DataFrame` (`dataFrame`) または `DynamicFrame` (`dynamicFrame`) が必要です。また、予想される S3 出力パスである `s3path` も必要になります。

設定: 関数オプションで `format="parquet"` を指定します。 `connection_options` で、 `paths` キーを使用して `s3path` を指定します。

ライターが S3 と対話する方法を、 `connection_options` でさらに詳しく変更することができます。詳細については、AWS Glue: [S3 接続パラメータ](#) の「ETL の接続タイプとオプション」を参照してください。操作によってファイルの内容が `format_options` にどのように書き込まれるかを設定できます。詳細については、「[Parquet 設定リファレンス](#)」を参照してください。

次の AWS Glue ETL スクリプトは、S3 に Parquet ファイルとフォルダを書き込むプロセスを示しています。

`useGlueParquetWriter` 設定キーを通して、`DynamicFrame` のパフォーマンスを最適化するカスタム Parquet ライターを提供します。このライターがワークロードに適しているかどうかを判断するには、「[Glue Parquet Writer](#)」を参照してください。

Python

この例では [write\\_dynamic\\_frame\\_from\\_options](#) メソッドを使用します。

```
Example: Write Parquet to S3
Consider whether useGlueParquetWriter is right for your workflow.

from pyspark.context import SparkContext
from awsglue.context import GlueContext
```

```
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

glueContext.write_dynamic_frame.from_options(
 frame=dynamicFrame,
 connection_type="s3",
 format="parquet",
 connection_options={
 "path": "s3://s3path",
 },
 format_options={
 # "useGlueParquetWriter": True,
 },
)
```

スクリプト (`pyspark.sql.DataFrame`) では `DataFrame` を使用することもできます。

```
df.write.parquet("s3://s3path/")
```

## Scala

この例では [getSinkWithFormat](#) メソッドを使用します。

```
// Example: Write Parquet to S3
// Consider whether useGlueParquetWriter is right for your workflow.

import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.{DynamicFrame, GlueContext}
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 glueContext.getSinkWithFormat(
 connectionType="s3",
 options=JsonOptions("""{"path": "s3://s3path"}"""),
 format="parquet"
).writeDynamicFrame(dynamicFrame)
 }
}
```

スクリプト (`org.apache.spark.sql.DataFrame`) では `DataFrame` を使用することもできます。

```
df.write.parquet("s3://s3path/")
```

## Parquet 設定リファレンス

AWS Glue ライブラリが `format="parquet"` を指定している場合には、以下の `format_options` を使用することができます。

- `useGlueParquetWriter - DynamicFrame` ワークフローのパフォーマンスを最適化するカスタム Parquet ライターの使用を指定します。使用方法の詳細については、「[Glue Parquet Writer](#)」を参照してください。
  - タイプ: ブール値、デフォルト: `false`
- `compression` - 使用する圧縮コーデックを指定します。値は `org.apache.parquet.hadoop.metadata.CompressionCodecName` との完全な互換性があります。
  - タイプ: 列挙型テキスト、デフォルト: `"snappy"`
  - 値: `"uncompressed"`、`"snappy"`、`"gzip"`、`"lzo"`
- `blockSize` - メモリにバッファされる行グループのサイズを、バイト数で指定します。これはパフォーマンスのチューニングに使用します。サイズは正確に、メガバイト数で分割する必要があります。
  - タイプ: 数値、デフォルト: `134217728`
  - デフォルト値は 128 MB です。
- `pageSize` - ページのサイズをバイト数で指定します。これはパフォーマンスのチューニングに使用します。ページは、単一のレコードにアクセスするために完全に読み取らなければならない最小単位です。
  - タイプ: 数値、デフォルト: `1048576`
  - デフォルト値は 1 MB です。

### Note

加えて、基盤となる SparkSQL コードで受け入れられるすべてのオプションは、`connection_options` マップパラメータを介し、この形式に渡されます。例え

ば、[mergeSchema](#) などの Spark 設定を AWS Glue Spark リーダーに行うことで、すべてのファイルのスキーマをマージすることができます。

## 書き込みパフォーマンスを AWS Glue Parquet ライターで最適化する

### Note

AWS Glue Parquet ライターには、これまで glueparquet 形式タイプでアクセスしていました。このアクセス方法は非推奨となりました。代わりに、`useGlueParquetWriter` を有効にして `parquet` タイプを使用してください。

AWS Glue Parquet ライターには、Parquet ファイルの書き込みを高速化できるパフォーマンス強化が組み込まれています。これまでのライターでは、書き込む前にスキーマを計算していました。Parquet 形式では、すぐに取得できる形でスキーマを保存しないため、時間がかかることがありました。AWS Glue Parquet ライターを使用すれば、事前に計算済みのスキーマは必要なくなります。データが到着次第、ライターがスキーマを動的に計算して変更します。

`useGlueParquetWriter` を指定する場合、以下の制限に注意してください。

- ライターは、列の追加や削除などのスキーマの進化のみをサポートし、`ResolveChoice` のように列タイプの変更はサポートしません。
- ライターは、空の `DataFrame` の書き込み (スキーマのみのファイルを書き込む場合など) をサポートしていません。`enableUpdateCatalog=True` の設定によって AWS Glue データカタログと統合する場合、空の `DataFrame` を書き込もうとしてもデータカタログは更新されません。これにより、スキーマなしでデータカタログにテーブルが作成されます。

変換時にこれらの制限が必要とならない場合、AWS Glue Parquet ライターをオンにすることで、パフォーマンスが向上するはずですが。

## AWS Glue で XML 形式を使用する

AWS Glue はソースからデータを取得し、さまざまなデータ形式で保存および転送されたターゲットにデータを書き込みます。このドキュメントでは、データが XML データ形式で保存または転送される場合に、AWS Glue でデータを使用する際に利用できる機能について説明します。

AWS Glue は XML 形式の使用をサポートしています。この形式は、行ベースでも列ベースでもなく、非常に設定可能性が高く、厳密に定義されたデータ構造を表します。XML は高度に標準化され

ています。標準局から発行されている形式の概要については、「[XML エッセンシャル](#)」を参照してください。

AWS Glue を使用して、Amazon S3 から XML ファイルを読み込んだり、XML ファイルが含まれる bzip および gzip アーカイブを読み込んだりすることができます。このページで説明する設定ではなく、[S3 接続パラメータ](#) 上で圧縮動作を設定します。

次の表は、XML 形式オプションをサポートする一般的な AWS Glue 機能を示しています。

読み込み	書き込み	ストリーミングの読み取り	小さなファイルのグループ化	ジョブのブックマーク
サポート	サポートされていません	サポートされていません	サポート	サポート

例: S3 から XML を読み込む

XML リーダーは XML タグ名を取得します。入力内にあるそのタグを持つ要素を調べてスキーマを推測し、対応する値を DynamicFrame に入力します。AWS Glue XML の機能は、[Apache Spark の XML データソース](#)に類似した動作をします。このリーダーをそのプロジェクトのドキュメントと比較することで、基本的な動作に関する洞察が得られる可能性があります。

前提条件: 読み取りたい XML ファイルまたはフォルダへの S3 パス (s3path) と、XML ファイルに関するいくつかの情報が必要です。また、読み込みたい XML 要素のタグである xmlTag が必要です。

設定: 関数オプションで format="xml" を指定します。connection\_options で、paths キーを使用して s3path を指定します。リーダーが S3 とやり取りする方法は、connection\_options でさらに詳しく設定できます。詳細については、AWS Glue: [S3 接続パラメータ](#) の「ETL の接続タイプとオプション」を参照してください。format\_options で、rowTag キーを使用して xmlTag を指定します。リーダーが XML ファイルを解釈する方法は、format\_options でさらに詳しく設定できます。詳細については、「[XML 設定リファレンス](#)」を参照してください。

次の AWS Glue ETL スクリプトは、S3 から XML ファイルまたはフォルダを読み取るプロセスを示しています。

Python

この例では [create\\_dynamic\\_frame.from\\_options](#) メソッドを使用します。

```
Example: Read XML from S3
Set the rowTag option to configure the reader.

from awsglue.context import GlueContext
from pyspark.context import SparkContext

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

dynamicFrame = glueContext.create_dynamic_frame.from_options(
 connection_type="s3",
 connection_options={"paths": ["s3://s3path"]},
 format="xml",
 format_options={"rowTag": "xmlTag"},
)
```

スクリプト (`pyspark.sql.DataFrame`) では `DataFrame` を使用することもできます。

```
dataFrame = spark.read\
 .format("xml")\
 .option("rowTag", "xmlTag")\
 .load("s3://s3path")
```

## Scala

この例では [getSourceWithFormat](#) 操作を使用します。

```
// Example: Read XML from S3
// Set the rowTag option to configure the reader.

import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.GlueContext
import org.apache.spark.sql.SparkSession

val glueContext = new GlueContext(SparkContext.getOrCreate())
val sparkSession: SparkSession = glueContext.getSparkSession

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val dynamicFrame = glueContext.getSourceWithFormat(
 formatOptions=JsonOptions("""{"rowTag": "xmlTag"}"""),
 connectionType="s3",
)
 }
}
```

```
format="xml",
options=JsonOptions("""{"paths": ["s3://s3path"], "recurse": true}""")
).getDynamicFrame()
}
```

スクリプト (`org.apache.spark.sql.DataFrame`) では `DataFrame` を使用することもできます。

```
val dataframe = spark.read
 .option("rowTag", "xmlTag")
 .format("xml")
 .load("s3://s3path")
```

## XML 設定リファレンス

AWS Glue ライブラリが `format="xml"` を指定している場合には、以下の `format_options` を使用することができます。

- `rowTag` - 行として扱うファイル内の XML タグを指定します。行のタグを自己終了型にすることはできません。
  - タイプ: テキスト、必須
- `encoding` - 文字エンコードを指定します。これは、ランタイム環境がサポートする [Charset](#) の名前またはエイリアスにすることができます。エンコーディングのサポートに関しては特に保証していませんが、主要なエンコーディングは動作するはずです。
  - タイプ: テキスト、デフォルト: "UTF-8"
- `excludeAttribute` - 要素の属性を除外するかどうかを指定します。
  - タイプ: ブール値、デフォルト: `false`
- `treatEmptyValuesAsNulls` - 空白文字を `null` 値として扱うかどうかをします。
  - タイプ: ブール値、デフォルト: `false`
- `attributePrefix` - 子要素テキストから区別するために属性に付加するプレフィックス。このプレフィックスをフィールド名として使用します。
  - タイプ: テキスト、デフォルト: "\_"
- `valueTag` - 要素内に子を持たない属性がある場合、値に使用するタグ。
  - タイプ: テキスト、デフォルト: "\_VALUE"

- `ignoreSurroundingSpaces` - 値を囲む空白文字を無視するかどうかを指定します。
  - タイプ: ブール値、デフォルト: `false`
- `withSchema` - 推論されたスキーマを上書きしたい場合に、期待されるスキーマを含みます。このオプションを使用しない場合は、AWS Glue で XML データからスキーマを推定します。
  - タイプ: テキスト、デフォルト: 該当なし
  - 値は `StructType` を表す JSON オブジェクトである必要があります。

## XML スキーマを手動で指定する

### 手動 XML スキーマの例

これは、XML データ用にスキーマを指定するために、`withSchema` 形式オプションを使用する場合の例です。

```
from awsglue.gluetypes import *

schema = StructType([
 Field("id", IntegerType()),
 Field("name", StringType()),
 Field("nested", StructType([
 Field("x", IntegerType()),
 Field("y", StringType()),
 Field("z", ChoiceType([IntegerType(), StringType()]))
]))
])

datasource0 = create_dynamic_frame_from_options(
 connection_type,
 connection_options={"paths": ["s3://xml_bucket/someprefix"]},
 format="xml",
 format_options={"withSchema": json.dumps(schema.jsonValue())},
 transformation_ctx = ""
)
```

## AWS Glue で Avro 形式を使用する

AWS Glue はソースからデータを取得し、さまざまなデータ形式で保存および転送されたターゲットにデータを書き込みます。このドキュメントでは、データが Avro データ形式で保存または転送される場合に、AWS Glue でデータを使用する際に利用できる機能について説明します。

AWS Glue は Avro 形式の使用をサポートしています。この形式は、パフォーマンス指向の行ベースのデータ形式です。標準局による形式の概要については、「[Apache Avro 1.8.2 Documentation](#)」(Apache Avro ドキュメント 1.8.2) を参照してください。

AWS Glue を使用して、Amazon S3 およびストリーミングソースから Avro ファイルを読み取り、Avro ファイルを Amazon S3 に書き込むことができます。S3 から、Avro ファイルが含まれる bzip2 および gzip アーカイブを読み書きすることができます。さらに、Avro ファイルが含まれる deflate、snappy、xz アーカイブを書き込むこともできます。このページで説明する設定ではなく、[S3 接続パラメータ](#) 上で圧縮動作を設定します。

次の表は、Avro 形式のオプションをサポートする一般的な AWS Glue オペレーションを示しています。

読み込み	書き込み	ストリーミングの読み取り	小さなファイルのグループ化	ジョブのブックマーク
サポート	サポート	サポート対象*	サポートされていません	サポート

制限付きで、\* がサポートされています。詳細については、「[the section called “Avro ストリーミングソースに関する注意事項と制約事項”](#)」を参照してください。

例: S3 から Avro ファイルまたはフォルダを読み取る

前提条件: 読み取る Avro ファイルまたはフォルダへの S3 パス (s3path) が必要です。

設定: 関数オプションで format="avro" を指定します。connection\_options で、paths キーを使用して s3path を指定します。リーダーが S3 とやり取りする方法は、connection\_options で設定できます。詳細については、AWS Glue: [the section called “S3 接続パラメータ”](#) の「ETL の入力および出力のデータ形式オプション」を参照してください。リーダーが Avro ファイルを解釈する方法は、format\_options で設定できます。詳細については、[Avro 設定リファレンス](#)を参照してください。

次の AWS Glue ETL スクリプトは、S3 から Avro ファイルまたはフォルダを読み取るプロセスを示しています。

Python

この例では [create\\_dynamic\\_frame.from\\_options](#) メソッドを使用します。

```
from pyspark.context import SparkContext
from awsglue.context import GlueContext

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

dynamicFrame = glueContext.create_dynamic_frame.from_options(
 connection_type="s3",
 connection_options={"paths": ["s3://s3path"]},
 format="avro"
)
```

## Scala

この例では [getSourceWithFormat](#) 操作を使用します。

```
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.GlueContext
import org.apache.spark.sql.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 val dynamicFrame = glueContext.getSourceWithFormat(
 connectionType="s3",
 format="avro",
 options=JsonOptions("""{"paths": ["s3://s3path"]}""")
).getDynamicFrame()
 }
}
```

例: Avro ファイルおよびフォルダを S3 に書き込む

前提条件: 初期化された DataFrame (dataFrame) または DynamicFrame (dynamicFrame) が必要です。また、予想される S3 出力パスである s3path も必要になります。

設定: 関数オプションで format="avro" を指定します。connection\_options で、paths キーを使用して s3path を指定します。ライターが S3 と対話する方法を、connection\_options でさらに詳しく変更することができます。詳細については、AWS Glue: [the section called “S3 接続パラメータ”](#) の「ETL の入力および出力のデータ形式オプション」を参照してください。リーダーが

Avro ファイルを解釈する方法は、`format_options` で変更できます。詳細については、[Avro 設定リファレンス](#)を参照してください。

次の AWS Glue ETL スクリプトは、S3 に Avro ファイルまたはフォルダを書き込むプロセスを示しています。

## Python

この例では [write\\_dynamic\\_frame.from\\_options](#) メソッドを使用します。

```
from pyspark.context import SparkContext
from awsglue.context import GlueContext

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

glueContext.write_dynamic_frame.from_options(
 frame=dynamicFrame,
 connection_type="s3",
 format="avro",
 connection_options={
 "path": "s3://s3path"
 }
)
```

## Scala

この例では [getSinkWithFormat](#) メソッドを使用します。

```
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.{DynamicFrame, GlueContext}
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 glueContext.getSinkWithFormat(
 connectionType="s3",
 options=JsonOptions("""{"path": "s3://s3path"}"""),
 format="avro"
).writeDynamicFrame(dynamicFrame)
 }
}
```

```
}
}
```

## Avro 設定リファレンス

AWS Glue ライブラリが `format="avro"` を指定する場合は、以下の `format_options` を使用することができます。

- `version` – サポートする Apache Avro リーダー/ライター形式のバージョンを指定します。デフォルト値は「1.7」です。Avro 論理型の読み取りと書き込みを有効にするには、`format_options={"version": "1.8"}` を指定します。詳細については、「[Apache Avro 1.7.7 の仕様](#)」および「[Apache Avro 1.8.2 の仕様](#)」を参照してください。

Apache Avro 1.8 コネクタは、次の論理型変換をサポートしています。

リーダーの場合: この表に、Avro リーダー 1.7 と 1.8 について、Avro データ型 (論理型と Avro プリミティブ型) と AWS GlueDynamicFrame データ型の間での型変換を示します。

Avro データ型: 論理型	Avro データ型: Avro プリミティブ型	GlueDynamicFrame データ型: Avro リーダー 1.7	GlueDynamicFrame データ型: Avro リーダー 1.8
10 進数	bytes	BINARY	10 進数
10 進数	固定	BINARY	10 進数
日付	整数	INT	日付
時間 (ミリ秒)	整数	INT	INT
時間 (マイクロ秒)	long	LONG	LONG
タイムスタンプ (ミリ秒)	long	LONG	タイムスタンプ
タイムスタンプ (マイクロ秒)	long	LONG	LONG

Avro データ型: 論理型	Avro データ型: Avro プリミティブ型	GlueDynamicFrame データ型: Avro リーダー 1.7	GlueDynamicFrame データ型: Avro リーダー 1.8
所要時間 (論理型では ない)	12 で固定	BINARY	BINARY

ライターの場合: この表に、Avro ライター 1.7 と 1.8 について、データ型と Avro データ型 AWS GlueDynamicFrame の間での型変換を示します。

AWS GlueDynamicFrame データ型	Avro データ型: Avro ライター 1.7	Avro データ型: Avro タイター 1.8
10 進数	文字列	decimal
日付	文字列	date
タイムスタンプ	文字列	timestamp-micros

## Avro Spark DataFrame サポート

Spark DataFrame API から Avro を使用するには、対応する Spark のバージョン用の Spark Avro プラグインをインストールする必要があります。ジョブで使用できる Spark のバージョンは、AWS Glue のバージョンによって決まります。Spark のバージョンの詳細については、「[the section called “AWS Glue バージョン”](#)」を参照してください。このプラグインは Apache によって管理されており、特定のサポートを保証するものではありません。

AWS Glue 2.0 - Spark Avro プラグインのバージョン 2.4.3 を使用します。この JAR は Maven Central にあります。[org.apache.spark:spark-avro\\_2.12:2.4.3](#) を参照してください。

AWS Glue 3.0 - Spark Avro プラグインのバージョン 3.1.1 を使用します。この JAR は Maven Central にあります。[org.apache.spark:spark-avro\\_2.12:3.1.1](#) を参照してください。

AWS Glue の ETL ジョブに追加の JAR を含めるには、`--extra-jars` ジョブパラメータを使用します。ジョブパラメータについては、「[the section called “ジョブのパラメータ”](#)」を参照してください。AWS Management Console で、このパラメータを設定することもできます。

## AWS Glue で grokLog 形式を使用する

AWS Glue はソースからデータを取得し、さまざまなデータ形式で保存および転送されたターゲットにデータを書き込みます。このドキュメントでは、データが大まかに構造化されたプレーンテキスト形式で保存または転送される場合に、AWS Glue でデータを使用する際に利用できる機能を紹介します。

AWS Glue は、Grok パターンの使用をサポートしています。Grok パターンは、正規表現のキャプチャグループに似ています。プレーンテキストファイル内の文字シーケンスのパターンを認識し、タイプと目的を与えます。AWS Glue での主な目的はログを読み取ることです。作成者による Grok の概要については、「[Logstash Reference: Grok filter plugin](#)」(Logstash リファレンス: Grok フィルタープラグイン) を参照してください。

読み込み	書き込み	ストリーミングの読み取り	小さなファイルのグループ化	ジョブのブックマーク
サポート	該当しません	サポート	サポート	サポートされていません

## grokLog 設定リファレンス

`format="grokLog"` には、以下の `format_options` 値を使用できます。

- `logFormat` – ログの形式と一致する Grok パターンを指定します。
- `customPatterns` – ここで使用する追加の Grok パターンを指定します。
- `MISSING` – 欠落した値の識別に使用するシグナルを指定します。デフォルト: `'-'`。
- `LineCount` – 各ログレコードの行数を指定します。デフォルト値は `'1'` です。現在 1 行のレコードのみがサポートされています。
- `StrictMode` – `strict` モードを有効にするかどうかを指定するブール値。厳格モードでは、リーダーは自動的な型変換や復旧を行いません。デフォルト値は `"false"` です。

## AWS Glue で Ion 形式を使用する

AWS Glue はソースからデータを取得し、さまざまなデータ形式で保存および転送されたターゲットにデータを書き込みます。このドキュメントでは、データが Ion データ形式で保存または転送される場合に、AWS Glue でデータを使用する際に利用できる機能を紹介します。

AWS Glue は Ion 形式の使用をサポートしています。この形式は、データ構造 (行または列ベースではない) を交換可能なバイナリおよびプレーンテキスト形式で表します。作成者による形式の概要については、「[Amazon Ion](#)」を参照してください。詳細については、[Amazon Ion の仕様](#)に関するドキュメントを参照してください。

AWS Glue を使用して、Amazon S3 から Ion ファイルを読み取ることができます。S3 から、Ion ファイルを含む bzip および gzip アーカイブを読み取ることができます。このページで説明する設定ではなく、[S3 接続パラメータ](#) 上で圧縮動作を設定します。

次の表は、Ion 形式のオプションをサポートする一般的な AWS Glue オペレーションを示しています。

読み込み	書き込み	ストリーミングの読み取り	小さなファイルのグループ化	ジョブのブックマーク
サポート	サポートされていません	サポートされていません	サポート	サポートされていません

例: S3 から Ion ファイルとフォルダを読み取る

前提条件: 読み取る Ion ファイルまたはフォルダへの S3 パス (s3path) が必要です。

設定: 関数オプションで format="json" を指定します。connection\_options で、paths キーを使用して s3path を指定します。リーダーが S3 とやり取りする方法は、connection\_options で設定できます。詳細については、AWS Glue: [the section called “S3 接続パラメータ”](#) の「ETL の接続タイプとオプション」を参照してください。

次の AWS Glue ETL スクリプトは、S3 から Ion ファイルまたはフォルダを読み取るプロセスを示しています。

Python

この例では [create\\_dynamic\\_frame.from\\_options](#) メソッドを使用します。

```
Example: Read ION from S3

from pyspark.context import SparkContext
from awsglue.context import GlueContext
```

```
sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

dynamicFrame = glueContext.create_dynamic_frame.from_options(
 connection_type="s3",
 connection_options={"paths": ["s3://s3path"]},
 format="ion"
)
```

## Scala

この例では [getSourceWithFormat](#) 操作を使用します。

```
// Example: Read ION from S3

import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.GlueContext
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 val dynamicFrame = glueContext.getSourceWithFormat(
 connectionType="s3",
 format="ion",
 options=JsonOptions("""{"paths": ["s3://s3path"], "recurse": true}""")
).getDynamicFrame()
 }
}
```

## Ion 設定リファレンス

`format="ion"` の `format_options` 値はありません。

## AWS Glue での JSON フォーマットの使用

AWS Glue はソースからデータを取得し、さまざまなデータ形式で保存および転送されるターゲットにデータを書き込みます。データが JSON データ形式で保存または転送される場合、このドキュメントでは AWS Glue でデータを使用するために利用できる機能を紹介します。

AWS Glue は JSON 形式の使用をサポートしています。この形式は、行ベースまたは列ベースではなく、形状は一貫しているが内容は柔軟なデータ構造を表します。JSON は、複数の機関が発行する並列した規格で定義されており、そのうちの 1 つが ECMA-404 です。一般的に参照されるソースによる形式の概要については、「[Introducing JSON](#)」(JSON の概要) を参照してください。

AWS Glue を使用して、Amazon S3 から JSON ファイルを読み取るだけでなく bzip、gzip 圧縮された JSON ファイルも読み取ることができます。このページで説明する設定ではなく、[S3 接続パラメータ](#) 上で圧縮動作を設定します。

読み込み	書き込み	ストリーミングの読み取り	小さなファイルのグループ化	ジョブのブックマーク
サポート	サポート対象	サポート対象	サポート対象	サポート

例: S3 から JSON ファイルまたはフォルダを読み取る

前提条件: 読み取る JSON ファイルまたはフォルダへの S3 パス (s3path) が必要です。

設定: 関数オプションで format="json" を指定します。connection\_options で、paths キーを使用して s3path を指定します。さらに、接続オプションで、読み取り操作が S3 を通過する方法を変更することができます。詳細については、「[the section called “S3 接続パラメータ”](#)」を参照してください。リーダーが JSON ファイルを解釈する方法は、format\_options で設定できます。詳細については、[JSON 設定リファレンス](#)を参照してください。

次の AWS Glue ETL スクリプトは、S3 から JSON ファイルまたはフォルダを読み取るプロセスを示しています。

Python

この例では [create\\_dynamic\\_frame.from\\_options](#) メソッドを使用します。

```
Example: Read JSON from S3
For show, we handle a nested JSON file that we can limit with the JsonPath
parameter
For show, we also handle a JSON where a single entry spans multiple lines
Consider whether optimizePerformance is right for your workflow.

from pyspark.context import SparkContext
```

```
from awsglue.context import GlueContext

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)
spark = glueContext.spark_session

dynamicFrame = glueContext.create_dynamic_frame.from_options(
 connection_type="s3",
 connection_options={"paths": ["s3://s3path"]},
 format="json",
 format_options={
 "jsonPath": "$.id",
 "multiline": True,
 # "optimizePerformance": True, -> not compatible with jsonPath, multiline
 }
)
```

スクリプト (`pyspark.sql.DataFrame`) `DataFrames` でも使用できます。

```
dataFrame = spark.read\
 .option("multiLine", "true")\
 .json("s3://s3path")
```

## Scala

この例では、[getSourceWithFormat](#) オペレーションを使用します。

```
// Example: Read JSON from S3
// For show, we handle a nested JSON file that we can limit with the JsonPath
// parameter
// For show, we also handle a JSON where a single entry spans multiple lines
// Consider whether optimizePerformance is right for your workflow.

import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.{DynamicFrame, GlueContext}
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 val dynamicFrame = glueContext.getSourceWithFormat(
```

```
formatOptions=JsonOptions("""{"jsonPath": "$.id", "multiline": true,
"optimizePerformance":false}"""),
connectionType="s3",
format="json",
options=JsonOptions("""{"paths": ["s3://s3path"], "recurse": true}""")
).getDynamicFrame()
}
}
```

スクリプト (`pyspark.sql.DataFrame`) `DataFrames` でも使用できます。

```
val dataframe = spark.read
 .option("multiLine", "true")
 .json("s3://s3path")
```

例: JSON ファイルおよびフォルダを S3 に書き込む

前提条件:初期化された `DataFrame` (`dataFrame`) または `()` が必要です。 `DynamicFrame` `dynamicFrame` また、予想される S3 出力パスである `s3path` も必要になります。

設定: 関数オプションで `format="json"` を指定します。 `connection_options` で、 `paths` キーを使用して `s3path` を指定します。ライターが S3 と対話する方法を、 `connection_options` でさらに詳しく変更することができます。詳細については、AWS Glue: [the section called “S3 接続パラメータ”](#) の ETL 入力と出力のデータ形式オプションを参照してください。ライターが JSON ファイルを解釈する方法は、 `format_options` で設定できます。詳細については、[JSON 設定リファレンス](#)を参照してください。

次の AWS Glue ETL スクリプトは、S3 から JSON ファイルまたはフォルダを書き込むプロセスを示しています。

Python

この例では [write\\_dynamic\\_frame\\_from\\_options](#) メソッドを使用します。

```
Example: Write JSON to S3

from pyspark.context import SparkContext
from awsglue.context import GlueContext

sc = SparkContext.getOrCreate()
```

```
glueContext = GlueContext(sc)

glueContext.write_dynamic_frame.from_options(
 frame=dynamicFrame,
 connection_type="s3",
 connection_options={"path": "s3://s3path"},
 format="json"
)
```

スクリプト (pyspark.sql.DataFrame) DataFrames でも使用できます。

```
df.write.json("s3://s3path/")
```

## Scala

この例では、[getSinkWithFormat](#) メソッドを使用します。

```
// Example: Write JSON to S3

import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.{DynamicFrame, GlueContext}
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 glueContext.getSinkWithFormat(
 connectionType="s3",
 options=JsonOptions("""{"path": "s3://s3path"}"""),
 format="json"
).writeDynamicFrame(dynamicFrame)
 }
}
```

スクリプト (pyspark.sql.DataFrame) DataFrames でも使用できます。

```
df.write.json("s3://s3path")
```

## JSON 設定リファレンス

format="json" には、以下の format\_options 値を使用できます。

- jsonPath— [JsonPath](#)レコードに読み込まれるオブジェクトを識別する式。これは、ファイルが外部配列内にネストされたレコードを含む場合に役立ちます。たとえば、JsonPath 次の式は JSON id オブジェクトのフィールドを対象としています。

```
format="json", format_options={"jsonPath": "$.id"}
```

- multiLine — 単一のレコードが複数行にまたがるかどうかを指定するブール値。これが発生するのは、フィールドに引用符で囲まれた改行文字がある場合などです。複数行にまたがるレコードがある場合は、このオプションを "true" に設定する必要があります。デフォルト値は "false" であり、解析時によりアグレッシブなファイル分割を可能にします。
- optimizePerformance – 高度な SIMD JSON リーダーで、Apache Arrow ベースの列指向メモリ形式を使用するかどうかを指定するブール値。これは、AWS Glue 3.0 でのみ使用できます。multiLine または jsonPath と互換性がありません。これらのオプションのいずれかを指定すると、AWS Glue は標準リーダーにフォールバックするよう指示します。
- withSchema — 「[the section called “XML スキーマを指定する”](#)」で説明されている形式でテーブルスキーマを指定する文字列の値。非カタログ接続から読み取る場合のみ、optimizePerformance で使用されます。

### Apache Arrow 列指向形式によりベクトル化された SIMD JSON リーダーの使用

AWS Glue バージョン 3.0 では、JSON データ用のベクター化されたリーダーが追加されています。特定の条件下では、標準のリーダーと比較して 2 倍高速に動作します。このリーダーには、このセクションに記載されているように、ユーザーが使用前に知っておく必要のある特定の制限があります。

最適化されたリーダーを使用するには、format\_options またはテーブルプロパティで "optimizePerformance" を True に設定します。カタログを読み取らない限り、withSchema を指定する必要もあります。withSchema は、「[the section called “XML スキーマを指定する”](#)」で説明されている入力を期待します

```
// Read from S3 data source
glueContext.create_dynamic_frame.from_options(
 connection_type = "s3",
```

```
connection_options = {"paths": ["s3://s3path"]},
format = "json",
format_options={
 "optimizePerformance": True,
 "withSchema": SchemaString
})

// Read from catalog table
glueContext.create_dynamic_frame.from_catalog(
 database = database,
 table_name = table,
 additional_options = {
 // The vectorized reader for JSON can read your schema from a catalog table
 property.
 "optimizePerformance": True,
 })
```

AWS Glue *SchemaString* ライブラリでビルドする方法の詳細については、[を参照してください](#) [the section called “型”](#)。

## ベクトル化された CSV リーダーでの制限事項

以下の制限事項に留意してください。

- ネストされたオブジェクトまたは配列値を持つ JSON 要素はサポートされていません。指定した場合、AWS Glue は標準リーダーにフォールバックします。
- カタログから、または withSchema で、スキーマを指定する必要があります。
- multiLine または jsonPath と互換性がありません。これらのオプションのいずれかを指定すると、AWS Glue は標準リーダーにフォールバックするよう指示します。
- 入力スキーマと一致しない入力レコードを指定すると、リーダーは失敗します。
- [エラーレコード](#) は作成されません。
- マルチバイト文字 (日本語や中国語の文字など) を含む JSON ファイルはサポートされていません。

## AWS Glue で ORC 形式を使用する

AWS Glue はソースからデータを取得し、さまざまなデータ形式で保存および転送されたターゲットにデータを書き込みます。このドキュメントでは、データが ORC データ形式で保存または転送される場合に、AWS Glue でデータを使用する際に利用できる機能について説明します。

AWS Glue は、ORC 形式の使用をサポートしています。この形式は、パフォーマンス指向の列ベースのデータ形式です。標準局による形式の概要については、「[Apache Orc](#)」を参照してください。

AWS Glue を使用して、Amazon S3 およびストリーミングソースから ORC ファイルを読み取り、Amazon S3 に ORC ファイルを書き込むことができます。S3 から、ORC ファイルを含む bzip および gzip アーカイブを読み書きすることができます。このページで説明する設定ではなく、[S3 接続パラメータ](#) 上で圧縮動作を設定します。

次の表は、ORC 形式のオプションをサポートする一般的な AWS Glue の機能を示しています。

読み込み	書き込み	ストリーミングの読み取り	小さなファイルのグループ化	ジョブのブックマーク
サポート	サポート対象	サポート	サポートされていません	サポート対象*

\* AWS Glue バージョン 1.0+ でサポート

例: S3 から ORC ファイルまたはフォルダを読み取る

前提条件: 読み取る ORC ファイルまたはフォルダへの S3 パス (s3path) が必要です。

設定: 関数オプションで format="orc" を指定します。connection\_options で、paths キーを使用して s3path を指定します。リーダーが S3 とやり取りする方法は、connection\_options で設定できます。詳細については、AWS Glue: [the section called "S3 接続パラメータ"](#) の「ETL の接続タイプとオプション」を参照してください。

次の AWS Glue ETL スクリプトは、S3 から ORC ファイルまたはフォルダを読み取るプロセスを示しています。

Python

この例では [create\\_dynamic\\_frame.from\\_options](#) メソッドを使用します。

```
from pyspark.context import SparkContext
from awsglue.context import GlueContext

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)
```

```
dynamicFrame = glueContext.create_dynamic_frame.from_options(
 connection_type="s3",
 connection_options={"paths": ["s3://s3path"]},
 format="orc"
)
```

スクリプト (pyspark.sql.DataFrame) では DataFrame を使用することもできます。

```
dataFrame = spark.read\
 .orc("s3://s3path")
```

## Scala

この例では [getSourceWithFormat](#) 操作を使用します。

```
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.GlueContext
import org.apache.spark.sql.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 val dynamicFrame = glueContext.getSourceWithFormat(
 connectionType="s3",
 format="orc",
 options=JsonOptions("""{"paths": ["s3://s3path"]}""")
).getDynamicFrame()
 }
}
```

スクリプト (pyspark.sql.DataFrame) では DataFrame を使用することもできます。

```
val dataFrame = spark.read\
 .orc("s3://s3path")
```

例: ORC ファイルおよびフォルダを S3 に書き込む

前提条件: 初期化された DataFrame (dataFrame) または DynamicFrame (dynamicFrame) が必要です。また、予想される S3 出力パスである s3path も必要になります。

設定: 関数オプションで `format="orc"` を指定します。接続オプションでは、`s3path` を指定するための `paths` キーを使用します。ライターが S3 と対話する方法を、`connection_options` でさらに詳しく変更することができます。詳細については、AWS Glue: [the section called “S3 接続パラメータ”](#) の「ETL の入力および出力のデータ形式オプション」を参照してください。次のコード例は、プロセスを示しています。

## Python

この例では [write\\_dynamic\\_frame.from\\_options](#) メソッドを使用します。

```
from pyspark.context import SparkContext
from awsglue.context import GlueContext

sc = SparkContext.getOrCreate()
glueContext = GlueContext(sc)

glueContext.write_dynamic_frame.from_options(
 frame=dynamicFrame,
 connection_type="s3",
 format="orc",
 connection_options={
 "path": "s3://s3path"
 }
)
```

スクリプト (`pyspark.sql.DataFrame`) では `DataFrame` を使用することもできます。

```
df.write.orc("s3://s3path/")
```

## Scala

この例では [getSinkWithFormat](#) メソッドを使用します。

```
import com.amazonaws.services.glue.util.JsonOptions
import com.amazonaws.services.glue.{DynamicFrame, GlueContext}
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 }
}
```

```
glueContext.getSinkWithFormat(
 connectionType="s3",
 options=JsonOptions("""{"path": "s3://s3path"}"""),
 format="orc"
).writeDynamicFrame(dynamicFrame)
}
}
```

スクリプト (`pyspark.sql.DataFrame`) では `DataFrame` を使用することもできます。

```
df.write.orc("s3://s3path/")
```

## ORC 設定リファレンス

`format="orc"` の `format_options` 値はありません。ただし、基になる SparkSQL コードで受け入れられるオプションは、`connection_options` マップパラメータを介して渡すことができます。

## AWS Glue ETL ジョブでのデータレイクフレームワークの使用

オープンソースのデータレイクフレームワークを使用するいことで、Amazon S3 上に構築されたデータレイクに保存するファイルの増分データ処理を簡素化できます。AWS Glue 3.0 以降で、以下のオープンソースデータレイクフレームワークがサポートされています。

- Apache Hudi
- Linux Foundation Delta Lake
- Apache Iceberg

これらのフレームワークはネイティブにサポートされているので、Amazon S3 に保存するデータに対する読み書きが、一貫性のあるトランザクションにより実行できます。これらのフレームワークを AWS Glue ETL ジョブで使用する場合にも、別のコネクタをインストールしたり、追加の構成手順を実行したりする必要はありません。

AWS Glue Data Catalog を介してデータセットを管理していれば、Spark DataFrames を使用してデータレイクテーブルを読み書きする場合に、AWS Glue メソッドを利用できます。また、Spark DataFrame API を使用して、Amazon S3 データを読み書きすることもできます。

このビデオでは、Apache Hudi、Apache Iceberg、Delta Lake の仕組みの基本について学ぶことができます。データレイクにデータを挿入、更新、削除する方法と、これらの各フレームワークの仕組みについて説明します。

## トピック

- [制限事項](#)
- [AWS Glue での Hudi フレームワークの使用](#)
- [AWS Glue での Delta Lake フレームワークの使用](#)
- [AWS Glue での Iceberg フレームワークの使用](#)

## 制限事項

でデータレイクフレームワークを使用する前に、次の制限を考慮してください。AWS Glue

- AWS Glue GlueContextの以下のメソッドは、DynamicFrame データレイクフレームワーク テーブルの読み取りと書き込みをサポートしていません。代わりに、DataFrame または Spark DataFrame API GlueContext のメソッドを使用してください。
  - Lake Formation の権限コントロールでは、GlueContext DynamicFrame 以下のメソッドはサポートされていません。
    - `create_dynamic_frame.from_catalog`
    - `write_dynamic_frame.from_catalog`
    - `getDynamicFrame`
    - `writeDynamicFrame`
  - Lake Formation の権限コントロールでは、GlueContext DataFrame 以下のメソッドがサポートされています。
    - `create_data_frame.from_catalog`
    - `write_data_frame.from_catalog`
    - `getDataFrame`
    - `writeDataFrame`
- [小さなファイルのグループ化](#)は、サポートされません。
- [ジョブのブックマーク](#)はサポートされません。
- AWS Glue 3.0 用 Apache Hudi 0.10.1 は Hudi マージオンリー (MoR) テーブルをサポートしていません。
- `ALTER TABLE ... RENAME TO` Apache Iceberg 0.13.1 for 3.0 では使用できません。AWS Glue

## Lake Formation の許可によって管理されるデータレイク形式のテーブルの制限事項

データレイクフォーマットは、Lake Formation AWS Glue の権限を通じてETLと統合されています。DynamicFrame create\_dynamic\_frameユーザー作成はサポートされていません。詳細については、次の例を参照してください。

- [例: Lake Formation の許可のコントロールを使用した Iceberg テーブルの読み取りおよび書き込み](#)
- [例: Lake Formation の許可のコントロールを使用した Hudi テーブルの読み取りおよび書き込み](#)
- [例: Lake Formation の許可のコントロールを使用した Delta Lake テーブルの読み取りおよび書き込み](#)

### Note

Apache Hudi、Apache Iceberg、Delta AWS Glue LakeのLake Formation 権限によるETLとの統合は、バージョン4.0でのみサポートされています。AWS Glue

Apache Icebergは、Lake Formation AWS Glue の権限を通じてETLと最もよく統合されています。ほぼすべてのオペレーションをサポートしており、SQL サポートも含まれています。

Hudi は、管理オペレーションを除くほとんどの基本オペレーションをサポートします。これは、これらのオプションが通常、データフレームの書き込みを介して実行され、additional\_options を介して指定されるためです。SparkSQL はサポートされていないため、DataFrames オペレーションの作成には AWS Glue API を使用する必要があります。

Delta Lake は、テーブルデータの読み取り、付加、および上書きのみをサポートします。Delta Lake では、更新などのさまざまなタスクを実行できるように独自のライブラリを使用する必要があります。

次の機能は、Lake Formation の許可によって管理される Iceberg テーブルでは使用できません。

- ETL を使用したコンパクション AWS Glue
- ETL 経由のスパーク SQL サポート AWS Glue

Lake Formation の許可によって管理される Hudi テーブルの制限は次のとおりです。

- 孤立したファイルの削除

Lake Formation の許可によって管理される Delta Lake テーブルの制限は次のとおりです。

- Delta Lake テーブルの挿入と読み取り以外のすべての機能。

## AWS Glue での Hudi フレームワークの使用

AWS Glue 3.0 以降では、データレイク向けに Apache Hudi フレームワークが利用できます。Hudi は、増分データ処理とデータパイプラインの開発を簡素化する、オープンソースのデータレイク用ストレージフレームワークです。このトピックでは、AWS Glue 内でデータを Hudi テーブルに転送または保存する際に利用可能な、各機能について説明します。Hudi の詳細については、公式の [Apache Hudi ドキュメント](#) を参照してください。

AWS Glue により、Amazon S3 内にある Hudi テーブルの読み取りおよび書き込み操作を実行できます。あるいは、AWS Glue データカタログを使用して、Hudi テーブルを操作することも可能です。挿入、更新、および、すべての [Apache Spark オペレーション](#) を含む操作も、追加でサポートされています。

### Note

Apache Hudi 0.10.1 for AWS Glue 3.0 では、Hudi Merge on Read (MoR) テーブルはサポートされません。

次の表に、AWS Glue の各バージョンに含まれている、Hudi のバージョンを一覧で示します。

AWS Glue のバージョン	サポートされる Hudi バージョン
4.0	0.12.1
3.0	0.10.1

AWS Glue がサポートするデータレイクフレームワークの詳細については、「[AWS Glue ETL ジョブでのデータレイクフレームワークの使用](#)」を参照してください。

## Hudi の有効化

AWS Glue で Hudi を有効化するには、以下のタスクを実行します。

- `hudi` を `--datalake-formats` のジョブパラメータの値として指定します。詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。
- AWS Glue ジョブ用に、`--conf` という名前でキーを作成し、それに次の値を設定します。または、スクリプトで `SparkConf` を使用して、次の構成を設定することもできます。これらの設定は、Apache Spark が Hudi テーブルを適切に処理するために役立ちます。

```
spark.serializer=org.apache.spark.serializer.KryoSerializer --conf
spark.sql.hive.convertMetastoreParquet=false
```

- Hudi に関する Lake Formation の許可のサポートは、AWS Glue 4.0 のためにデフォルトで有効になっています。Lake Formation に登録された Hudi テーブルの読み取り/書き込みに追加の設定は必要ありません。登録された Hudi テーブルを読み取るには、AWS Glue ジョブの IAM ロールに `SELECT` 許可が必要です。登録された Hudi テーブルに書き込むには、AWS Glue ジョブの IAM ロールに `SUPER` 許可が必要です。Lake Formation の許可の管理の詳細については、「[Data Catalog リソースに対する許可の付与と取り消し](#)」を参照してください。

## 別の Hudi バージョンの使用

AWS Glue でサポートされないバージョンの Hudi を使用するには、`--extra-jars job` パラメータにより独自の Hudi JAR ファイルを指定します。`--datalake-formats` ジョブパラメータの値として、`hudi` は含めないでください。

例: Hudi テーブルを Amazon S3 に書き込み、そのテーブルを AWS Glue データカタログに登録する

このスクリプト例では、Hudi テーブルを Amazon S3 に書き込み、そのテーブルを AWS Glue データカタログに登録する方法を示します。この例では、テーブルに登録するために、Hudi [Hive Sync ツール](#)を使用します。

### Note

この例では、AWS Glue データカタログを Apache Spark Hive のメタストアとして使用するために、`--enable-glue-datacatalog` ジョブパラメータの設定が必要となります。詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。

## Python

```
Example: Create a Hudi table from a DataFrame
```

```
and register the table to Glue Data Catalog

additional_options={
 "hoodie.table.name": "<your_table_name>",
 "hoodie.datasource.write.storage.type": "COPY_ON_WRITE",
 "hoodie.datasource.write.operation": "upsert",
 "hoodie.datasource.write.recordkey.field": "<your_recordkey_field>",
 "hoodie.datasource.write.precombine.field": "<your_precombine_field>",
 "hoodie.datasource.write.partitionpath.field": "<your_partitionkey_field>",
 "hoodie.datasource.write.hive_style_partitioning": "true",
 "hoodie.datasource.hive_sync.enable": "true",
 "hoodie.datasource.hive_sync.database": "<your_database_name>",
 "hoodie.datasource.hive_sync.table": "<your_table_name>",
 "hoodie.datasource.hive_sync.partition_fields": "<your_partitionkey_field>",
 "hoodie.datasource.hive_sync.partition_extractor_class":
"org.apache.hudi.hive.MultiPartKeyValueExtractor",
 "hoodie.datasource.hive_sync.use_jdbc": "false",
 "hoodie.datasource.hive_sync.mode": "hms",
 "path": "s3://<s3Path/>"
}

dataFrame.write.format("hudi") \
 .options(**additional_options) \
 .mode("overwrite") \
 .save()
```

## Scala

```
// Example: Example: Create a Hudi table from a DataFrame
// and register the table to Glue Data Catalog

val additionalOptions = Map(
 "hoodie.table.name" -> "<your_table_name>",
 "hoodie.datasource.write.storage.type" -> "COPY_ON_WRITE",
 "hoodie.datasource.write.operation" -> "upsert",
 "hoodie.datasource.write.recordkey.field" -> "<your_recordkey_field>",
 "hoodie.datasource.write.precombine.field" -> "<your_precombine_field>",
 "hoodie.datasource.write.partitionpath.field" -> "<your_partitionkey_field>",
 "hoodie.datasource.write.hive_style_partitioning" -> "true",
 "hoodie.datasource.hive_sync.enable" -> "true",
 "hoodie.datasource.hive_sync.database" -> "<your_database_name>",
 "hoodie.datasource.hive_sync.table" -> "<your_table_name>",
 "hoodie.datasource.hive_sync.partition_fields" -> "<your_partitionkey_field>",
```

```
"hoodie.datasource.hive_sync.partition_extractor_class" ->
"org.apache.hudi.hive.MultiPartKeyValueExtractor",
"hoodie.datasource.hive_sync.use_jdbc" -> "false",
"hoodie.datasource.hive_sync.mode" -> "hms",
"path" -> "s3://<s3Path/>")

dataFrame.write.format("hudi")
 .options(additionalOptions)
 .mode("append")
 .save()
```

例: AWS Glue データカタログを使用した Amazon S3 からの Hudi テーブルの読み取り

この例では、「[例: Hudi テーブルを Amazon S3 に書き込み、そのテーブルを AWS Glue データカタログに登録する](#)」で作成した Hudi テーブルを Amazon S3 から読み取ります。

#### Note

この例では、AWS Glue データカタログを Apache Spark Hive のメタストアとして使用するために、`--enable-glue-datacatalog` ジョブパラメータの設定が必要となります。詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。

## Python

この例では、[GlueContext.create\\_data\\_frame.from\\_catalog\(\)](#) メソッドを使用します。

```
Example: Read a Hudi table from Glue Data Catalog

from awsglue.context import GlueContext
from pyspark.context import SparkContext

sc = SparkContext()
glueContext = GlueContext(sc)

dataFrame = glueContext.create_data_frame.from_catalog(
 database = "<your_database_name>",
 table_name = "<your_table_name>"
)
```

## Scala

この例では [getCatalogSource](#) メソッドを使用します。

```
// Example: Read a Hudi table from Glue Data Catalog

import com.amazonaws.services.glue.GlueContext
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 val dataframe = glueContext.getCatalogSource(
 database = "<your_database_name>",
 tableName = "<your_table_name>"
).getDataFrame()
 }
}
```

例: **DataFrame** を更新して、Amazon S3 内にある Hudi テーブルに挿入する

この例では、「[例: Hudi テーブルを Amazon S3 に書き込み、そのテーブルを AWS Glue データカタログに登録する](#)」で作成した Hudi テーブルに対し、AWS Glue データカタログを使用して DataFrame を挿入します。

### Note

この例では、AWS Glue データカタログを Apache Spark Hive のメタストアとして使用するために、`--enable-glue-datacatalog` ジョブパラメータの設定が必要となります。詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。

## Python

この例では、[GlueContext.write\\_data\\_frame\\_from\\_catalog\(\)](#) メソッドを使用します。

```
Example: Upsert a Hudi table from Glue Data Catalog

from awsglue.context import GlueContext
```

```
from pyspark.context import SparkContext

sc = SparkContext()
glueContext = GlueContext(sc)

glueContext.write_data_frame.from_catalog(
 frame = dataframe,
 database = "<your_database_name>",
 table_name = "<your_table_name>",
 additional_options={
 "hoodie.table.name": "<your_table_name>",
 "hoodie.datasource.write.storage.type": "COPY_ON_WRITE",
 "hoodie.datasource.write.operation": "upsert",
 "hoodie.datasource.write.recordkey.field": "<your_recordkey_field>",
 "hoodie.datasource.write.precombine.field": "<your_precombine_field>",
 "hoodie.datasource.write.partitionpath.field": "<your_partitionkey_field>",
 "hoodie.datasource.write.hive_style_partitioning": "true",
 "hoodie.datasource.hive_sync.enable": "true",
 "hoodie.datasource.hive_sync.database": "<your_database_name>",
 "hoodie.datasource.hive_sync.table": "<your_table_name>",
 "hoodie.datasource.hive_sync.partition_fields": "<your_partitionkey_field>",
 "hoodie.datasource.hive_sync.partition_extractor_class":
"org.apache.hudi.hive.MultiPartKeyValueExtractor",
 "hoodie.datasource.hive_sync.use_jdbc": "false",
 "hoodie.datasource.hive_sync.mode": "hms"
 }
)
```

## Scala

この例では [getCatalogSink](#) メソッドを使用します。

```
// Example: Upsert a Hudi table from Glue Data Catalog

import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.util.JsonOptions
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 glueContext.getCatalogSink("<your_database_name>", "<your_table_name>",
 additionalOptions = JsonOptions(Map(
```

```

"hoodie.table.name" -> "<your_table_name>",
"hoodie.datasource.write.storage.type" -> "COPY_ON_WRITE",
"hoodie.datasource.write.operation" -> "upsert",
"hoodie.datasource.write.recordkey.field" -> "<your_recordkey_field>",
"hoodie.datasource.write.precombine.field" -> "<your_precombine_field>",
"hoodie.datasource.write.partitionpath.field" ->
"<your_partitionkey_field>",
"hoodie.datasource.write.hive_style_partitioning" -> "true",
"hoodie.datasource.hive_sync.enable" -> "true",
"hoodie.datasource.hive_sync.database" -> "<your_database_name>",
"hoodie.datasource.hive_sync.table" -> "<your_table_name>",
"hoodie.datasource.hive_sync.partition_fields" ->
"<your_partitionkey_field>",
"hoodie.datasource.hive_sync.partition_extractor_class" ->
"org.apache.hudi.hive.MultiPartKeyValueExtractor",
"hoodie.datasource.hive_sync.use_jdbc" -> "false",
"hoodie.datasource.hive_sync.mode" -> "hms"
)))
.writeDataFrame(dataFrame, glueContext)
}
}

```

例: Spark を使用した Amazon S3 からの Hudi テーブルの読み取り

この例では、Spark DataFrame API を使用して Amazon S3 から Hudi テーブルを読み取ります。

Python

```

Example: Read a Hudi table from S3 using a Spark DataFrame
dataFrame = spark.read.format("hudi").load("s3://<s3path/>")

```

Scala

```

// Example: Read a Hudi table from S3 using a Spark DataFrame
val dataFrame = spark.read.format("hudi").load("s3://<s3path/>")

```

例: スパークを使用した Amazon S3 への Hudi テーブルの書き込み

この例では、Spark を使用して Amazon S3 に Hudi テーブルを書き込みます。

## Python

```
Example: Write a Hudi table to S3 using a Spark DataFrame

dataFrame.write.format("hudi") \
 .options(**additional_options) \
 .mode("overwrite") \
 .save("s3://<s3Path/>")
```

## Scala

```
// Example: Write a Hudi table to S3 using a Spark DataFrame

dataFrame.write.format("hudi")
 .options(additionalOptions)
 .mode("overwrite")
 .save("s3://<s3path/>")
```

例: Lake Formation の許可のコントロールを使用した Hudi テーブルの読み取りおよび書き込み

この例では、Lake Formation の許可のコントロールを使用して Hudi テーブルの読み取りと書き込みを行います。

1. Hudi テーブルを作成して Lake Formation に登録します。

- Lake Formation の許可のコントロールを有効にするには、まずテーブルの Amazon S3 パスを Lake Formation に登録する必要があります。詳細については、「[Amazon S3 ロケーションの登録](#)」を参照してください。Lake Formation コンソールから、または AWS CLI を使用して登録できます。

```
aws lakeformation register-resource --resource-arn arn:aws:s3:::<s3-bucket>/<s3-
folder> --use-service-linked-role --region <REGION>
```

Amazon S3 の場所が登録されると、その場所 (またはその子である場所) をポイントするすべての AWS Glue テーブルが、GetTable 呼び出しで IsRegisteredWithLakeFormation パラメータの値を true として返します。

- Spark データフレーム API を介して登録された Amazon S3 パスをポイントする Hudi テーブルを作成します。

```
hudi_options = {
```

```

'hoodie.table.name': table_name,
'hoodie.datasource.write.storage.type': 'COPY_ON_WRITE',
'hoodie.datasource.write.recordkey.field': 'product_id',
'hoodie.datasource.write.table.name': table_name,
'hoodie.datasource.write.operation': 'upsert',
'hoodie.datasource.write.precombine.field': 'updated_at',
'hoodie.datasource.write.hive_style_partitioning': 'true',
'hoodie.upsert.shuffle.parallelism': 2,
'hoodie.insert.shuffle.parallelism': 2,
'path': <S3_TABLE_LOCATION>,
'hoodie.datasource.hive_sync.enable': 'true',
'hoodie.datasource.hive_sync.database': database_name,
'hoodie.datasource.hive_sync.table': table_name,
'hoodie.datasource.hive_sync.use_jdbc': 'false',
'hoodie.datasource.hive_sync.mode': 'hms'
}

df_products.write.format("hudi") \
 .options(**hudi_options) \
 .mode("overwrite") \
 .save()

```

2. AWS Glue ジョブ IAM ロールに Lake Formation の許可を付与します。Lake Formation コンソールから、または AWS CLI を使用して許可を付与できます。詳細については、「[Lake Formation コンソールと名前付きリソース方式を使用したテーブル許可の付与](#)」を参照してください。
3. Lake Formation に登録されている Hudi テーブルを読み取ります。このコードは、未登録の Hudi テーブルを読み取る場合と同じです。読み取りを成功させるには、AWS Glue ジョブの IAM ロールに SELECT 許可が必要であることに留意してください。

```

val dataframe = glueContext.getCatalogSource(
 database = "<your_database_name>",
 tableName = "<your_table_name>"
).getDataFrame()

```

4. Lake Formation に登録されている Hudi テーブルに書き込みます。このコードは、未登録の Hudi テーブルに書き込む場合と同じです。書き込みを成功させるには、AWS Glue ジョブの IAM ロールに SUPER 許可が必要であることに留意してください。

```

glueContext.getCatalogSink("<your_database_name>", "<your_table_name>",
 additionalOptions = JsonOptions(Map(
 "hoodie.table.name" -> "<your_table_name>",
 "hoodie.datasource.write.storage.type" -> "COPY_ON_WRITE",

```

```

"hoodie.datasource.write.operation" -> "<write_operation>",
"hoodie.datasource.write.recordkey.field" -> "<your_recordkey_field>",
"hoodie.datasource.write.precombine.field" -> "<your_precombine_field>",
"hoodie.datasource.write.partitionpath.field" -> "<your_partitionkey_field>",
"hoodie.datasource.write.hive_style_partitioning" -> "true",
"hoodie.datasource.hive_sync.enable" -> "true",
"hoodie.datasource.hive_sync.database" -> "<your_database_name>",
"hoodie.datasource.hive_sync.table" -> "<your_table_name>",
"hoodie.datasource.hive_sync.partition_fields" ->
"<your_partitionkey_field>",
"hoodie.datasource.hive_sync.partition_extractor_class" ->
"org.apache.hudi.hive.MultiPartKeyValueExtractor",
"hoodie.datasource.hive_sync.use_jdbc" -> "false",
"hoodie.datasource.hive_sync.mode" -> "hms"
)))
.writeDataFrame(dataFrame, glueContext)

```

## AWS Glue での Delta Lake フレームワークの使用

AWS Glue 3.0 以降では、Linux Foundation Delta Lake フレームワークを利用できます。Delta Lake は、ACID トランザクションの実行、メタデータ処理のスケーリング、さらにストリーミングとバッチデータ処理の統合を支援する、オープンソースのデータレイクストレージフレームワークです。このトピックでは、Delta Lake テーブルに転送または保存するデータに対して、AWS Glue 内で利用可能な機能について説明します。Delta Lake の詳細については、公式の [Delta Lake のドキュメント](#) を参照してください。

Amazon S3 内にある Delta Lake テーブルに対する読み取りおよび書き込みの操作は、AWS Glue を使用して実行できます。あるいは、AWS Glue データカタログにより Delta Lake テーブルを操作することもできます。挿入、更新、および [テーブルに対する一括での読み取りと書き込み](#) など、その他の操作もサポートされています。Delta Lake テーブルを使用する場合、Delta Lake Python ライブラリに備わったメソッド (例えば、DeltaTable.forPath) も使用可能です。Delta Lake Python ライブラリの詳細については、「Delta Lake の Python ドキュメント」を参照してください。

次の表に、AWS Glue の各バージョンに含まれている Delta Lake のバージョンを一覧で示します。

AWS Glue のバージョン	サポートされている Delta Lake バージョン
4.0	2.1.0
3.0	1.0.0

AWS Glue がサポートするデータレイクフレームワークの詳細については、「[AWS Glue ETL ジョブでのデータレイクフレームワークの使用](#)」を参照してください。

## AWS Glue のための Delta Lake の有効化

AWS Glue で Delta Lake を有効化するには、以下のタスクを実行します。

- `delta` を `--datalake-formats` のジョブパラメータの値として指定します。詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。
- AWS Glue ジョブ用に、`--conf` という名前でキーを作成し、それに次の値を設定します。または、スクリプトで `SparkConf` を使用して、次の構成を設定することもできます。これらの設定は、Apache Spark が Delta Lake テーブルを適切に処理するために役立ちます。

```
spark.sql.extensions=io.delta.sql.DeltaSparkSessionExtension --conf
spark.sql.catalog.spark_catalog=org.apache.spark.sql.delta.catalog.DeltaCatalog --conf
spark.delta.logStore.class=org.apache.spark.sql.delta.storage.S3SingleDriverLogStore
```

- Delta テーブルに関する Lake Formation の許可のサポートは、AWS Glue 4.0 のためにデフォルトで有効になっています。Lake Formation に登録された Delta テーブルの読み取り/書き込みに追加の設定は必要ありません。登録された Delta テーブルを読み取るには、AWS Glue ジョブの IAM ロールに `SELECT` 許可が必要です。登録された Delta テーブルに書き込むには、AWS Glue ジョブの IAM ロールに `SUPER` 許可が必要です。Lake Formation の許可の管理の詳細については、「[Data Catalog リソースに対する許可の付与と取り消し](#)」を参照してください。

## 別のバージョンの Delta Lake を使用する

AWS Glue でサポートされないバージョンの Delta Lake を使用するには、`--extra-jars` ジョブパラメータにより、独自の Delta Lake JAR ファイルを指定します。`--datalake-formats` ジョブパラメータの値として、`delta` は含めないでください。この場合、Delta Lake Python ライブラリを使用するには、`--extra-py-files` ジョブパラメータにより、ライブラリの JAR ファイルを指定する必要があります。Python ライブラリは、Delta Lake JAR ファイルにパッケージ化されています。

例: Delta Lake テーブルを Amazon S3 に書き込んで、そのテーブルを AWS Glue データカタログに登録する

次の AWS Glue ETL スクリプトは、Delta Lake テーブルを Amazon S3 に書き込み、そのテーブルを AWS Glue データカタログに登録する方法を示しています。

## Python

```
Example: Create a Delta Lake table from a DataFrame
and register the table to Glue Data Catalog

additional_options = {
 "path": "s3://<s3Path>"
}
dataFrame.write \
 .format("delta") \
 .options(**additional_options) \
 .mode("append") \
 .partitionBy("<your_partitionkey_field>") \
 .saveAsTable("<your_database_name>.<your_table_name>")
```

## Scala

```
// Example: Example: Create a Delta Lake table from a DataFrame
// and register the table to Glue Data Catalog

val additional_options = Map(
 "path" -> "s3://<s3Path>"
)
dataFrame.write.format("delta")
 .options(additional_options)
 .mode("append")
 .partitionBy("<your_partitionkey_field>")
 .saveAsTable("<your_database_name>.<your_table_name>")
```

例: AWS Glue データカタログを使用して Amazon S3 から Delta Lake テーブルを読み取る

次の AWS Glue ETL スクリプトでは、「[例: Delta Lake テーブルを Amazon S3 に書き込んで、そのテーブルを AWS Glue データカタログに登録する](#)」で作成した Delta Lake テーブルの読み取りを実行します。

## Python

この例では [create\\_data\\_frame\\_from\\_catalog](#) メソッドを使用します。

```
Example: Read a Delta Lake table from Glue Data Catalog
```

```
from awsglue.context import GlueContext
from pyspark.context import SparkContext

sc = SparkContext()
glueContext = GlueContext(sc)

df = glueContext.create_data_frame.from_catalog(
 database="<your_database_name>",
 table_name="<your_table_name>",
 additional_options=additional_options
)
```

## Scala

この例では [getCatalogSource](#) メソッドを使用します。

```
// Example: Read a Delta Lake table from Glue Data Catalog

import com.amazonaws.services.glue.GlueContext
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 val df = glueContext.getCatalogSource("<your_database_name>",
 "<your_table_name>",
 additionalOptions = additionalOptions)
 .getDataFrame()
 }
}
```

例: AWS Glue データカタログを使用して Amazon S3 内にある Delta Lake テーブルに **DataFrame** を挿入する

この例では、「[例: Delta Lake テーブルを Amazon S3 に書き込んで、そのテーブルを AWS Glue データカタログに登録する](#)」で作成した Delta Lake テーブルにデータを挿入します。

**Note**

この例では、AWS Glue データカタログを Apache Spark Hive のメタストアとして使用するために、`--enable-glue-datacatalog` ジョブパラメータの設定が必要となります。詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。

**Python**

この例では、[write\\_data\\_frame.from\\_catalog](#) メソッドを使用します。

```
Example: Insert into a Delta Lake table in S3 using Glue Data Catalog

from awsglue.context import GlueContext
from pyspark.context import SparkContext

sc = SparkContext()
glueContext = GlueContext(sc)

glueContext.write_data_frame.from_catalog(
 frame=dataFrame,
 database="<your_database_name>",
 table_name="<your_table_name>",
 additional_options=additional_options
)
```

**Scala**

この例では [getCatalogSink](#) メソッドを使用します。

```
// Example: Insert into a Delta Lake table in S3 using Glue Data Catalog

import com.amazonaws.services.glue.GlueContext
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 glueContext.getCatalogSink("<your_database_name>", "<your_table_name>",
 additionalOptions = additionalOptions)
 .writeDataFrame(dataFrame, glueContext)
 }
}
```

```
}
}
```

例: Spark API を使用して Amazon S3 から Delta Lake テーブルを読み取る

この例では、Spark API を使用して Amazon S3 から Delta Lake を読み取ります。

Python

```
Example: Read a Delta Lake table from S3 using a Spark DataFrame

dataFrame = spark.read.format("delta").load("s3://<s3path/>")
```

Scala

```
// Example: Read a Delta Lake table from S3 using a Spark DataFrame

val dataFrame = spark.read.format("delta").load("s3://<s3path/>")
```

例: Spark を使用した Amazon S3 への Delta Lake テーブルの書き込み

この例では、Spark を使用して Amazon S3 に Delta Lake テーブルを書き込みます。

Python

```
Example: Write a Delta Lake table to S3 using a Spark DataFrame

dataFrame.write.format("delta") \
 .options(**additional_options) \
 .mode("overwrite") \
 .partitionBy("<your_partitionkey_field>") \
 .save("s3://<s3Path>")
```

Scala

```
// Example: Write a Delta Lake table to S3 using a Spark DataFrame

dataFrame.write.format("delta") \
 .options(additionalOptions) \
 .mode("overwrite")
```

```
.partitionBy("<your_partitionkey_field>")
.save("s3://<s3path/>")
```

例: Lake Formation の許可のコントロールを使用した Delta Lake テーブルの読み取りおよび書き込み

この例では、Lake Formation の許可のコントロールを使用して Delta Lake テーブルの読み取りと書き込みを行います。

## 1. Delta テーブルを作成して Lake Formation に登録する

- a. Lake Formation の許可のコントロールを有効にするには、まずテーブルの Amazon S3 パスを Lake Formation に登録する必要があります。詳細については、「[Amazon S3 ロケーションの登録](#)」を参照してください。Lake Formation コンソールから、または AWS CLI を使用して登録できます。

```
aws lakeformation register-resource --resource-arn arn:aws:s3:::<s3-bucket>/<s3-
folder> --use-service-linked-role --region <REGION>
```

Amazon S3 の場所が登録されると、その場所 (またはその子である場所) をポイントするすべての AWS Glue テーブルが、GetTable 呼び出しで IsRegisteredWithLakeFormation パラメータの値を true として返します。

- b. Spark を介して登録された Amazon S3 パスをポイントする Delta テーブルを作成します。

### Note

次に Python の例を示します。

```
dataFrame.write \
.format("delta") \
.mode("overwrite") \
.partitionBy("<your_partitionkey_field>") \
.save("s3://<the_s3_path>")
```

データが Amazon S3 に書き込まれた後、AWS Glue クローラーを使用して新しい Delta カタログテーブルを作成します。詳細については、「[AWS Glue クローラーによるネイティブデータレイクテーブルサポートの紹介](#)」を参照してください。

AWS Glue CreateTable API を通じてテーブルを手動で作成することもできます。

2. AWS Glue ジョブ IAM ロールに Lake Formation の許可を付与します。Lake Formation コンソールから、または AWS CLI を使用して許可を付与できます。詳細については、「[Lake Formation コンソールと名前付きリソース方式を使用したテーブル許可の付与](#)」を参照してください。
3. Lake Formation に登録されている Delta テーブルを読み取ります。このコードは、未登録の Delta テーブルを読み取る場合と同じです。読み取りを成功させるには、AWS Glue ジョブの IAM ロールに SELECT 許可が必要であることを留意してください。

```
Example: Read a Delta Lake table from Glue Data Catalog

df = glueContext.create_data_frame.from_catalog(
 database="<your_database_name>",
 table_name="<your_table_name>",
 additional_options=additional_options
)
```

4. Lake Formation に登録されている Delta テーブルに書き込みます。このコードは、未登録の Delta テーブルに書き込む場合と同じです。書き込みを成功させるには、AWS Glue ジョブの IAM ロールに SUPER 許可が必要であることを留意してください。

デフォルトでは、AWS Glue は saveMode として Append を使用します。これを変更するには、additional\_options で saveMode オプションを設定します。Delta テーブルでの saveMode サポートの詳細については、「[テーブルに書き込む](#)」を参照してください。

```
glueContext.write_data_frame.from_catalog(
 frame=dataFrame,
 database="<your_database_name>",
 table_name="<your_table_name>",
 additional_options=additional_options
)
```

## AWS Glue での Iceberg フレームワークの使用

AWS Glue 3.0 以降では、データレイク向けに Apache Iceberg フレームワークがサポートされています。Iceberg により、SQL テーブルと同様な機能を持つ、高性能なテーブルフォーマットが提供されます。このトピックでは、AWS Glue 内でデータを Iceberg テーブルに転送または保存する際に利用可能な、機能について説明します。Iceberg の詳細については、公式の[Apache Iceberg ドキュメント](#)を参照してください。

Amazon S3 内の Iceberg テーブルに対する読み取りと書き込みの操作は、AWS Glue を使用して実行することができます。あるいは、AWS Glue データカタログを使用して Iceberg テーブルを操作することも可能です。挿入、更新に加え、[Spark Queries](#) や [Spark Writes](#) のすべての操作も、追加でサポートされています。

### Note

Apache Iceberg 0.13.1 for AWS Glue 3.0 では、ALTER TABLE ... RENAME TO はサポートされません。

次の表に、AWS Glue の各バージョンに含まれる Iceberg のバージョンを一覧で示します。

AWS Glue のバージョン	サポートされる Iceberg バージョン
4.0	1.0.0
3.0	0.13.1

AWS Glue がサポートするデータレイクフレームワークの詳細については、「[AWS Glue ETL ジョブでのデータレイクフレームワークの使用](#)」を参照してください。

### Iceberg フレームワークの有効化

AWS Glue で Iceberg を有効化するには、以下のタスクを実行します。

- iceberg を --datalake-formats のジョブパラメータの値として指定します。詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。
- AWS Glue ジョブ用に、--conf という名前でキーを作成し、それに次の値を設定します。または、スクリプトで SparkConf を使用して、次の構成を設定することもできます。これらの設定は、Apache Spark が Iceberg テーブルを適切に処理する際に役立ちます。

```
spark.sql.extensions=org.apache.iceberg.spark.extensions.IcebergSparkSessionExtensions
--conf spark.sql.catalog.glue_catalog=org.apache.iceberg.spark.SparkCatalog
--conf spark.sql.catalog.glue_catalog.warehouse=s3://<your-warehouse-dir>/
--conf spark.sql.catalog.glue_catalog.catalog-impl=org.apache.iceberg.aws.glue.GlueCatalog
--conf spark.sql.catalog.glue_catalog.io-impl=org.apache.iceberg.aws.s3.S3FileIO
```

Lake Formation に登録されている Iceberg テーブルの読み取りまたは書き込みを行う場合は、次の設定を追加して Lake Formation のサポートを有効にします。Lake Formation に登録された Iceberg テーブルをサポートしているのは AWS Glue 4.0 のみであることに留意してください。

```
--conf spark.sql.catalog.glue_catalog.glue.lakeformation-enabled=true
--conf spark.sql.catalog.glue_catalog.glue.id=<table-catalog-id>
```

AWS Glue 3.0 で Iceberg 0.13.1 を使用する場合には、Amazon DynamoDB ロックマネージャーを使用して確実なアトミックトランザクションを行うために、次の追加事項を設定する必要があります。AWS Glue 4.0 では、デフォルトで楽観的ロックが使用されます。詳細については、公式の Apache Iceberg ドキュメントで「[Iceberg と AWS の統合](#)」を参照してください。

```
--conf spark.sql.catalog.glue_catalog.lock-
impl=org.apache.iceberg.aws.glue.DynamoLockManager
--conf spark.sql.catalog.glue_catalog.lock.table=<your-dynamodb-table-name>
```

## 別バージョンの Iceberg の使用

AWS Glue でサポートされないバージョンの Iceberg を使用するには、`--extra-jars` ジョブパラメータを使用して独自の Iceberg JAR ファイルを指定します。`--datalake-formats` パラメータの値として、`iceberg` を含めないようにしてください。

## Iceberg テーブルの暗号化を有効にする

### Note

Iceberg テーブルにはサーバー側の暗号化を有効にする独自のメカニズムがあります。AWS Glue のセキュリティ設定に加えて、これらの設定を有効にする必要があります。

Iceberg テーブルでサーバー側の暗号化を有効にするには、「[Iceberg ドキュメント](#)」のガイダンスを確認してください。

例: Amazon S3 に Iceberg テーブルを書き込み、そのテーブルを AWS Glue データカタログに登録する

このスクリプト例は、Amazon S3 に Iceberg テーブルを書き込む方法を示しています。この例では、[Iceberg と AWS の統合](#)を使用して、テーブルを AWS Glue データカタログに登録しています。

## Python

```
Example: Create an Iceberg table from a DataFrame
and register the table to Glue Data Catalog

dataFrame.createOrReplaceTempView("tmp_<your_table_name>")

query = f"""
CREATE TABLE glue_catalog.<your_database_name>.<your_table_name>
USING iceberg
TBLPROPERTIES ("format-version"="2")
AS SELECT * FROM tmp_<your_table_name>
"""
spark.sql(query)
```

## Scala

```
// Example: Example: Create an Iceberg table from a DataFrame
// and register the table to Glue Data Catalog

dataFrame.createOrReplaceTempView("tmp_<your_table_name>")

val query = """CREATE TABLE glue_catalog.<your_database_name>.<your_table_name>
USING iceberg
TBLPROPERTIES ("format-version"="2")
AS SELECT * FROM tmp_<your_table_name>
"""
spark.sql(query)
```

または、Spark メソッドを使用して Amazon S3 とデータカタログに Iceberg テーブルを書き込むこともできます。

前提条件:Iceberg ライブラリが使用するカタログを用意する必要があります。AWS Glue データカタログを使う場合、AWS Glue を使えば簡単です。AWS Glue データカタログは、`glue_catalog` として Spark ライブラリで使用できるように事前に設定されています。データカタログテーブルは、`atabaseName` と `tableName` で識別されます。AWS Glue データカタログの詳細については、「[データ検出とカタログ化](#)」を参照してください

AWS Glue データカタログを使用していない場合は、Spark API を使用してカタログをプロビジョニングする必要があります。詳細については、Spark ドキュメントの「[Spark の設定](#)」を参照してください。

この例では、Spark を使用して Amazon S3 と データカタログに Iceberg テーブルを書き込みます。

## Python

```
Example: Write an Iceberg table to S3 on the Glue Data Catalog

Create (equivalent to CREATE TABLE AS SELECT)
dataFrame.writeTo("glue_catalog.databaseName.tableName") \
 .tableProperty("format-version", "2") \
 .create()

Append (equivalent to INSERT INTO)
dataFrame.writeTo("glue_catalog.databaseName.tableName") \
 .tableProperty("format-version", "2") \
 .append()
```

## Scala

```
// Example: Write an Iceberg table to S3 on the Glue Data Catalog

// Create (equivalent to CREATE TABLE AS SELECT)
dataFrame.writeTo("glue_catalog.databaseName.tableName")
 .tableProperty("format-version", "2")
 .create()

// Append (equivalent to INSERT INTO)
dataFrame.writeTo("glue_catalog.databaseName.tableName")
 .tableProperty("format-version", "2")
 .append()
```

例: AWS Glue データカタログを使用した Amazon S3 からの Iceberg テーブルの読み込み

この例では、「[例: Amazon S3 に Iceberg テーブルを書き込み、そのテーブルを AWS Glue データカタログに登録する](#)」で作成した Iceberg テーブルを読み取っています。

## Python

この例では、[GlueContext.create\\_data\\_frame\\_from\\_catalog\(\)](#) メソッドを使用します。

```
Example: Read an Iceberg table from Glue Data Catalog
```

```
from aws glue.context import GlueContext
from pyspark.context import SparkContext

sc = SparkContext()
glueContext = GlueContext(sc)

df = glueContext.create_data_frame.from_catalog(
 database="<your_database_name>",
 table_name="<your_table_name>",
 additional_options=additional_options
)
```

## Scala

この例では [getCatalogSource](#) メソッドを使用します。

```
// Example: Read an Iceberg table from Glue Data Catalog

import com.amazonaws.services.glue.GlueContext
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 val df = glueContext.getCatalogSource("<your_database_name>",
"<your_table_name>",
 additionalOptions = additionalOptions)
 .getDataFrame()
 }
}
```

例: AWS Glue データカタログを使用した Amazon S3 にある Iceberg テーブルへの **DataFrame** の挿入

この例では、「[例: Amazon S3 に Iceberg テーブルを書き込み、そのテーブルを AWS Glue データカタログに登録する](#)」で作成した Iceberg テーブルにデータを挿入します。

**Note**

この例では、AWS Glue データカタログを Apache Spark Hive のメタストアとして使用するために、`--enable-glue-datacatalog` ジョブパラメータの設定が必要となります。詳細については、「[AWS Glue ジョブのパラメータ](#)」を参照してください。

## Python

この例では、[GlueContext.write\\_data\\_frame.from\\_catalog\(\)](#) メソッドを使用します。

```
Example: Insert into an Iceberg table from Glue Data Catalog

from awsglue.context import GlueContext
from pyspark.context import SparkContext

sc = SparkContext()
glueContext = GlueContext(sc)

glueContext.write_data_frame.from_catalog(
 frame=dataFrame,
 database="<your_database_name>",
 table_name="<your_table_name>",
 additional_options=additional_options
)
```

## Scala

この例では [getCatalogSink](#) メソッドを使用します。

```
// Example: Insert into an Iceberg table from Glue Data Catalog

import com.amazonaws.services.glue.GlueContext
import org.apache.spark.SparkContext

object GlueApp {
 def main(sysArgs: Array[String]): Unit = {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 glueContext.getCatalogSink("<your_database_name>", "<your_table_name>",
 additionalOptions = additionalOptions)
 .writeDataFrame(dataFrame, glueContext)
 }
}
```

```
}
}
```

例: Spark を使用した Iceberg テーブルの Amazon S3 からの読み取り

前提条件:Iceberg ライブラリが使用するカタログを用意する必要があります。AWS Glue データカタログを使う場合、AWS Glue を使えば簡単です。AWS Glue データカタログは、`glue_catalog` として Spark ライブラリで使用できるように事前に設定されています。データカタログテーブルは、`atabaseName` と `tableName` で識別されます。AWS Glue データカタログの詳細については、「[データ検出とカタログ化](#)」を参照してください

AWS Glue データカタログを使用していない場合は、Spark API を使用してカタログをプロビジョニングする必要があります。詳細については、Spark ドキュメントの「[Spark の設定](#)」を参照してください。

この例では、Spark を使用してデータカタログから Amazon S3 の Iceberg テーブルを読み取ります。

Python

```
Example: Read an Iceberg table on S3 as a DataFrame from the Glue Data Catalog

dataFrame = spark.read.format("iceberg").load("glue_catalog.databaseName.tableName")
```

Scala

```
// Example: Read an Iceberg table on S3 as a DataFrame from the Glue Data Catalog

val dataFrame =
 spark.read.format("iceberg").load("glue_catalog.databaseName.tableName")
```

例: Lake Formation の許可のコントロールを使用した Iceberg テーブルの読み取りおよび書き込み

この例では、Lake Formation の許可のコントロールを使用して Iceberg テーブルの読み取りと書き込みを行います。

1. Iceberg テーブルを作成して Lake Formation に登録します:

- a. Lake Formation の許可のコントロールを有効にするには、まずテーブルの Amazon S3 パスを Lake Formation に登録する必要があります。詳細については、「[Amazon S3 ロケーションの](#)

[登録](#)」を参照してください。Lake Formation コンソールから、または AWS CLI を使用して登録できます。

```
aws lakeformation register-resource --resource-arn arn:aws:s3:::<s3-bucket>/<s3-folder> --use-service-linked-role --region <REGION>
```

Amazon S3 の場所が登録されると、その場所 (またはその子である場所) をポイントするすべての AWS Glue テーブルが、GetTable 呼び出しで IsRegisteredWithLakeFormation パラメータの値を true として返します。

- b. Spark SQL を介して登録されたパスをポイントする Iceberg テーブルを作成します。

 Note

次に Python の例を示します。

```
dataFrame.createOrReplaceTempView("tmp_<your_table_name>")

query = f"""
CREATE TABLE glue_catalog.<your_database_name>.<your_table_name>
USING iceberg
AS SELECT * FROM tmp_<your_table_name>
"""
spark.sql(query)
```

AWS Glue CreateTable API を通じてテーブルを手動で作成することもできます。詳細については、「[Apache Iceberg テーブルの作成](#)」を参照してください。

2. ジョブ IAM ロールに Lake Formation の許可を付与します。Lake Formation コンソールから、または AWS CLI を使用して許可を付与できます。詳細については、次を参照してください: <https://docs.aws.amazon.com/lake-formation/latest/dg/granting-table-permissions.html>
3. Lake Formation に登録されている Iceberg テーブルを読み取ります。このコードは、未登録の Iceberg テーブルを読み取る場合と同じです。読み取りを成功させるには、AWS Glue ジョブの IAM ロールに SELECT 許可が必要であることを留意してください。

```
Example: Read an Iceberg table from the AWS Glue Data Catalog
from awsglue.context import GlueContext
from pyspark.context import SparkContext

sc = SparkContext()
```

```
glueContext = GlueContext(sc)

df = glueContext.create_data_frame.from_catalog(
 database="<your_database_name>",
 table_name="<your_table_name>",
 additional_options=additional_options
)
```

4. Lake Formation に登録されている Iceberg テーブルに書き込みます。このコードは、未登録の Iceberg テーブルに書き込む場合と同じです。書き込みを成功させるには、AWS Glue ジョブの IAM ロールに SUPER 許可が必要であることを留意してください。

```
glueContext.write_data_frame.from_catalog(
 frame=dataFrame,
 database="<your_database_name>",
 table_name="<your_table_name>",
 additional_options=additional_options
)
```

## 共有設定リファレンス

形式タイプには次の `format_options` 値を使用できます。

- `attachFilename` — 列名として使用する適切な形式の文字列。このオプションを指定すると、レコードのソースファイルの名前がレコードに追加されます。パラメータ値は列名として使用されません。
- `attachTimestamp` — 列名として使用する適切な形式の文字列。このオプションを指定すると、レコードのソースファイルの変更時刻がレコードに追加されます。パラメータ値は列名として使用されます。

## Spark SQL ジョブ用の AWS Glue Data Catalog のサポート

AWS Glue の Data Catalog は、Apache Hive メタストア互換のカタログです。Data Catalog を外部の Apache Hive メタストアとして使用するよう、ジョブと開発エンドポイント AWS Glue を設定できます。その後、Data Catalog に格納されたテーブルに対して Apache Spark SQL クエリを直接実行することができます。AWS Glue ダイナミックフレームは、デフォルトで Data Catalog と統合されています。また、この機能により Spark SQL ジョブは、Data Catalog を外部の Hive メタストアとして使用できるようになります。

この機能には、AWS Glue API エンドポイントへのネットワークアクセスが必要です。ネットワークアクセスを提供する場合、プライベートサブネットに接続がある AWS Glue ジョブでは VPC エンドポイントまたは NAT ゲートウェイを構成する必要があります。VPC エンドポイントの設定については、「[データストアへのネットワークアクセスを設定する](#)」を参照してください。NAT ゲートウェイを作成するには、Amazon VPC ユーザーガイドの「[NAT ゲートウェイ](#)」を参照してください。

ジョブ引数と開発エンドポイント引数にそれぞれ "--enable-glue-datacatalog": "" 引数を追加することで、AWS Glue ジョブと開発エンドポイントを設定できます。この引数を渡すことで、Spark が外部 Hive メタストアとして Data Catalog にアクセスできるようにする、特定の設定が構成されます。同時に、SparkSession オブジェクトにある AWS Glue ジョブまたは開発エンドポイント内で作成された [Hive のサポートを有効化](#)します。

データ カタログへのアクセスを有効にするには、コンソールの [ ジョブの追加 ] または [ エンドポイントの追加 ] ページのカタログオプショングループで [ Hive メタストアとして AWS Glue データ カタログを使用 ] のチェックボックスをオンにします。ジョブまたは開発エンドポイントに使用される IAM ロールには glue:CreateDatabase アクセス許可が必要です。存在しない場合、「default」というデータベースが Data Catalog に作成されます。

Spark SQL ジョブでこの機能の使用法の例を見てみましょう。次の例では、s3://awsglue-datasets/examples/us-legislators にある米国国会議員のデータセットをクローリングしたと想定しています。

AWS Glue Data Catalog で定義されたテーブルからデータをシリアル化/逆シリアル化するには、Spark ジョブのクラスパスにある AWS Glue Data Catalog で定義された形式の [Hive SerDe](#) クラスが、Spark SQL に必要です。

特定の一般的なフォーマットの SerDes は AWS Glue によって配布されています。以下はこれらへの Amazon S3 リンクです。

- [JSON](#)
- [XML](#)
- [Grok](#)

JSON SerDe を [追加の JAR として開発エンドポイント](#) に追加します。ジョブでは、引数フィールドで --extra-jars 引数を使用して SerDe を追加できます。詳細については、[を参照してください](#)。 [AWS Glue ジョブのパラメータ](#)

ここに、Spark SQL 用に Data Catalog を有効にして開発エンドポイントを作成するための入力 JSON の例を示します。

```
{
 "EndpointName": "Name",
 "RoleArn": "role_ARN",
 "PublicKey": "public_key_contents",
 "NumberOfNodes": 2,
 "Arguments": {
 "--enable-glue-datacatalog": ""
 },
 "ExtraJarsS3Path": "s3://crawler-public/json/serde/json-serde.jar"
}
```

Spark SQL を使用して、米国国会議員のデータセットから作成されたテーブルにクエリを実行します。

```
>>> spark.sql("use legislators")
DataFrame[]
>>> spark.sql("show tables").show()
+-----+-----+-----+
| database| tableName|isTemporary|
+-----+-----+-----+
|legislators| areas_json| false|
|legislators| countries_json| false|
|legislators| events_json| false|
|legislators| memberships_json| false|
|legislators| organizations_json| false|
|legislators| persons_json| false|
+-----+-----+-----+
>>> spark.sql("describe memberships_json").show()
+-----+-----+-----+
| col_name|data_type| comment|
+-----+-----+-----+
| area_id| string|from deserializer|
| on_behalf_of_id| string|from deserializer|
| organization_id| string|from deserializer|
| role| string|from deserializer|
| person_id| string|from deserializer|
|legislative_perio...| string|from deserializer|
| start_date| string|from deserializer|
| end_date| string|from deserializer|
+-----+-----+-----+
```

形式の SerDe クラスがジョブのクラスパスで使用できない場合は、以下に示すようなエラーが表示されます。

```
>>> spark.sql("describe memberships_json").show()

Caused by: MetaException(message:java.lang.ClassNotFoundException Class
org.openx.data.jsonserde.JsonSerDe not found)
 at
org.apache.hadoop.hive.metastore.MetaStoreUtils.getDeserializer(MetaStoreUtils.java:399)
 at
org.apache.hadoop.hive.ql.metadata.Table.getDeserializerFromMetaStore(Table.java:276)
 ... 64 more
```

memberships テーブルとは異なる organization\_id だけを表示するには、次の SQL クエリを実行します。

```
>>> spark.sql("select distinct organization_id from memberships_json").show()
+-----+
| organization_id|
+-----+
|d56acebe-8fdc-47b...|
|8fa6c3d2-71dc-478...|
+-----+
```

動的フレームで同じ手順を実行する必要がある場合は、以下を実行してください。

```
>>> memberships = glueContext.create_dynamic_frame.from_catalog(database="legislators",
table_name="memberships_json")
>>> memberships.toDF().createOrReplaceTempView("memberships")
>>> spark.sql("select distinct organization_id from memberships").show()
+-----+
| organization_id|
+-----+
|d56acebe-8fdc-47b...|
|8fa6c3d2-71dc-478...|
+-----+
```

DynamicFrame は ETL オペレーション用に最適化されていますが、Spark SQL が Data Catalog に直接アクセスできるようにすることで、複雑な SQL ステートメントの実行や、既存のアプリケーションの移植が簡素化されます。

## ジョブのブックマークを使用する

Spark 用の AWS Glue はジョブのブックマークを使用して処理済みのデータを追跡します。ジョブのブックマーク機能の概要とサポート内容については、「[the section called “ジョブのブックマークを使用した処理済みデータの追跡”](#)」を参照してください。ブックマークを使用して AWS Glue ジョブをプログラミングすると、ビジュアルジョブでは実現できない柔軟性が得られます。

- JDBC から読み込む場合、AWS Glue スクリプトのブックマークキーとして使用する列を指定できます。
- 各メソッド呼び出しに適用する `transformation_ctx` を選択できます。

常に適切に設定されたパラメータでスクリプトの先頭に `job.init` を呼び出して、末尾に `job.commit` を呼び出します。これら 2 つの関数は、ブックマークサービスを初期化し、サービスの状態変化を更新します。ブックマークは呼び出さないと機能しません。

### ブックマークキーの指定

JDBC ワークフローの場合、ブックマークはキーフィールドの値をブックマークされた値と比較することで、ジョブがどの行を読み込んだかを追跡します。これは Amazon S3 ワークフローには必要ありませんし、適用可能でもありません。ビジュアルエディタを使用せずに AWS Glue スクリプトを記述する場合、どの列をブックマークで追跡するかを指定できます。複数の列を指定することもできます。ユーザー定義のブックマークキーを指定する場合、値の順序にギャップがあってもかまいません。

#### Warning

ユーザー定義のブックマークキーを使用する場合、キーはそれぞれ厳密に一定間隔で増減する必要があります。複合キーに追加のフィールドを選択しても、「マイナーバージョン」や「リビジョン番号」などの概念のフィールドは、値がデータセット全体で再利用されるため、この基準を満たしません。

`jobBookmarkKeys` および `jobBookmarkKeysSortOrder` は以下の方法で指定できます。

- `create_dynamic_frame.from_catalog` — `additional_options` を使用する。
- `create_dynamic_frame.from_options` — `connection_options` を使用する。

## 変換コンテキスト

AWS Glue PySpark の動的フレームの多くのメソッドには、`transformation_ctx` というオプションのパラメータが含まれています。このパラメータは ETL 演算子インスタンスの一意の識別子です。`transformation_ctx` パラメータは指定された演算子に対するジョブのブックマーク内の状態情報を識別するために使用されます。具体的には、AWS Glue では `transformation_ctx` を使用してブックマーク状態に対するキーにインデックスを付けます。

### Warning

`transformation_ctx` は、スクリプト内の特定のソースのブックマーク状態を検索するためのキーとして機能します。ブックマークを正しく機能させるためには、ソースと関連する `transformation_ctx` の一貫性を常に保つ必要があります。ソースプロパティを変更したり、`transformation_ctx` の名前を変更したりすると、前のブックマークが無効になり、タイムスタンプベースのフィルタリングで正しい結果が得られない場合があります。

ジョブのブックマークが正しく機能するように、ジョブのブックマークのパラメータを有効にし、`transformation_ctx` パラメータを設定します。`transformation_ctx` パラメータを渡さない場合、メソッドで使用されている動的フレームやテーブルに対してジョブのブックマークは有効になりません。例えば、ETL ジョブで 2 つの Amazon S3 ソースを読み取って結合する場合、ブックマークを有効にするメソッドに対してのみ `transformation_ctx` パラメータを渡すことができます。1 つのジョブについてジョブのブックマークをリセットした場合、`transformation_ctx` の使用には関係なく、このジョブに関連付けられているすべての変換がリセットされます。

`DynamicFrameReader` クラスの詳細については、「[DynamicFrameReader クラス](#)」を参照してください。PySpark 拡張の詳細については、「[AWS Glue PySpark 拡張機能リファレンス](#)」を参照してください

### 例

#### Example

Amazon S3 データソース用に生成されたスクリプトの例を次に示します。ジョブのブックマークを使用するために必要なスクリプトの部分は、斜体で表示されています。これらの要素の詳細については、[GlueContext クラス](#) API および [DynamicFrameWriter クラス](#) API を参照してください。

```
Sample Script
import sys
```

```
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job

args = getResolvedOptions(sys.argv, ['JOB_NAME'])
sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'], args)

datasource0 = glueContext.create_dynamic_frame.from_catalog(
 database = "database",
 table_name = "relatedqueries_csv",
 transformation_ctx = "datasource0"
)

applymapping1 = ApplyMapping.apply(
 frame = datasource0,
 mappings = [("col0", "string", "name", "string"), ("col1", "string", "number",
"string")],
 transformation_ctx = "applymapping1"
)

datasink2 = glueContext.write_dynamic_frame.from_options(
 frame = applymapping1,
 connection_type = "s3",
 connection_options = {"path": "s3://input_path"},
 format = "json",
 transformation_ctx = "datasink2"
)

job.commit()
```

## Example

JDBC ソース用に生成されたスクリプトの例を以下に示します。ソーステーブルは、empno 列がプライマリキーである従業員テーブルです。デフォルトでは、ブックマークキーが指定されていない場合、ジョブはブックマークキーとしてシーケンシャルプライマリキーを使用しますが、empno は必ずしもシーケンシャルではなく、値にギャップがある可能性があるため、デフォルトのブックマーク

キーとして適切ではありません。したがって、スクリプトは `empno` を明示的にブックマークキーとして指定します。コードのその部分は、斜体で表示されます。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job

args = getResolvedOptions(sys.argv, ['JOB_NAME'])

sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args['JOB_NAME'], args)

datasource0 = glueContext.create_dynamic_frame.from_catalog(
 database = "hr",
 table_name = "emp",
 transformation_ctx = "datasource0",
 additional_options = {"jobBookmarkKeys":["empno"],"jobBookmarkKeysSortOrder":"asc"}
)

applymapping1 = ApplyMapping.apply(
 frame = datasource0,
 mappings = [("ename", "string", "ename", "string"), ("hrly_rate", "decimal(38,0)",
"hrly_rate", "decimal(38,0)"), ("comm", "decimal(7,2)", "comm", "decimal(7,2)"),
("hiredate", "timestamp", "hiredate", "timestamp"), ("empno", "decimal(5,0)", "empno",
"decimal(5,0)"), ("mgr", "decimal(5,0)", "mgr", "decimal(5,0)"), ("photo", "string",
"photo", "string"), ("job", "string", "job", "string"), ("deptno", "decimal(3,0)",
"deptno", "decimal(3,0)"), ("ssn", "decimal(9,0)", "ssn", "decimal(9,0)"), ("sal",
"decimal(7,2)", "sal", "decimal(7,2)"]],
 transformation_ctx = "applymapping1"
)

datasink2 = glueContext.write_dynamic_frame.from_options(
 frame = applymapping1,
 connection_type = "s3",
 connection_options = {"path": "s3://hr/employees"},
 format = "csv",
 transformation_ctx = "datasink2"
```

```
)

job.commit()
```

## AWS Glue Studio 外でのセンシティブデータ検出の使用

AWS Glue Studio では機密データを検出できますが、AWS Glue Studio の外部で機密データ検出機能を使用することもできます。

マネージド機密データタイプの全リストについては、「[Managed data types](#)」を参照してください。

### AWS 管理 PII タイプを使用した機密データ検出の検出

AWS Glue 1 つの AWS Glue ETL ジョブに 2 つの API を提供します。これらは、`detect()` および `classifyColumns()` です。

```
detect(frame: DynamicFrame,
 entityTypesToDetect: Seq[String],
 outputColumnName: String = "DetectedEntities",
 detectionSensitivity: String = "LOW"): DynamicFrame

detect(frame: DynamicFrame,
 detectionParameters: JsonOptions,
 outputColumnName: String = "DetectedEntities",
 detectionSensitivity: String = "LOW"): DynamicFrame

classifyColumns(frame: DynamicFrame,
 entityTypesToDetect: Seq[String],
 sampleFraction: Double = 0.1,
 thresholdFraction: Double = 0.1,
 detectionSensitivity: String = "LOW")
```

`detect()` API を使用して、AWS 管理対象 PII タイプとカスタムエンティティタイプを識別できます。検出結果が含まれる新しい列が自動的に作成されます。`classifyColumns()` API は、キーが列名で値が検出されたエンティティタイプのリストになっているマップを返します。SampleFraction は、PII エンティティをスキャンしているときにサンプリングするためのデータの割合を示し、ThresholdFraction は、列が PII データとして識別されるために満たされる必要があるデータの割合を示します。

## 行レベルの検出

この例では、ジョブが `detect()` と `classifyColumns()` API を使用して以下のアクションを実行しています。

- Amazon S3 バケットからデータを読み込んで `DynamicFrame` に変換します。
- `dynamicFrame` 内にある「Email」と「Credit Card」のインスタンスを検出する
- 元の値と、各行の検出結果が網羅された 1 つの列が含まれる `dynamicFrame` を返す
- 返された `DynamicFrame` を別のパスに書き込みます。Amazon S3

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.MappingSpec
import com.amazonaws.services.glue.errors.CallSite
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._
import com.amazonaws.services.glue.ml.EntityDetector

object GlueApp {
 def main(sysArgs: Array[String]) {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)
 val frame=
glueContext.getSourceWithFormat(formatOptions=JsonOptions("""{"quoteChar": "\"",
"withHeader": true, "separator": ","}"""), connectionType="s3", format="csv",
options=JsonOptions("""{"paths": ["s3://pathToSource"], "recurse": true}"""),
transformationContext="AmazonS3_node1650160158526").getDynamicFrame()

 val frameWithDetectedPII = EntityDetector.detect(frame, Seq("EMAIL",
"CREDIT_CARD"))

 glueContext.getSinkWithFormat(connectionType="s3",
options=JsonOptions("""{"path": "s3://pathToOutput/", "partitionKeys": []}"""),
transformationContext="someCtx",
format="json").writeDynamicFrame(frameWithDetectedPII)
```

```
 Job.commit()
 }
}
```

## 詳細なアクションによる行レベルの検出

この例では、ジョブが `detect()` API を使用して次のアクションを実行しています。

- Amazon S3 バケットからデータを読み取り、`dynamicFrame` に変換する
- `dynamicFrame` 内の「USA\_PTIN」、「BANK\_ACCOUNT」、「USA\_SSN」、「USA\_PASSPORT\_NUMBER」、および「PHONE\_NUMBER」の機密データタイプの検出
- 変更およびマスキングされた値と、各行の検出結果が網羅された 1 つの列が含まれる `dynamicFrame` を返す
- 返された `dynamicFrame` を別の Amazon S3 パスに書き込む

上記の `detect()` API とは対照的に、これはエンティティタイプを検出するための詳細なアクションを使用します。詳細については、「[detect\(\) を使用する検出パラメータ](#)」を参照してください。

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.MappingSpec
import com.amazonaws.services.glue.errors.CallSite
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._
import com.amazonaws.services.glue.ml.EntityDetector

object GlueApp {
 def main(sysArgs: Array[String]) {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)
 val frame =
glueContext.getSourceWithFormat(formatOptions=JsonOptions("""{"quoteChar": "\"",
"withHeader": true, "separator": ","}"""), connectionType="s3", format="csv",
options=JsonOptions("""{"paths": ["s3://pathToSource"], "recurse": true}"""),
transformationContext="AmazonS3_node_source").getDynamicFrame()
```

```
val detectionParameters = JsonOptions(
 """
 {
 "USA_DRIVING_LICENSE": [{
 "action": "PARTIAL_REDACT",
 "sourceColumns": ["Driving License"],
 "actionOptions": {
 "matchPattern": "[0-9]",
 "redactChar": "*"
 }
 }
]],
 "BANK_ACCOUNT": [{
 "action": "DETECT",
 "sourceColumns": ["*"]
 }],
 "USA_SSN": [{
 "action": "SHA256_HASH",
 "sourceColumns": ["SSN"]
 }],
 "IP_ADDRESS": [{
 "action": "REDACT",
 "sourceColumns": ["IP Address"],
 "actionOptions": {"redactText": "*****"}
 }],
 "PHONE_NUMBER": [{
 "action": "PARTIAL_REDACT",
 "sourceColumns": ["Phone Number"],
 "actionOptions": {
 "numLeftCharsToExclude": 1,
 "numRightCharsToExclude": 0,
 "redactChar": "*"
 }
 }
]
}
 """
)

val frameWithDetectedPII = EntityDetector.detect(frame, detectionParameters,
"DetectedEntities", "HIGH")

glueContext.getSinkWithFormat(connectionType="s3", options=JsonOptions("""{"path":
"s3://pathToOutput/", "partitionKeys": []}""")),
```

```
transformationContext="AmazonS3_node_target",
format="json").writeDynamicFrame(frameWithDetectedPII)

 Job.commit()
}
}
```

## 列レベルの検出

この例では、ジョブが `classifyColumns()` API を使用して次のアクションを実行しています。

- Amazon S3 バケットからデータを読み取り、`dynamicFrame` に変換する
- `dynamicFrame` 内にある「Email」と「Credit Card」のインスタンスを検出する
- 列の 100% をサンプリングし、エンティティがセルの 10% にある場合に検出済みとしてマークするとともに、感度が [低] になるようにパラメータを設定します
- キーが列名、値が検出されたエンティティタイプのリストであるマップを返します
- 返された `dynamicFrame` を別の Amazon S3 パスに書き込む

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.MappingSpec
import com.amazonaws.services.glue.errors.CallSite
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._
import com.amazonaws.services.glue.DynamicFrame
import com.amazonaws.services.glue.ml.EntityDetector

object GlueApp {
 def main(sysArgs: Array[String]) {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)
 val frame =
glueContext.getSourceWithFormat(formatOptions=JsonOptions("""{"quoteChar":
"\\"", "withHeader": true, "separator": ",", "optimizePerformance": false}"""),
```

```
connectionType="s3", format="csv", options=JsonOptions("""{"paths": ["s3://
pathToSource"], "recurse": true}"""), transformationContext="frame").getDynamicFrame()

import glueContext.sparkSession.implicits._

val detectedDataFrame = EntityDetector.classifyColumns(
 frame,
 entityTypeToDetect = Seq("CREDIT_CARD", "PHONE_NUMBER"),
 sampleFraction = 1.0,
 thresholdFraction = 0.1,
 detectionSensitivity = "LOW"
)
val detectedDF = (detectedDataFrame).toSeq.toDF("columnName", "entityTypes")
val DetectSensitiveData_node = DynamicFrame(detectedDF, glueContext)

glueContext.getSinkWithFormat(connectionType="s3", options=JsonOptions("""{"path":
"s3://pathToOutput", "partitionKeys": []}"""), transformationContext="someCtx",
format="json").writeDynamicFrame(DetectSensitiveData_node)

Job.commit()
}
```

## 機密データの検出、PII タイプを使用した検出 AWS CustomEntityType

Studio を使用してカスタムエンティティを定義できます。AWS ただし、AWS Studio 以外でこの機能を使用するには、まずカスタムエンティティタイプを定義し、定義したカスタムエンティティタイプをのリストに追加する必要がありますentityTypesToDetect。

データ内に特定の機密データタイプ (「Employee Id」など) がある場合は、CreateCustomEntityType() API を呼び出してカスタムエンティティを作成することができます。以下の例では、リクエストパラメータを使用して CreateCustomEntityType() API にカスタムエンティティタイプ「EMPLOYEE\_ID」を定義します。

```
{
 "name": "EMPLOYEE_ID",
 "regexString": "\\d{4}-\\d{3}",
 "contextWords": ["employee"]
}
```

次に、カスタムエンティティタイプ (EMPLOYEE\_ID) をEntityDetector() API に追加することで、新しいカスタム機密データタイプを使用するようにジョブを変更します。

```
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.MappingSpec
import com.amazonaws.services.glue.errors.CallSite
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._
import com.amazonaws.services.glue.ml.EntityDetector

object GlueApp {
 def main(sysArgs: Array[String]) {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)
 val frame=
glueContext.getSourceWithFormat(formatOptions=JsonOptions("""{"quoteChar": "\"",
"withHeader": true, "separator": ","}"""), connectionType="s3", format="csv",
options=JsonOptions("""{"paths": ["s3://pathToSource"], "recurse": true}"""),
transformationContext="AmazonS3_node1650160158526").getDynamicFrame()

 val frameWithDetectedPII = EntityDetector.detect(frame, Seq("EMAIL",
"CREDIT_CARD", "EMPLOYEE_ID"))

 glueContext.getSinkWithFormat(connectionType="s3",
options=JsonOptions("""{"path": "s3://pathToOutput/", "partitionKeys": []}"""),
transformationContext="someCtx",
format="json").writeDynamicFrame(frameWithDetectedPII)

 Job.commit()
 }
}
```

**Note**

カスタム機密データタイプが既存マネージドエンティティタイプと同じ名前で定義されている場合は、カスタム機密データタイプが優先され、マネージドエンティティタイプのロジックが上書きされます。

**detect()** を使用する検出パラメータ

このメソッドは、内のエンティティを検出するために使用されます。DynamicFrame DataFrame 元の値を含む新しい列と、PII outputColumnName 検出メタデータを含む追加の列を返します。DynamicFrame カスタムマスキングは、AWS Glueスクリプト内で返された後で行うことも、代わりに detect () によるきめ細かいアクション API を使用することもできます。

```
detect(frame: DynamicFrame,
 entityTypeToDetect: Seq[String],
 outputColumnName: String = "DetectedEntities",
 detectionSensitivity: String = "LOW"): DynamicFrame
```

## パラメータ:

- frame — (type:DynamicFrame) DynamicFrame 処理するデータを含む入力。
- entityTypeToDetect — (type:[Seq[String]) 検出するエンティティタイプのリスト。マネージドエンティティタイプまたはカスタムエンティティタイプのいずれかになります。
- outputColumnName— (type: String、デフォルト:"DetectedEntities ") 検出されたエンティティが保存される列の名前。指定しない場合、デフォルトの列名は "DetectedEntities" です。
- detectionSensitivity – (タイプ: String、オプション: [低] または [高]、デフォルト: [低]) 検出プロセスの感度を指定します。有効なオプションは [低] または [高] です。指定しない場合、デフォルトの感度は [低] に設定されます。

## outputColumnName の設定:

検出されたエンティティが保存される列の名前。指定しない場合、デフォルトの列名は "DetectedEntities" です。出力列の各行について、補足列には、次の key-value ペアを含む、検出されたエンティティメタデータに対する列名のマップが含まれます。

- entityType – 検出されたエンティティタイプ。

- start – 元のデータ内の検出されたエンティティの開始場所。
- end – 元のデータ内の検出されたエンティティの終了場所。
- actionUsed – 検出されたエンティティに対して実行されたアクション (例: 「DETECT」、 「REDACT」、 「PARTIAL\_REDACT」、 「SHA256\_HASH」)。

例 :

```
{
 "DetectedEntities":{
 "SSN Col":[
 {
 "entityType":"USA_SSN",
 "actionUsed":"DETECT",
 "start":4,
 "end":15
 }
],
 "Random Data col":[
 {
 "entityType":"BANK_ACCOUNT",
 "actionUsed":"PARTIAL_REDACT",
 "start":4,
 "end":13
 },
 {
 "entityType":"IP_ADDRESS",
 "actionUsed":"REDACT",
 "start":4,
 "end":13
 }
]
 }
}
```

### 詳細なアクションの **detect()** の検出パラメータ

このメソッドは、DynamicFrame 指定されたパラメータを使用してエンティティを検出するために使用されます。DataFrame 元の値をマスクされた機密データに置き換えた新しい列と、PII outputColumnName 検出メタデータを含む追加の列を返します。

```
detect(frame: DynamicFrame,
 detectionParameters: JsonOptions,
 outputColumnName: String = "DetectedEntities",
 detectionSensitivity: String = "LOW"): DynamicFrame
```

## パラメータ:

- `frame` — (type:DynamicFrame): DynamicFrame 処理するデータを含む入力。
- `detectionParameters` – (タイプ: JsonOptions): 検出プロセスのパラメータを指定する JSON オプション。
- `outputColumnName`— (type: String、デフォルト:"DetectedEntities「」): 検出されたエンティティが保存される列の名前。指定しない場合、デフォルトの列名は "DetectedEntities" です。
- `detectionSensitivity` – (タイプ: String、オプション: [低] または [高]、デフォルト: [低]): 検出プロセスの感度を指定します。有効なオプションは [低] または [高] です。指定しない場合、デフォルトの感度は [低] に設定されます。

## detectionParameters の設定

設定が含まれていない場合は、デフォルト値が使用されます。

- `action` – (タイプ: String、オプション: 「DETECT」、「REDACT」、「PARTIAL\_REDACT」、「SHA256\_HASH」) エンティティに対して実行するアクションを指定します。必須。マスキングを実行するアクション (「DETECT」を除くすべて) は、列ごとに 1 つのアクションしか実行できないことに留意してください。これは、合体したエンティティをマスキングするための予防措置です。
- `sourceColumns` – (タイプ: List[String]、デフォルト: ["\*"]) エンティティの検出を実行するソース列名のリスト。存在しない場合は、デフォルトで ["\*"] に設定されます。無効な列名が使用された場合に `IllegalArgumentException` を発生させます。
- `sourceColumnsTo除外` — (タイプ: List[String]) エンティティの検出対象となるソース列名のリスト。 `sourceColumns` または `sourceColumnsToExclude` のいずれかを使用します。無効な列名が使用された場合に `IllegalArgumentException` を発生させます。
- `actionOptions` – 指定されたアクションに基づく追加オプション:
  - 「DETECT」および「SHA256\_HASH」の場合、いかなるオプションも使用できません。
  - 「REDACT」の場合:

- `redactText` – (タイプ: `String`、デフォルト: 「\*\*\*\*\*」) 検出されたエンティティを置き換えるテキスト。
- 「PARTIAL\_REDACT」の場合:
  - `redactChar` – (タイプ: `String`、デフォルト: 「\*」) エンティティ内で検出された各文字を置き換える文字。
  - `matchPattern` – (タイプ: `String`) 部分的なマスキングの正規表現パターン。 `numLeftCharsToExclude` またはと組み合わせることはできません `numRightCharsToExclude`。
  - `numLeftCharsToExclude` — (タイプ: `String`, `integer`) 除外する残りの文字の数。 `matchPattern` と組み合わせることはできませんが、 `numRightCharsToExclude` と組み合わせて使用することはできます。
  - `numRightCharsToExclude` — (タイプ: `String`, `integer`) 除外する右側の文字の数。 `matchPattern` と組み合わせることはできませんが、 `numRightCharsToExclude` と組み合わせて使用することはできます。

## outputColumnName の設定

[「outputColumnName 設定」を参照してください。](#)

## classifyColumns() の検出パラメータ

このメソッドは、内のエンティティを検出するために使用されます。 `DynamicFrame` キーが列名、値が検出されたエンティティタイプのリストであるマップを返します。これが AWS Glue スクリプト内で返された後、カスタムマスキングを実行できます。

```
classifyColumns(frame: DynamicFrame,
 entityTypeToDetect: Seq[String],
 sampleFraction: Double = 0.1,
 thresholdFraction: Double = 0.1,
 detectionSensitivity: String = "LOW")
```

### パラメータ:

- `frame` — (type: `DynamicFrame`) `DynamicFrame` 処理するデータを含む入力。
- `entityTypesToDetect` — (type: `Seq[String]`) 検出するエンティティタイプのリスト。マネージドエンティティタイプまたはカスタムエンティティタイプのいずれかになります。

- `sampleFraction` – (タイプ: `Double`、デフォルト: 10%) PII エンティティをスキャンする際にサンプリングするデータの割合。
- `thresholdFraction` – (タイプ: `Double`、デフォルト: 10%): 列が PII データとして識別されるために満たす必要があるデータの割合。
- `detectionSensitivity` – (タイプ: `String`、オプション: [低] または [高]、デフォルト: [低]) 検出プロセスの感度を指定します。有効なオプションは [低] または [高] です。指定しない場合、デフォルトの感度は [低] に設定されます。

## 管理対象の機密データタイプ

### グローバルエンティティ

データタイプ	カテゴリ	説明
PERSON_NAME	Universal	その人物の名前
EMAIL	個人	E メールアドレス
IP_ADDRESS	[コンピュータ]	IP アドレス
MAC_ADDRESS	個人	MAC アドレス

### 米国のデータタイプ

データタイプ	説明
BANK_ACCOUNT	銀行口座番号。国や地域に固有のものではありませんが、米国とカナダの口座形式のみが検出されます。
CREDIT_CARD	クレジットカード番号。
PHONE_NUMBER	電話番号。国や地域に固有のものではありませんが、現時点では米国とカナダの電話番号のみが検出されます。

データタイプ	説明
USA_ATIN	内国歳入庁 (IRS、Internal Revenue Service) が発行する米国養子縁組納税者番号 (ATIN、Adoption Taxpayer Identification Number)。
USA_CPT_CODE	CPT コード (米国のみ)。
USA_DEA_NUMBER	DEA 番号 (米国のみ)。
USA_DRIVING_LICENSE	運転免許証番号 (米国のみ)。
USA_HCPCS_CODE	HCPCS コード (米国のみ)。
USA_HEALTH_INSURANCE_CLAIM_NUMBER	健康保険請求番号 (米国のみ)。
USA_ITIN	ITIN (米国人または法人の場合)。
USA_MEDICARE_BENEFICIARY_IDENTIFIER	メディケア受益者識別子 (米国のみ)。
USA_NATIONAL_DRUG_CODE	NDC コード (米国のみ)。
USA_NATIONAL_PROVIDER_IDENTIFIER	米国医療従事者識別子番号 (米国のみ)。
USA_PASSPORT_NUMBER	パスポート番号 (米国人の場合)。
USA_PTIN	内国歳入庁が発行する米国申告書作成者番号 (PTIN、Preparer Tax Identification Number)。
USA_SSN	社会保障番号 (米国人の場合)。

### アルゼンチンのデータタイプ

データタイプ	説明
ARGENTINA_TAX_IDENTIFICATION_NUMBER	アルゼンチンの納税者番号。CUIT または CUIL と呼ばれます。

## オーストラリアのデータタイプ

データタイプ	説明
AUSTRALIA_BUSINESS_NUMBER	オーストラリア事業者番号 (ABN、Australia Business Number)。オーストラリア事業者登記所 (ABR、Australian Business Register) が発行する固有の識別番号で、政府や地域社会が事業者を識別します。
AUSTRALIA_COMPANY_NUMBER	オーストラリア企業番号 (ACN、Australia Company Number)。オーストラリア証券投資委員会 (ASIC、Australian Securities and Investments Commission) が発行する固有の識別番号。
AUSTRALIA_DRIVING_LICENSE	オーストラリアの運転免許証番号。
AUSTRALIA_MEDICARE_NUMBER	オーストラリアの健康保険番号。オーストラリア健康保険委員会 (Australian Health Insurance Commission) が発行する個人識別番号。
AUSTRALIA_PASSPORT_NUMBER	オーストラリアのパスポート番号。
AUSTRALIA_TAX_FILE_NUMBER	オーストラリアのタックスファイルナンバー (TFN、Tax File Number)。オーストラリア税務局 (ATO、Australian Taxation Office) が納税者 (個人、会社など) に税務取引を行うために発行します。

## オーストリアのデータタイプ

データタイプ	説明
AUSTRIA_DRIVING_LICENSE	運転免許証番号 (オーストリアのみ)。
AUSTRIA_PASSPORT_NUMBER	パスポート番号 (オーストリアのみ)。

データタイプ	説明
AUSTRIA_SSN	社会保障番号 (オーストリア人の場合)。
AUSTRIA_TAX_IDENTIFICATION_NUMBER	納税者番号 (オーストリアのみ)。
AUSTRIA_VALUE_ADDED_TAX	付加価値税 (オーストリアのみ)。

### バルカン諸国のデータタイプ

データタイプ	説明
BOSNIA_UNIQUE_MASTER_CITIZEN_NUMBER	ボスニア・ヘルツェゴビナ市民の固有マスター市民番号 (JMBG、Unique master citizen number)。
KOSOVO_UNIQUE_MASTER_CITIZEN_NUMBER	コソボの固有マスター市民番号 (JMBG)。
MACEDONIA_UNIQUE_MASTER_CITIZEN_NUMBER	マケドニアの固有マスター市民番号。
MONTENEGRO_UNIQUE_MASTER_CITIZEN_NUMBER	モンテネグロの固有マスター市民番号 (JMBG)。
SERBIA_UNIQUE_MASTER_CITIZEN_NUMBER	セルビアの固有マスター市民番号 (JMBG)。
SERBIA_VALUE_ADDED_TAX	付加価値税 (セルビアのみ)。
VOJVODINA_UNIQUE_MASTER_CITIZEN_NUMBER	ヴォイヴォディナの固有マスター市民番号 (JMBG)。

### ベルギーのデータタイプ

データタイプ	説明
BELGIUM_DRIVING_LICENSE	運転免許証番号 (ベルギーのみ)。
BELGIUM_NATIONAL_IDENTIFICATION_NUMBER	ベルギー国民番号 (BNN、Belgian National Number)。
BELGIUM_PASSPORT_NUMBER	パスポート番号 (ベルギーのみ)。
BELGIUM_TAX_IDENTIFICATION_NUMBER	納税者番号 (ベルギーのみ)。
BELGIUM_VALUE_ADDED_TAX	付加価値税 (ベルギーのみ)。

### ブラジルのデータタイプ

データタイプ	説明
BRAZIL_BANK_ACCOUNT	銀行口座番号 (ブラジルのみ)。
BRAZIL_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (ブラジルのみ)。
BRAZIL_NATIONAL_REGISTRY_OF_LEGAL_ENTITIES_NUMBER	企業に発行される識別番号 (ブラジルのみ)。CNPJとも呼ばれます。
BRAZIL_NATURAL_PERSON_REGISTRY_NUMBER	CPFとも呼ばれる自然人登録番号 (Natural Person Registry Number)。

### ブルガリアのデータタイプ

データタイプ	説明
BULGARIA_DRIVING_LICENSE	運転免許証番号 (ブルガリアのみ)。
BULGARIA_UNIFORM_CIVIL_NUMBER	国民識別番号として機能する統一市民番号 (EGN、Unified Civil Number)。
BULGARIA_VALUE_ADDED_TAX	付加価値税 (ブルガリアのみ)。

## カナダのデータタイプ

データタイプ	説明
CANADA_DRIVING_LICENSE	運転免許証番号 (カナダのみ)。
CANADA_GOVERNMENT_IDENTIFICATION_CARD_NUMBER	国民識別番号 (カナダのみ)。
CANADA_PASSPORT_NUMBER	パスポート番号 (カナダのみ)。
CANADA_PERMANENT_RESIDENCE_NUMBER	永住者番号 (PR カード番号)。
CANADA_PERSONAL_HEALTH_NUMBER	ヘルスケアの固有識別番号 (PHN 番号)。
CANADA_SOCIAL_INSURANCE_NUMBER	カナダの社会保険番号 (SIN)。

## チリのデータタイプ

データタイプ	説明
CHILE_DRIVING_LICENSE	運転免許証番号 (チリのみ)。
CHILE_NATIONAL_IDENTIFICATION_NUMBER	チリの国民識別番号。RUT または RUN と呼ばれます。

## 中国、香港、マカオ、台湾のデータタイプ

データタイプ	説明
CHINA_IDENTIFICATION	中国の識別子。
CHINA_LICENSE_PLATE_NUMBER	運転免許証番号 (中国のみ)。
CHINA_MAINLAND_TRAVEL_PERMIT_ID_HONG_KONG_MACAU	香港とマカオの居住者向けの本土旅行許可証。

データタイプ	説明
CHINA_MAINLAND_TRAVEL_PERMIT_ID_TAIWAN	中華人民共和国 (PRC) 政府が発行する台湾居住者向けの本土旅行許可証。
CHINA_PASSPORT_NUMBER	パスポート番号 (中国のみ)。
CHINA_PHONE_NUMBER	電話番号 (中国のみ)。
HONG_KONG_IDENTITY_CARD	香港入国管理局が発行する公式の身分証明書。
MACAU_RESIDENT_IDENTITY_CARD	マカオ居住者 ID カード (BIR) は、マカオの ID サービス局 (Identification Services Bureau) が発行する公式の ID カードです。
TAIWAN_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (台湾のみ)。
TAIWAN_PASSPORT_NUMBER	パスポート番号 (台湾のみ)。

### コロンビアのデータタイプ

データタイプ	説明
COLOMBIA_PERSONAL_IDENTIFICATION_NUMBER	出生時にコロンビア人に割り当てられる固有の識別番号。
COLOMBIA_TAX_IDENTIFICATION_NUMBER	納税者番号 (コロンビアのみ)。

### クロアチアのデータタイプ

データタイプ	説明
CROATIA_DRIVING_LICENSE	運転免許証番号 (クロアチアのみ)。
CROATIA_IDENTITY_NUMBER	国民識別番号 (クロアチアのみ)。
CROATIA_PASSPORT_NUMBER	パスポート番号 (クロアチアのみ)。

データタイプ	説明
CROATIA_PERSONAL_IDENTIFICATION_NUMBER	個人識別番号 (OIB)。

### キプロスのデータタイプ

データタイプ	説明
CYPRUS_DRIVING_LICENSE	運転免許証番号 (キプロスのみ)。
CYPRUS_NATIONAL_IDENTIFICATION_NUMBER	キプロスの身分証明書。
CYPRUS_PASSPORT_NUMBER	パスポート番号 (キプロスのみ)。
CYPRUS_TAX_IDENTIFICATION_NUMBER	納税者番号 (キプロスのみ)。
CYPRUS_VALUE_ADDED_TAX	付加価値税 (キプロスのみ)。

### チェコのデータタイプ

データタイプ	説明
CZECHIA_DRIVING_LICENSE	運転免許証番号 (チェコのみ)。
CZECHIA_PERSONAL_IDENTIFICATION_NUMBER	個人識別番号 (チェコのみ)。
CZECHIA_VALUE_ADDED_TAX	付加価値税 (チェコのみ)。

### デンマークのデータタイプ

データタイプ	説明
DENMARK_DRIVING_LICENSE	運転免許証番号 (デンマークのみ)。

データタイプ	説明
DENMARK_PERSONAL_IDENTIFICATION_NUMBER	個人識別番号 (デンマークのみ)。
DENMARK_TAX_IDENTIFICATION_NUMBER	納税者番号 (デンマークのみ)。
DENMARK_VALUE_ADDED_TAX	付加価値税 (デンマークのみ)。

### エストニアのデータタイプ

データタイプ	説明
ESTONIA_DRIVING_LICENSE	運転免許証番号 (エストニアのみ)。
ESTONIA_PASSPORT_NUMBER	パスポート番号 (エストニアのみ)。
ESTONIA_PERSONAL_IDENTIFICATION_CODE	個人識別番号 (エストニアのみ)。
ESTONIA_VALUE_ADDED_TAX	付加価値税 (エストニアのみ)。

### フィンランドのデータタイプ

データタイプ	説明
FINLAND_DRIVING_LICENSE	運転免許証番号 (フィンランドのみ)。
FINLAND_HEALTH_INSURANCE_NUMBER	健康保険番号 (フィンランドのみ)。
FINLAND_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (フィンランドのみ)。
FINLAND_PASSPORT_NUMBER	パスポート番号 (フィンランドのみ)。
FINLAND_VALUE_ADDED_TAX	付加価値税 (フィンランドのみ)。

## フランスのデータタイプ

データタイプ	説明
FRANCE_BANK_ACCOUNT	銀行口座番号 (フランスのみ)。
FRANCE_DRIVING_LICENSE	運転免許証番号 (フランスのみ)。
FRANCE_HEALTH_INSURANCE_NUMBER	フランスの健康保険番号。
FRANCE_INSEE_CODE	フランスの社会保障 (SSN または NIR) 番号。
FRANCE_NATIONAL_IDENTIFICATION_NUMBER	フランス国民識別番号 (CNI)。
FRANCE_PASSPORT_NUMBER	パスポート番号 (フランスのみ)。
FRANCE_TAX_IDENTIFICATION_NUMBER	納税者番号 (フランスのみ)。
FRANCE_VALUE_ADDED_TAX	付加価値税 (フランスのみ)。

## ドイツのデータタイプ

データタイプ	説明
GERMANY_BANK_ACCOUNT	銀行口座番号 (ドイツのみ)。
GERMANY_DRIVING_LICENSE	運転免許証番号 (ドイツのみ)。
GERMANY_PASSPORT_NUMBER	パスポート番号 (ドイツのみ)。
GERMANY_PERSONAL_IDENTIFICATION_NUMBER	個人識別番号 (ドイツのみ)。
GERMANY_TAX_IDENTIFICATION_NUMBER	納税者番号 (ドイツのみ)。
GERMANY_VALUE_ADDED_TAX	付加価値税 (ドイツのみ)。

## ギリシャのデータタイプ

データタイプ	説明
GREECE_DRIVING_LICENSE	運転免許証番号 (ギリシャのみ)。
GREECE_PASSPORT_NUMBER	パスポート番号 (ギリシャのみ)。
GREECE_SSN	社会保障番号 (ギリシャ人の場合)。
GREECE_TAX_IDENTIFICATION_NUMBER	納税者番号 (ギリシャのみ)。
GREECE_VALUE_ADDED_TAX	付加価値税 (ギリシャのみ)。

### ハンガリーのデータタイプ

データタイプ	説明
HUNGARY_DRIVING_LICENSE	運転免許証番号 (ハンガリーのみ)。
HUNGARY_PASSPORT_NUMBER	パスポート番号 (ハンガリーのみ)。
HUNGARY_SSN	社会保障番号 (ハンガリー人の場合)。
HUNGARY_TAX_IDENTIFICATION_NUMBER	納税者番号 (ハンガリーのみ)。
HUNGARY_VALUE_ADDED_TAX	付加価値税 (ハンガリーのみ)。

### アイスランドのデータタイプ

データタイプ	説明
ICELAND_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (アイスランドのみ)。
ICELAND_PASSPORT_NUMBER	パスポート番号 (アイスランドのみ)。
ICELAND_VALUE_ADDED_TAX	付加価値税 (アイスランドのみ)。

## インドのデータタイプ

データタイプ	説明
INDIA_AADHAAR_NUMBER	インドの固有 ID 局 (Unique Identification Authority) が発行する Aadhaar 識別番号。
INDIA_PERMANENT_ACCOUNT_NUMBER	インドの納税者番号 (PAN、Permanent Account Number)。

## インドネシアのデータタイプ

データタイプ	説明
INDONESIA_IDENTITY_CARD_NUMBER	国民識別番号 (インドネシアのみ)。

## アイルランドのデータタイプ

データタイプ	説明
IRELAND_DRIVING_LICENSE	運転免許証番号 (アイルランドのみ)。
IRELAND_PASSPORT_NUMBER	パスポート番号 (アイルランドのみ)。
IRELAND_PERSONAL_PUBLIC_SERVICE_NUMBER	アイルランド個人公共サービス (PPS、Personal Public Service) 番号。
IRELAND_TAX_IDENTIFICATION_NUMBER	納税者番号 (アイルランドのみ)。
IRELAND_VALUE_ADDED_TAX	付加価値税 (アイルランドのみ)。

## イスラエルのデータタイプ

データタイプ	説明
ISRAEL_IDENTIFICATION_NUMBER	国民識別番号 (イスラエルのみ)。

## イタリアのデータタイプ

データタイプ	説明
ITALY_BANK_ACCOUNT	銀行口座番号 (イタリアのみ)。
ITALY_DRIVING_LICENSE	運転免許証番号 (イタリアのみ)。
ITALY_FISCAL_CODE	識別番号。Italian Codice Fiscale (イタリアの納税者コード) としても知られています。
ITALY_PASSPORT_NUMBER	パスポート番号 (イタリアのみ)。
ITALY_VALUE_ADDED_TAX	付加価値税 (イタリアのみ)。

## 日本のデータタイプ

データタイプ	説明
JAPAN_BANK_ACCOUNT	日本の銀行口座。
JAPAN_DRIVING_LICENSE	日本の運転免許証番号。
JAPAN_MY_NUMBER	国税行政、年金行政、災害対応に使用される日本の国民または法人のための固有の識別子
JAPAN_PASSPORT_NUMBER	日本のパスポート番号。

## 韓国のデータタイプ

データタイプ	説明
KOREA_PASSPORT_NUMBER	パスポート番号 (韓国のみ)。
KOREA_RESIDENCE_REGISTRATION_NUMBER_FOR_CITIZENS	居住者の韓国居住登録者番号。

データタイプ	説明
KOREA_RESIDENCE_REGISTRATION_NUMBER_FOR_FOREIGNERS	外国人の韓国居住登録者番号。

### ラトビアのデータタイプ

データタイプ	説明
LATVIA_DRIVING_LICENSE	運転免許証番号 (ラトビアのみ)。
LATVIA_PASSPORT_NUMBER	パスポート番号 (ラトビアのみ)。
LATVIA_PERSONAL_IDENTIFICATION_NUMBER	個人識別番号 (ラトビアのみ)。
LATVIA_VALUE_ADDED_TAX	付加価値税 (ラトビアのみ)。

### リヒテンシュタインのデータタイプ

データタイプ	説明
LIECHTENSTEIN_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (リヒテンシュタインのみ)。
LIECHTENSTEIN_PASSPORT_NUMBER	パスポート番号 (リヒテンシュタインのみ)。
LIECHTENSTEIN_TAX_IDENTIFICATION_NUMBER	納税者番号 (リヒテンシュタインのみ)。

### リトアニアのデータタイプ

データタイプ	説明
LITHUANIA_DRIVING_LICENSE	運転免許証番号 (リトアニアのみ)。

データタイプ	説明
LITHUANIA_PERSONAL_IDENTIFICATION_NUMBER	個人識別番号 (リトアニアのみ)。
LITHUANIA_TAX_IDENTIFICATION_NUMBER	納税者番号 (リトアニアのみ)。
LITHUANIA_VALUE_ADDED_TAX	付加価値税 (リトアニアのみ)。

### ルクセンブルクのデータタイプ

データタイプ	説明
LUXEMBOURG_DRIVING_LICENSE	運転免許証番号 (ルクセンブルクのみ)。
LUXEMBOURG_NATIONAL_INDIVIDUAL_NUMBER	国民識別番号 (ルクセンブルクのみ)。
LUXEMBOURG_PASSPORT_NUMBER	パスポート番号 (ルクセンブルクのみ)。
LUXEMBOURG_TAX_IDENTIFICATION_NUMBER	納税者番号 (ルクセンブルクのみ)。
LUXEMBOURG_VALUE_ADDED_TAX	付加価値税 (ルクセンブルクのみ)。

### マレーシアのデータタイプ

データタイプ	説明
MALAYSIA_MYKAD_NUMBER	国民識別番号 (マレーシアのみ)。
MALAYSIA_PASSPORT_NUMBER	パスポート番号 (マレーシアのみ)。

### マルタのデータタイプ

データタイプ	説明
MALTA_DRIVING_LICENSE	運転免許証番号 (マルタのみ)。
MALTA_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (マルタのみ)。
MALTA_TAX_IDENTIFICATION_NUMBER	納税者番号 (マルタのみ)。
MALTA_VALUE_ADDED_TAX	付加価値税 (マルタのみ)。

### メキシコのデータタイプ

データタイプ	説明
MEXICO_CLABE_NUMBER	メキシコ CLABE (Clave Bancaria Estandarizada) 銀行番号。
MEXICO_DRIVING_LICENSE	運転免許証番号 (メキシコのみ)。
MEXICO_PASSPORT_NUMBER	パスポート番号 (メキシコのみ)。
MEXICO_TAX_IDENTIFICATION_NUMBER	納税者番号 (メキシコのみ)。
MEXICO_UNIQUE_POPULATION_REGISTRY_CODE	メキシコの Clave Única de Registro de Población (CURP) 固有識別コード。

### オランダのデータタイプ

データタイプ	説明
NETHERLANDS_CITIZEN_SERVICE_NUMBER	オランダ市民番号 (BSN、burgerservice nummer)。
NETHERLANDS_DRIVING_LICENSE	運転免許証番号 (オランダのみ)。
NETHERLANDS_PASSPORT_NUMBER	パスポート番号 (オランダのみ)。

データタイプ	説明
NETHERLANDS_TAX_IDENTIFICATION_NUMBER	納税者番号 (オランダのみ)。
NETHERLANDS_VALUE_ADDED_TAX	付加価値税 (オランダのみ)。
NETHERLANDS_BANK_ACCOUNT	銀行口座番号 (オランダのみ)。

### ニュージーランドのデータタイプ

データタイプ	説明
NEW_ZEALAND_DRIVING_LICENSE	運転免許証番号 (ニュージーランドのみ)。
NEW_ZEALAND_NATIONAL_HEALTH_INDEX_NUMBER	ニュージーランドの国民健康指標番号 (National Health Index Number)。
NEW_ZEALAND_TAX_IDENTIFICATION_NUMBER	納税者番号。内国歳入番号 (Inland Revenue Number) と呼ばれます (ニュージーランドのみ)。

### ノルウェーのデータタイプ

データタイプ	説明
NORWAY_BIRTH_NUMBER	ノルウェーの国民識別番号。
NORWAY_DRIVING_LICENSE	運転免許証番号 (ノルウェーのみ)。
NORWAY_HEALTH_INSURANCE_NUMBER	ノルウェーの健康保険番号。
NORWAY_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (ノルウェーのみ)。
NORWAY_VALUE_ADDED_TAX	付加価値税 (ノルウェーのみ)。

## フィリピンのデータタイプ

データタイプ	説明
PHILIPPINES_DRIVING_LICENSE	運転免許証番号 (フィリピンのみ)。
PHILIPPINES_PASSPORT_NUMBER	パスポート番号 (フィリピンのみ)。

## ポーランドのデータタイプ

データタイプ	説明
POLAND_DRIVING_LICENSE	運転免許証番号 (ポーランドのみ)。
POLAND_IDENTIFICATION_NUMBER	ポーランドの識別番号。
POLAND_PASSPORT_NUMBER	パスポート番号 (ポーランドのみ)。
POLAND_REGON_NUMBER	REGON 識別番号。統計識別番号 (Statistical Identification Number) と呼ばれます。
POLAND_SSN	社会保障番号 (ポーランド人の場合)。
POLAND_TAX_IDENTIFICATION_NUMBER	納税者番号 (ポーランドのみ)。
POLAND_VALUE_ADDED_TAX	付加価値税 (ポーランドのみ)。

## ポルトガルのデータタイプ

データタイプ	説明
PORTUGAL_DRIVING_LICENSE	運転免許証番号 (ポルトガルのみ)。
PORTUGAL_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (ポルトガルのみ)。
PORTUGAL_PASSPORT_NUMBER	パスポート番号 (ポルトガルのみ)。

データタイプ	説明
PORTUGAL_TAX_IDENTIFICATION_NUMBER	納税者番号 (ポルトガルのみ)。
PORTUGAL_VALUE_ADDED_TAX	付加価値税 (ポルトガルのみ)。

### ルーマニアのデータタイプ

データタイプ	説明
ROMANIA_DRIVING_LICENSE	運転免許証番号 (ルーマニアのみ)。
ROMANIA_NUMERICAL_PERSONAL_CODE	個人識別番号 (ルーマニアのみ)。
ROMANIA_PASSPORT_NUMBER	パスポート番号 (ルーマニアのみ)。
ROMANIA_VALUE_ADDED_TAX	付加価値税 (ルーマニアのみ)。

### シンガポールのデータタイプ

データタイプ	説明
SINGAPORE_DRIVING_LICENSE	運転免許証番号 (シンガポールのみ)。
SINGAPORE_NATIONAL_REGISTRY_IDENTIFICATION_NUMBER	シンガポールの国民登録 ID カード。
SINGAPORE_PASSPORT_NUMBER	パスポート番号 (シンガポールのみ)。
SINGAPORE_UNIQUE_ENTITY_NUMBER	シンガポールの個別企業登録番号 (UEN)。

### スロバキアのデータタイプ

データタイプ	説明
SLOVAKIA_DRIVING_LICENSE	運転免許証番号 (スロバキアのみ)。

データタイプ	説明
SLOVAKIA_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (スロバキアのみ)。
SLOVAKIA_PASSPORT_NUMBER	パスポート番号 (スロバキアのみ)。
SLOVAKIA_VALUE_ADDED_TAX	付加価値税 (スロバキアのみ)。

### スロベニアのデータタイプ

データタイプ	説明
SLOVENIA_DRIVING_LICENSE	運転免許証番号 (スロベニアのみ)。
SLOVENIA_PASSPORT_NUMBER	パスポート番号 (スロベニアのみ)。
SLOVENIA_TAX_IDENTIFICATION_NUMBER	納税者番号 (スロベニアのみ)。
SLOVENIA_UNIQUE_MASTER_CITIZEN_NUMBER	スロベニア市民の固有マスター市民番号 (JMBG、Unique master citizen number)。
SLOVENIA_VALUE_ADDED_TAX	付加価値税 (スロベニアのみ)。

### 南アフリカのデータタイプ

データタイプ	説明
SOUTH_AFRICA_PERSONAL_IDENTIFICATION_NUMBER	個人識別番号 (南アフリカのみ)。

### スペインのデータタイプ

データタイプ	説明
SPAIN_BANK_ACCOUNT	銀行口座番号 (スペインのみ)。

データタイプ	説明
SPAIN_DNI	スペインの国民 ID カード (Documento Nacional de Identidad)。
SPAIN_DRIVING_LICENSE	運転免許証番号 (スペインのみ)。
SPAIN_NIE	NIE と呼ばれる外国人 ID 番号 (スペインのみ)。
SPAIN_NIF	NIF と呼ばれる納税者番号 (スペインのみ)。
SPAIN_PASSPORT_NUMBER	パスポート番号 (スペインのみ)。
SPAIN_SSN	社会保障番号 (スペイン人の場合)。
SPAIN_VALUE_ADDED_TAX	付加価値税 (スペインのみ)。

### スリランカのデータタイプ

データタイプ	説明
SRI_LANKA_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (スリランカのみ)。

### スウェーデンのデータタイプ

データタイプ	説明
SWEDEN_DRIVING_LICENSE	運転免許証番号 (スウェーデンのみ)。
SWEDEN_PASSPORT_NUMBER	パスポート番号 (スウェーデンのみ)。
SWEDEN_PERSONAL_IDENTIFICATION_NUMBER	国民識別番号 (スウェーデンのみ)。
SWEDEN_TAX_IDENTIFICATION_NUMBER	スウェーデンの納税者番号 (personnummer)。

データタイプ	説明
SWEDEN_VALUE_ADDED_TAX	付加価値税 (スウェーデンのみ)。

### スイスのデータタイプ

データタイプ	説明
SWITZERLAND_AHV	スイス人の社会保障番号 (AHV)。
SWITZERLAND_HEALTH_INSURANCE_NUMBER	スイスの健康保険番号。
SWITZERLAND_PASSPORT_NUMBER	パスポート番号 (スイスのみ)。
SWITZERLAND_VALUE_ADDED_TAX	付加価値税 (スイスのみ)。

### タイのデータタイプ

データタイプ	説明
THAILAND_PASSPORT_NUMBER	パスポート番号 (タイのみ)。
THAILAND_PERSONAL_IDENTIFICATION_NUMBER	個人識別番号 (タイのみ)。

### トルコのデータタイプ

データタイプ	説明
TURKEY_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (トルコのみ)。
TURKEY_PASSPORT_NUMBER	パスポート番号 (トルコのみ)。
TURKEY_VALUE_ADDED_TAX	付加価値税 (トルコのみ)。

## ウクライナのデータタイプ

データタイプ	説明
UKRAINE_INDIVIDUAL_IDENTIFICATION_NUMBER	固有識別番号 (ウクライナのみ)。
UKRAINE_PASSPORT_NUMBER_DOMESTIC	国内パスポート番号 (ウクライナのみ)。
UKRAINE_PASSPORT_NUMBER_INTERNATIONAL	国際パスポート番号 (ウクライナのみ)。

## アラブ首長国連邦 (UAE) のデータタイプ

データタイプ	説明
UNITED_ARAB_EMIRATES_PERSONAL_NUMBER	個人識別番号 (UAE のみ)。

## 英国のデータタイプ

データタイプ	説明
UK_BANK_ACCOUNT	英国 (UK) の銀行口座。
UK_BANK_SORT_CODE	英国 (UK) の銀行ソートコード。ソートコードは、各国の銀行間で、それぞれの決済機関を通じて送金をルーティングするために使用される銀行コードです。
UK_DRIVING_LICENSE	グレートブリテンおよび北アイルランド連合王国の運転免許証番号 (英国固有)
UK_ELECTORAL_ROLL_NUMBER	選挙人名簿番号 (ERN) は、英国の選挙登録のために個人に発行される識別番号です。この番号の形式は、英国内閣府の英国政府基準によって定められています。

データタイプ	説明
UK_NATIONAL_HEALTH_SERVICE_NUMBER	国民保険サービス (NHS) 番号は、英国の公衆衛生サービスの登録ユーザーに割り当てられる固有の番号です。
UK_NATIONAL_INSURANCE_NUMBER	国民保険番号 (NINO) は、国民保険プログラムまたは社会保障システムで個人を識別するために英国で使用される番号です。この番号は、NI No または NINO と呼ばれることもあります。
UK_PASSPORT_NUMBER	英国 (UK) パスポート番号。
UK_UNIQUE_TAXPAYER_REFERENCE_NUMBER	英国 (UK) 納税者固有番号 (UTR) 番号。英国政府が課税システムを管理するために使用する識別子。
UK_VALUE_ADDED_TAX	VAT は最終消費者が負担する消費税です。VAT は、製造および流通プロセスの各取引に対して支払われます。英国の場合、VAT 番号は事業が設立された地域の VAT 事務所によって発行されます。
UK_PHONE_NUMBER	英国 (UK) の電話番号。

### ベネズエラのデータタイプ

データタイプ	説明
VENEZUELA_DRIVING_LICENSE	運転免許証番号 (ベネズエラのみ)。
VENEZUELA_NATIONAL_IDENTIFICATION_NUMBER	国民識別番号 (ベネズエラのみ)。
VENEZUELA_VALUE_ADDED_TAX	付加価値税 (ベネズエラのみ)。

## 詳細な機密データ検出の使用

### Note

詳細なアクションは AWS Glue 3.0 および 4.0 でのみ使用できます。これには AWS Glue Studio エクスペリエンスも含まれます。永続的な監査ログの変更も 2.0 では使用できません。

すべての AWS Glue Studio 3.0 および 4.0 のビジュアルジョブでは、詳細なアクション API を自動的に使用するスクリプトが作成されます。

Detect Sensitive Data は、定義したエンティティ、または AWS Glue によって事前定義されたエンティティを検出、マスキング、削除する機能を提供します。さらに、詳細なアクションを使用することで、エンティティごとに特定のアクションを適用できます。その他の利点には次が含まれます。

- データが検出されるとすぐにアクションが適用されるようになり、パフォーマンスが改善されました。
- 特定の列を含めるか、または除外するオプション。
- 部分的なマスキングを使用する機能。これにより、文字列全体をマスキングするのではなく、検出された機密データエンティティを部分的にマスキングできます。オフセットを含むシンプルなパラメータと正規表現の両方がサポートされます。

次のセクションで参照するサンプルジョブで使用される機密データ検出 API と詳細なアクションのコードスニペットを次に示します。

detect API – 詳細なアクションは、新しい `detectionParameters` パラメータを使用します。

```
def detect(
 frame: DynamicFrame,
 detectionParameters: JsonOptions,
 outputColumnName: String = "DetectedEntities",
 detectionSensitivity: String = "LOW"
): DynamicFrame = {}
```

## 詳細なアクションで機密データ検出 API を使用する

detect を使用する機密データ検出 API は、指定されたデータを分析し、行または列が機密データエンティティタイプであるかどうかを判断して、エンティティタイプごとにユーザーが指定したアクションを実行します。

### 詳細なアクションでの detect API の使用

detect API を使用して、outputColumnName と detectionParameters を指定します。

```
object GlueApp {
 def main(sysArgs: Array[String]) {

 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)

 // @params: [JOB_NAME]
 val args = GlueArgParser.getResolvedOptions(sysArgs, Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)

 // Script generated for node S3 bucket. Creates DataFrame from data stored in
 S3.
 val S3bucket_node1 =
glueContext.getSourceWithFormat(formatOptions=JsonOptions("""{"quoteChar":
"\\"", "withHeader": true, "separator": ",", "optimizePerformance": false}"""),
connectionType="s3", format="csv", options=JsonOptions("""{"paths":
["s3://189657479688-ddevansh-pii-test-bucket/tiny_pii.csv"], "recurse": true}"""),
transformationContext="S3bucket_node1").getDynamicFrame()

 // Script generated for node Detect Sensitive Data. Will run detect API for the
 DataFrame
 // detectionParameter contains information on which EntityType are being
 detected
 // and what actions are being applied to them when detected.
 val DetectSensitiveData_node2 = EntityDetector.detect(
 frame = S3bucket_node1,
 detectionParameters = JsonOptions(
 """
 {
 "PHONE_NUMBER": [
 {
 "action": "PARTIAL_REDACT",
```

```

 "actionOptions": {
 "numLeftCharsToExclude": "3",
 "numRightCharsToExclude": "4",
 "redactChar": "#"
 },
 "sourceColumnsToExclude": ["Passport No", "DL NO#"]
 }
],
"USA_PASSPORT_NUMBER": [
 {
 "action": "SHA256_HASH",
 "sourceColumns": ["Passport No"]
 }
],
"USA_DRIVING_LICENSE": [
 {
 "action": "REDACT",
 "actionOptions": {
 "redactText": "USA_DL"
 },
 "sourceColumns": ["DL NO#"]
 }
]
}
""
),
outputColumnName = "DetectedEntities"
)

// Script generated for node S3 bucket. Store Results of detect to S3 location
val S3bucket_node3 = glueContext.getSinkWithFormat(connectionType="s3",
options=JsonOptions("""{"path": "s3://189657479688-ddevansh-pii-test-bucket/
test-output/", "partitionKeys": []}"""), transformationContext="S3bucket_node3",
format="json").writeDynamicFrame(DetectSensitiveData_node2)

Job.commit()
}

```

上記のスク립トは、Amazon S3 内の場所から DataFrame を作成し、detect API を実行します。detect API ではフィールド detectionParameters (エンティティに使用されるすべてのア

クシオン設定のリストに対するエンティティ名のマップ) が AWS Glue の `JsonOptions` オブジェクトによって表されるため、API の機能を拡張することもできます。

エンティティごとに指定された各アクションについて、エンティティ/アクションの組み合わせを適用するすべての列名のリストを入力します。これにより、エンティティをカスタマイズしてデータセット内のすべての列を検出したり、特定の列に存在しないことがわかっているエンティティをスキップしたりできます。また、これらのエンティティを呼び出す不要な検出を実行しないことでジョブのパフォーマンスが改善し、各列とエンティティの組み合わせに固有のアクションを実行できるようになります。

`detectionParameters` を詳しく見てみると、サンプルジョブには 3 つのエンティティタイプがあることがわかります。これらは、`Phone Number`、`USA_PASSPORT_NUMBER`、`USA_DRIVING_LICENSE` です。これらの各エンティティタイプについて、AWS Glue は、異なるアクション (`PARTIAL_REDACT`、`SHA256_HASH`、`REDACT`、`DETECT` のいずれか) を実行します。各エンティティタイプには、適用する `sourceColumns` および/または `sourceColumnsToExclude` (検出された場合) があります。

#### Note

列ごとに使用できるその場編集アクション (`PARTIAL_REDACT`、`SHA256_HASH`、または `REDACT`) は 1 つだけですが、`DETECT` アクションはこれらのアクションのいずれかと併用できます。

`detectionParameters` フィールドのレイアウトは以下のとおりです。

```
ENTITY_NAME -> List[Actions]
{
 "ENTITY_NAME": [{
 Action, // required
 ColumnSpecs,
 ActionOptionsMap
 }],
 "ENTITY_NAME2": [{
 ...
 }]
}
```

タイプ actions および actionOptions は以下のとおりです。

```
DETECT
{
 # Required
 "action": "DETECT",
 # Optional, depending on action chosen
 "actionOptions": {
 // There are no actionOptions for DETECT
 },
 # 1 of below required, both can also used
 "sourceColumns": [
 "COL_1", "COL_2", ..., "COL_N"
],
 "sourceColumnsToExclude": [
 "COL_5"
]
}

SHA256_HASH
{
 # Required
 "action": "SHA256_HASH",
 # Required or optional, depending on action chosen
 "actionOptions": {
 // There are no actionOptions for SHA256_HASH
 },

 # 1 of below required, both can also used
 "sourceColumns": [
 "COL_1", "COL_2", ..., "COL_N"
],
 "sourceColumnsToExclude": [
 "COL_5"
]
}

REDACT
{
 # Required
 "action": "REDACT",
 # Required or optional, depending on action chosen
 "actionOptions": {
 // The text that is being replaced
```

```
 "redactText": "USA_DL"
 },

 # 1 of below required, both can also used
 "sourceColumns": [
 "COL_1", "COL_2", ..., "COL_N"
],
 "sourceColumnsToExclude": [
 "COL_5"
]
}

PARTIAL_REDACT
{
 # Required
 "action": "PARTIAL_REDACT",
 # Required or optional, depending on action chosen
 "actionOptions": {
 // number of characters to not redact from the left side
 "numLeftCharsToExclude": "3",
 // number of characters to not redact from the right side
 "numRightCharsToExclude": "4",
 // the partial redact will be made with this redacted character
 "redactChar": "#",
 // regex pattern for partial redaction
 "matchPattern": "[0-9]"
 },

 # 1 of below required, both can also used
 "sourceColumns": [
 "COL_1", "COL_2", ..., "COL_N"
],
 "sourceColumnsToExclude": [
 "COL_5"
]
}
```

スクリプトが実行されると、指定された Amazon S3 の場所に結果が出力されます。Amazon S3 でデータを表示することはできますが、選択したアクションに基づいて、選択したエンティティタイプの感度が高くなります。この場合、行は次のようになります。

```
{
```

```
"Name": "Colby Schuster",
"Address": "39041 Antonietta Vista, South Rodgerside, Nebraska 24151",
"Car Owned": "Fiat",
"Email": "Kitty46@gmail.com",
"Company": "O'Reilly Group",
"Job Title": "Dynamic Functionality Facilitator",
"ITIN": "991-22-2906",
"Username": "Cassandre.Kub43",
"SSN": "914-22-2906",
"DOB": "2020-08-27",
"Phone Number": "1-2#####1718",
"Bank Account No": "69741187",
"Credit Card Number": "6441-6289-6867-2162-2711",
"Passport No": "94f311e93a623c72ccb6fc46cf5f5b0265ccb42c517498a0f27fd4c43b47111e",
"DL NO#": "USA_DL"
}
```

上記のスク립トでは、Phone Number は部分的に # でマスキングされています。Passport No は SHA256 ハッシュに変更されました。DL NO# は米国の運転免許証番号として検出され、detectionParameters で記述されているように「USA\_DL」でマスキングされました。

#### Note

classifyColumns API は、API の性質上、詳細なアクションでは使用できません。この API は、列サンプリング (ユーザーによる調整が可能です) が、デフォルト値が設定されています) を実行して、検出をより迅速に実行します。このため、詳細なアクションでは、すべての値に対してイテレーションを実行する必要があります。

## 永続的な監査ログ

詳細なアクションで導入された新機能 (ただし、通常の API を使用する場合にも使用可能) は、永続的な監査ログです。現在、detect API を実行すると、PII 検出メタデータを含む追加の列 (DetectedEntities にデフォルト設定されていますが、outputColumnName を通じてカスタマイズ可能です) パラメータが追加されます。これには、「actionUsed」メタデータキーが追加されました。これは、DETECT、PARTIAL\_REDACT、SHA256\_HASH、REDACT のいずれかです。

```
"DetectedEntities": {
 "Credit Card Number": [
 {
```

```

 "entityType": "CREDIT_CARD",
 "actionUsed": "DETECT",
 "start": 0,
 "end": 19
 }
],
"Phone Number": [
 {
 "entityType": "PHONE_NUMBER",
 "actionUsed": "REDACT",
 "start": 0,
 "end": 14
 }
]
}

```

`detect(entityTypesToDetect, outputColumnName)` などの詳細なアクションなしで API を使用しているお客様でも、結果として得られるデータフレームでこの永続的な監査ログを確認できません。

詳細なアクションが設定されている API を使用しているお客様は、マスキングされているかどうかにかかわらず、すべてのアクションを表示できます。例：

```

+-----+-----+-----+-----+
+-----+-----+-----+-----+
+
| Credit Card Number | Phone Number |
+-----+-----+-----+-----+
+
+-----+-----+-----+-----+
+-----+-----+-----+-----+
+
| 622126741306XXXX | +12#####7890 | {"Credit Card Number":
[{"entityType":"CREDIT_CARD","actionUsed":"PARTIAL_REDACT","start":0,"end":16}], "Phone
Number":
[{"entityType":"PHONE_NUMBER","actionUsed":"PARTIAL_REDACT","start":0,"end":12}]} |
| 6221 2674 1306 XXXX | +12#####7890 | {"Credit Card Number":
[{"entityType":"CREDIT_CARD","actionUsed":"PARTIAL_REDACT","start":0,"end":19}], "Phone
Number":
[{"entityType":"PHONE_NUMBER","actionUsed":"PARTIAL_REDACT","start":0,"end":14}]} |
| 6221-2674-1306-XXXX | 22#####7890 | {"Credit Card Number":
[{"entityType":"CREDIT_CARD","actionUsed":"PARTIAL_REDACT","start":0,"end":19}], "Phone

```

```
Number":
[{"entityType":"PHONE_NUMBER","actionUsed":"PARTIAL_REDACT","start":0,"end":14}]]} |
+-----+-----
+-----
+
+
```

[DetectedEntities] 列を表示したくない場合、必要なのは、カスタムスクリプトに追加の列をドロップすることだけです。

## AWS Glue ビジュアルジョブ API

AWS Glue は ビジュアルステップワークフローを表す JSON オブジェクトから AWS Glue API を使用して、お客様がデータ統合ジョブを作成できるようにする API を提供します。お客様は AWS Glue Studio のビジュアルエディターを使って、これらのジョブを処理することができます。

Visual Job API データ型の詳細については、「[Visual Job API](#)」(ビジュアル Job API) を参照してください。

### トピック

- [API 設計と CRUD API](#)
- [開始](#)
- [ビジュアルジョブの制限](#)

### API 設計と CRUD API

CreateJobAPI と UpdateJob の [API](#) は、追加のオプションパラメータとして codeGenConfigurationNodes をサポートするようになりました。このフィールドに空でない JSON 構造を指定すると、作成されたジョブの AWS Glue Studio に対して DAG が登録され、関連するコードが生成されます。ジョブ作成時のこのフィールドの NULL 値または空の文字列は無視されません。

codegenConfigurationNodes フィールドの更新は、CreateJob と同様の方法で UpdateJob AWS Glue API で行われます。DAG が必要に応じて変更された UpdateJob でフィールド全体を指定する必要があります。指定された NULL 値は無視され、DAG への更新は実行されません。空の構造体または文字列を指定すると、CodeGenConfigurationNodes が空に設定され、以前の DAG が削除されます。GetJob API は DAG が存在する場合、DAG を返します。DeleteJob API では、関連する DAG も削除されます。

## 開始

ジョブを作成するには、[CreateJob アクション](#)を使用します。CreateJob リクエストの入力には、JSON で DAG オブジェクトを指定できる追加のフィールド「CodeGenConfigurationNodes」があります。

留意すべきこと:

- 「CodeGenConfigurationNodes」フィールドは、ノードへの nodeID のマップです。
- 各ノードは、ノードの種類を識別するキーで始まります。
- ノードには 1 つのタイプしか指定できないため、キーは 1 つしか指定できません。
- 入力フィールドには、現在のノードの親ノードが含まれます。

以下は、CreateJob 入力の JSON 表現です。

```
{
 "node-1": {
 "S3CatalogSource": {
 "Table": "csvFormattedTable",
 "PartitionPredicate": "",
 "Name": "S3 bucket",
 "AdditionalOptions": {},
 "Database": "myDatabase"
 }
 },
 "node-3": {
 "S3DirectTarget": {
 "Inputs": ["node-2"],
 "PartitionKeys": [],
 "Compression": "none",
 "Format": "json",
 "SchemaChangePolicy": { "EnableUpdateCatalog": false },
 "Path": "",
 "Name": "S3 bucket"
 }
 },
 "node-2": {
 "ApplyMapping": {
 "Inputs": ["node-1"],
 "Name": "ApplyMapping",
 "Mapping": [
```

```
{
 "FromType": "long",
 "ToType": "long",
 "Dropped": false,
 "ToKey": "myheader1",
 "FromPath": ["myheader1"]
},
{
 "FromType": "long",
 "ToType": "long",
 "Dropped": false,
 "ToKey": "myheader2",
 "FromPath": ["myheader2"]
},
{
 "FromType": "long",
 "ToType": "long",
 "Dropped": false,
 "ToKey": "myheader3",
 "FromPath": ["myheader3"]
}
]
}
}
```

## ジョブの更新と取得

[Updating jobs] (ジョブの更新) に 'codegenConfigurationNodes' フィールドがあるため、入力形式は同じになります。「[UpdateJob](#)」アクションを参照してください。

GetJob アクションは、同じ形式の「CodeGenConfigurationNodes」フィールドも返します。「[GetJob](#)」アクションを参照してください。

## ビジュアルジョブの制限

既存の API に「CodeGenConfigurationNodes」パラメータが追加されているため、これらの API の制限はすべて継承されます。さらに、CodeGenConfigurationNodes および一部のノードのサイズは制限されます。詳細については、次の「[ジョブ構造](#)」を参照してください。

## Ray スクリプトのプログラミング

AWS Glue を使用すると、Ray スクリプトの記述と実行が簡単になります。このセクションでは、AWS Glue for Ray でサポートされている Ray の機能について説明します。Ray スクリプトは、Python でプログラムします。

カスタムスクリプトは、ジョブ定義内の Runtime フィールドで定義されている Ray のバージョンと互換性がある必要があります。Jobs API の Runtime の詳細については、「[the section called “ジョブ”](#)」を参照してください。各ランタイム環境の詳細については、「[the section called “サポートされている Ray のランタイム環境”](#)」を参照してください。

### トピック

- [チュートリアル: AWS Glue for Ray で ETL スクリプトを記述する](#)
- [AWS Glue for Ray での Ray Core と Ray Data の使用](#)
- [ファイルと Python ライブラリを Ray ジョブに提供する](#)
- [Ray ジョブのデータに接続する](#)

## チュートリアル: AWS Glue for Ray で ETL スクリプトを記述する

Ray を使用すると、分散したタスクを Python でネイティブに記述し、スケールすることができます。AWS Glue for Ray には、ジョブとインタラクティブセッションの両方からアクセスできるサーバーレスの Ray 環境が用意されています (Ray のインタラクティブセッションは現在プレビュー中です)。AWS Glue のジョブシステムを使用すると、スケジュールに基づいて、トリガーにより、または AWS Glue コンソールを使用して、一貫した方法でタスクを管理し、実行できます。

これら AWS Glue のツールを組み合わせると、AWS Glue の一般的なユースケースである、抽出、変換、ロード (ETL) ワークロードに使用できる強力なツールチェーンができあがります。このチュートリアルでは、このソリューションの組み立ての基礎を学びます。

また、ETL ワークロードで AWS Glue for Spark を使用方法についても学習します。AWS Glue for Spark スクリプトの記述方法に関するチュートリアルは、「[the section called “チュートリアル: Spark スクリプトの記述”](#)」を参照してください。利用可能なエンジンの詳細については、「[the section called “AWS Glue for Spark と AWS Glue for Ray”](#)」を参照してください。Ray は、分析、機械学習 (ML)、アプリケーション開発におけるさまざまな種類のタスクに対応できます。

このチュートリアルでは、Amazon Simple Storage Service (Amazon S3) でホストされている CSV データセットの抽出、変換、ロードを行います。まず、パブリックな Amazon S3 バケットに保存さ

れている New York City Taxi and Limousine Commission (TLC) の乗車記録データのデータセットから始めます。このデータセットの詳細については、「[Registry of Open Data on AWS](#)」を参照してください。

Ray Data ライブラリで利用できる事前定義された変換を使用して、データを変換していきます。Ray Data は Ray が設計したデータセットの準備ライブラリで、Ray の環境用に AWS Glue にデフォルトで含まれています。デフォルトで含まれているライブラリの詳細については、「[the section called “Ray ジョブで提供されるモジュール”](#)」を参照してください。次に、変換したデータを、自分が管理する Amazon S3 バケットに書き込みます。

前提条件 — このチュートリアルを行うには、AWS Glue と Amazon S3 へのアクセス権限を持つ AWS アカウントが必要です。

## ステップ 1: Amazon S3 に出力データを保存するバケットを作成する

このチュートリアルで作成したデータのシンクとして機能する、ユーザーが管理する Amazon S3 バケットが必要です。このバケットは、以下の手順で作成します。

### Note

自分で管理している既存のバケットにデータを書き込む場合は、このステップは飛ばしてかまいません。その場合は、既存のバケット名である *yourBucketName* をメモしておきます。こちらは後のステップで使用します。

Ray ジョブの出力用バケットを作成するには

- 「Amazon S3 ユーザーガイド」の「[バケットの作成](#)」の手順に従ってバケットを作成します。
- バケット名を選択ぶときは *yourBucketName* をメモしておきます。こちらは後のステップで使用します。
- 他の設定については、Amazon S3 コンソールに表示された推奨の設定を使えばチュートリアルを問題なく進められるはずです。

例えば、Amazon S3 コンソールのバケット作成のダイアログボックスは、次のように表示されます。

Amazon S3 &gt; Buckets &gt; Create bucket

## Create bucket Info

Buckets are containers for data stored in S3. [Learn more](#)

### General configuration

Bucket name

Bucket name must be globally unique and must not contain spaces or uppercase letters. [See rules for bucket naming](#)

AWS Region

## ステップ 2: Ray ジョブ用の IAM ロールを作成する

ジョブには、以下を持つ AWS Identity and Access Management (IAM) ロールが必要です。

- `AWSGlueServiceRole` マネージドポリシーが付与したアクセス権限。これらは、AWS Glue ジョブを実行するために必要な、基本的な権限です。
- `nyc-tlc/*` Amazon S3 リソースの Read のアクセスレベル権限。
- `yourBucketName/*` Amazon S3 リソースの Write のアクセスレベル権限。
- `glue.amazonaws.com` プリンシパルにロールの引き受けを許可する信頼関係。

このロールは、以下の手順で作成します。

AWS Glue for Ray ジョブ用に IAM ロールを作成するには

### Note

IAM ロールは、さまざまな手順で作成できます。IAM リソースのプロビジョニング方法の詳細またはオプションについては、[AWS Identity and Access Management のドキュメント](#)を参照してください。

1. 上記の Amazon S3 アクセス権を定義するポリシーを、「IAM ユーザーガイド」の「[ビジュアルエディタを使用してポリシーを作成するには](#)」の手順に従って作成します。
  - サービスを選ぶときは、Amazon S3 を選択します。
  - ポリシーのアクセス権を選ぶときは、以下のリソースの、以下に挙げるアクションのセットを添付します (上記のとおり)。
    - `nyc-tlc/*` Amazon S3 リソースの Read のアクセスレベル権限。
    - `yourBucketName/*` Amazon S3 リソースの Write のアクセスレベル権限。
  - ポリシー名を選ぶときは、`YourPolicyName` をメモしておきます。こちらは後のステップで使用します。
2. AWS Glue for Ray ジョブ用のロールを、「IAM ユーザーガイド」の「[AWS サービスのロールを作成するには \(コンソール\)](#)」の手順に従って作成します。
  - 信頼された AWS サービスエンティティを選ぶときは、Glue を選択します。これにより、ジョブに必要な信頼関係が自動的に設定されます。
  - アクセス権ポリシーのポリシーを選ぶときは、次のポリシーを添付します。
    - `AWSGlueServiceRole`
    - `YourPolicyName`
  - ロール名を選ぶときは、`YourRoleName` をメモしておきます。こちらは後のステップで使用します。

### ステップ 3: AWS Glue for Ray ジョブを作成して実行する

このステップでは、AWS Management Console を使用して AWS Glue ジョブを作成し、これにサンプルスクリプトを提供して、ジョブを実行します。ジョブを作成すると、Ray スクリプトを保存、設定、編集するための場所がコンソールに作成されます。ジョブの作成の詳細については、「[the section called “コンソールにサインインする”](#)」をご参照ください。

このチュートリアルでは、次の ETL シナリオを取り上げます。New York City TLC の乗車記録データセットから 2022 年 1 月のレコードを読み取り、既存の列のデータを組み合わせてデータセットに新しい列 (`tip_rate`) を追加し、現在の分析に無関係のいくつかの列を削除してから、結果を `yourBucketName` に書き込みます。次の Ray スクリプトは、以下のステップを実行します。

```
import ray
import pandas
from ray import data
```

```
ray.init('auto')

ds = ray.data.read_csv("s3://nyc-tlc/opendata_repo/opendata_webconvert/yellow/
yellow_tripdata_2022-01.csv")

Add the given new column to the dataset and show the sample record after adding a new
column
ds = ds.add_column("tip_rate", lambda df: df["tip_amount"] / df["total_amount"])

Dropping few columns from the underlying Dataset
ds = ds.drop_columns(["payment_type", "fare_amount", "extra", "tolls_amount",
"improvement_surcharge"])

ds.write_parquet("s3://yourBucketName/ray/tutorial/output/")
```

AWS Glue for Ray ジョブを作成して実行するには

1. AWS Management Console で、AWS Glue のランディングページに進みます。
2. ナビゲーションペインで、[ETL Jobs] を選択します。
3. [ジョブを作成する] で、[Ray script editor] を選択し、[作成] をクリックします。以下の図を参照してください。

## AWS Glue Studio Info

**Create job** Info Create

<input type="radio"/> <b>Visual with a source and target</b> Start with a source, ApplyMapping transform, and target.	<input type="radio"/> <b>Visual with a blank canvas</b> Author using an interactive visual interface.	<input type="radio"/> <b>Spark script editor</b> Write or upload your own Spark code.
<input type="radio"/> <b>Python Shell script editor</b> Write or upload your own Python shell script.	<input type="radio"/> <b>Jupyter Notebook</b> Write your own code in a Jupyter Notebook for interactive development.	<input checked="" type="radio"/> <b>Ray script editor</b> Write your own code to run on Ray.

**Options** Info

- Create a new script with boilerplate code
- Upload and edit an existing script  
Choose a local file.

4. スクリプトの全文を [スクリプト] ペインに貼り付け、既存のテキストがあれば置き換えます。
5. [ジョブの詳細] に進み、[IAM ロール] のプロパティを *YourRoleName* に設定します。
6. [保存] をクリックし、[実行] をクリックします。

## ステップ 4: 出力を確認する

AWS Glue ジョブを実行したら、出力がこのシナリオの期待値と一致していることを確認する必要があります。確認は、次の手順で行います。

Ray ジョブが正常に実行されたことを確認するには

1. ジョブの詳細ページから [実行] に移動します。
2. 数分後、実行の [実行ステータス] が [成功] と表示されるはずですが。
3. Amazon S3 コンソール (<https://console.aws.amazon.com/s3/>) に移動し、*yourBucketName* をチェックします。出力バケットに書き込まれたファイルが表示されるはずですが。
4. Parquet ファイルを読み、内容を確認します。これは既存のツールで実行できます。Parquet ファイルを確認するプロセスがない場合は、Spark または Ray を使用して、AWS Glue コンソールで AWS Glue インタラクティブセッションを実行します (現在プレビュー中)。

インタラクティブセッションでは、Ray Data、Spark、pandas のいずれかのライブラリにアクセスします。これは (選択したエンジンに応じて) デフォルトで指定されています。ファイルの内容の確認には、これらのライブラリで利用できる一般的な確認方法、例えば、count、schema、show などを使用できます。コンソールのインタラクティブセッションの詳細については、「[AWS Glue Studio と AWS Glue を用いてノートブックの使用](#)」を参照してください。

ファイルがバケットに書き込まれていることを確認できたので、出力に問題があったとしても、それは IAM の設定とは無関係であると確実に言えます。*yourRoleName* を使用してセッションを設定するのは、関連するファイルにアクセスできるようにするためです。

期待した結果が得られない場合は、このガイドのトラブルシューティングのページを参照し、エラーの原因を特定して修正します。ジョブの実行エラー状態を解釈するには、「[the section called “ジョブ実行ステータス”](#)」を参照してください。トラブルシューティングのページは「[AWS Glue のトラブルシューティング](#)」にあります。Ray ジョブに関連する特定のエラーについては、トラブルシューティング内の「[the section called “Ray エラーをトラブルシューティングする”](#)」を参照してください。

## 次のステップ

これで、AWS Glue for Ray を使用した ETL のプロセスの全体を確認し、実行できました。次のリソースを使用すると、AWS Glue for Ray がデータを大規模に変換し解釈するためにどのようなツールを提供しているのかを理解できます。

- Ray のタスクのモデルに関する詳細は、「[the section called “AWS Glue for Ray での Ray Core と Ray Data の使用”](#)」を参照してください。Ray のタスクの使用に関する詳細は、Ray Core ドキュメントにある例を参照してください。Ray のドキュメントの「[Ray Core: Ray Tutorials and Examples \(2.4.0\)](#)」を参照してください。
- AWS Glue for Ray で利用可能なデータ管理ライブラリのガイダンスは、「[the section called “データへの接続”](#)」を参照してください。Ray Data を使用してデータセットを変換し記述する方法の詳細については、Ray Data ドキュメントにある例を参照してください。「[Ray Data: Examples \(2.4.0\)](#)」を参照してください。
- AWS Glue for Ray ジョブの設定に関する詳細は、「[the section called “Ray ジョブの使用”](#)」を参照してください。
- AWS Glue for Ray のスクリプトの記述に関する詳細は、本セクション内のドキュメントを引き続きお読みください。

## AWS Glue for Ray での Ray Core と Ray Data の使用

Ray は、クラスター全体に作業を分散することで Python スクリプトをスケールアップするためのフレームワークです。Ray はさまざまな問題の解決策として利用できます。つまり、特定のタスクを最適化するためのライブラリを提供しています。AWS Glue では、Ray を使用して大規模なデータセットを変換することに重点を置いています。AWS Glue は Ray Core の一部と Ray Data をサポートしているため、このタスクが容易になります。

### Ray Core とは

分散アプリケーションを構築する最初のステップは、同時に実行できる作業を特定して定義することです。Ray Core には、同時に実行できるタスクを定義するために使用する Ray の機能が一部含まれています。Ray では、提供されるツールを学習するのに役立つリファレンスおよびクイックスタート情報を入手できます。詳細については、「[What is Ray Core?](#)」と「[Ray Core Quick Start](#)」を参照してください。Ray で同時タスクを効果的に定義する方法の詳細については、「[Tips for first-time users](#)」を参照してください。

#### Ray のタスクとアクター

AWS Glue for Ray ドキュメントでは、Ray の中核的な概念であるタスクとアクターについて言及することがあります。

Ray では、分散コンピューティングシステムの構成要素として Python の関数とクラスを使用します。Python の関数や変数をクラスで使用すると「メソッド」や「属性」になるのと同様に、ワーカーにコードを送信するために Ray で使用すると、関数は「タスク」

になり、クラスは「アクター」になります。Ray で使用される可能性のある関数やクラスは、`@ray.remote` アノテーションで識別できます。

タスクとアクターは設定可能で、ライフサイクルがあり、そのライフサイクルを通じてコンピューティングリソースを消費します。エラーをスローするコードは、問題の根本原因を見つける場合に、タスクまたはアクターまでさかのぼることができます。したがって、AWS Glue for Ray ジョブの設定、モニタリング、デバッグの方法を学習するときに、これらの用語が出現する可能性があります。

タスクとアクターを効果的に使用して分散アプリケーションを構築する方法を学習するには、Ray のドキュメントの「[Key Concepts](#)」を参照してください。

## AWS Glue for Ray での Ray Core

AWS Glue for Ray 環境では、ログの収集と視覚化だけでなく、クラスターの形成とスケーリングも管理します。こうした懸念事項を管理しているため、当社は、オープンソースのクラスターでこれらの対処に使用される、Ray Core の API へのアクセスとサポートを制限しています。

マネージド Ray 2.4 ランタイム環境では、以下はサポートされていません。

- [Ray Core CLI](#)
- [Ray State CLI](#)
- `ray.util.metrics` Prometheus メトリクスユーティリティメソッド:
  - [Counter](#)
  - [Gauge](#)
  - [Histogram](#)
- その他のデバッグツール
  - [ray.util.pdb.set\\_trace](#)
  - [ray.util.inspect\\_serializability](#)
  - [ray.timeline](#)

## Ray Data とは

データソースや送信先に接続する場合、データセットを処理する場合、一般的な変換を開始する場合などに Ray Data を使用すれば、Ray データセットの変換に関する問題を Ray で簡単に解決できます。Ray Data の使用方法の詳細については、「[Ray Datasets: Distributed Data Preprocessing](#)」を参照してください。

Ray Data やその他のツールを使用すると、データにアクセスできます。Ray 内におけるデータへのアクセスの詳細については、「[the section called “データへの接続”](#)」を参照してください。

## AWS Glue for Ray での Ray Data

Ray Data では、マネージド Ray2.4 ランタイム環境がデフォルトでサポートおよび提供されています。提供されるモジュールの詳細については、「[the section called “Ray ジョブで提供されるモジュール”](#)」を参照してください。

## ファイルと Python ライブラリを Ray ジョブに提供する

このセクションでは、AWS Glue の Ray ジョブで Python ライブラリを使用するために必要な情報を提供します。すべての Ray ジョブにデフォルトで含まれている、特定の共通ライブラリを使用できます。また、Ray ジョブに独自の Python ライブラリを指定することもできます。

### Ray ジョブで提供されるモジュール

次のパッケージを使用して、Ray ジョブでデータ統合ワークフローを実行できます。これらのパッケージは Ray ジョブで、デフォルトで使用できます。

#### AWS Glue version 4.0

AWS Glue 4.0 では、Ray (Ray2.4 ランタイム) 環境で次のパッケージが利用可能です。

- boto3 == 1.26.133
- ray == 2.4.0
- pyarrow == 11.0.0
- pandas == 1.5.3
- numpy == 1.24.3
- fsspec == 2023.4.0

このリストには、ray[data] == 2.4.0 でインストールされるすべてのパッケージが含まれます。Ray データは追加設定なしにサポートされます。

### Ray ジョブへのファイルの提供

--working-dir パラメータを使用して Ray ジョブにファイルを提供できます。このパラメータには、Amazon S3 でホストされている .zip ファイルへのパスを指定します。.zip ファイル内では、

ファイルは最上位の 1 つのディレクトリに含まれている必要があります。最上位には、他のファイルがないようにします。

ファイルは、スクリプトの実行を開始する前に各 Ray ノードに配布されます。それにより、各 Ray ノードが利用できるディスク容量にどのような影響が出るかを考えます。利用できるディスク容量は、ジョブ構成で設定された WorkerType により決まります。ジョブに大量のデータを提供したい場合、このメカニズムは適切なソリューションではありません。データをジョブに提供する方法の詳細については、「[the section called “データへの接続”](#)」を参照してください。

あたかもディレクトリが `working_dir` パラメータで Ray に提供されたかのように、ファイルにアクセスできるようになります。例えば、`.zip` ファイルの最上位ディレクトリにある `sample.txt` という名前のファイルを読み取るには、次のように呼び出すことができます。

```
@ray.remote
def do_work():
 f = open("sample.txt", "r")
 print(f.read())
```

`working_dir` の詳細については、[Ray のドキュメント](#)を参照してください。この機能は、Ray のネイティブの機能と同じように動作します。

## Ray ジョブ用の Python モジュールを追加する

### PyPI の追加モジュール

Ray ジョブでは、Python Package Installer (pip3) を使用して、Ray スクリプトで使用するモジュールを追加でインストールします。`--pip-install` パラメータでコマンド区切りの Python モジュールのリストを指定することで、新しいモジュールを追加したり、既存のモジュールのバージョンを変更したりできます。

例えば、更新や新しい `scikit-learn` モジュールの追加には、次のキー値ペアを使用します。

```
"--pip-install", "scikit-learn==0.21.3"
```

カスタムモジュールまたはカスタムパッチがある場合は、`--s3-py-modules` パラメータを使用して Amazon S3 から自分のライブラリをディストリビューションできます。ディストリビューションは、アップロード前に再パッケージ化と構築が必要になる場合があります。[the section called “Ray ジョブに Python コードを含める”](#) のガイドラインに従います。

### Amazon S3 のカスタムディストリビューション

カスタムディストリビューションは、依存関係に関する Ray パッケージのガイドラインに従う必要があります。これらのディストリビューションの構築方法については、次のセクションをご覧ください。Ray での依存関係の設定方法に関する詳細は、Ray のドキュメントの「[Environment Dependencies](#)」(環境の依存関係)を参照してください。

コンテンツを評価した後にカスタムディストリビューションを含めるには、ジョブの IAM ロールで利用できるバケットにディストリビューションをアップロードします。パラメータの設定で、Python zip アーカイブへの Amazon S3 パスを指定します。複数のディストリビューションを指定する場合は、カンマで区切ります。例:

```
"--s3-py-modules", "s3://s3bucket/pythonPackage.zip"
```

### 制約事項

Ray ジョブは、ジョブ環境でのネイティブコードのコンパイルをサポートしていません。Python の依存関係がネイティブのコンパイル済みコードに推移的に依存している場合、これにより制限される可能性があります。Ray ジョブは提供されたバイナリを実行できますが、Linux on ARM64 用にコンパイルする必要があります。つまり、aarch64manylinux wheel の内容を使用できる場合があるということです。Ray の基準に合わせて wheel を再パッケージ化することで、ネイティブの依存関係をコンパイルされた形式で提供できます。通常、これは dist-info フォルダを削除して、アーカイブのルートにあるフォルダが 1 つだけになるようにすることを意味します。

このパラメータを使用して ray または ray[data] のバージョンをアップグレードすることはできません。Ray の新しいバージョンを使用するには、ジョブのランタイムフィールドを、そのサポートをリリースした後に変更する必要があります。サポートされている Ray バージョンの詳細については、「[the section called “AWS Glue バージョン”](#)」を参照してください。

### Ray ジョブに Python コードを含める

Python Software Foundation は、さまざまなランタイムで使用できるよう、Python ファイルをパッケージ化するための標準化された機能を提供しています。Ray では、注意すべきパッケージ基準に制限が設けられています。AWS Glue が、Ray に指定された以外のパッケージ基準を指定することはありません。次の手順では、シンプルな Python パッケージのパッケージ化に関する標準的なガイドンスを提供します。

ファイルを .zip アーカイブにパッケージ化します。1 つのディレクトリがアーカイブのルートにある必要があります。アーカイブのルートレベルにその他のファイルがあると、予期しない動作が発生する可能性があります。ルートディレクトリはパッケージで、インポート時に Python コードを参照するためにその名前が使用されます。

--s3-py-modules を使用してこの形式のディストリビューションを Ray ジョブに提供すると、Ray スクリプト内のパッケージから Python コードをインポートできます。

パッケージは、単一の Python モジュールにいくつかの Python ファイルを提供することも、多数のモジュールと一緒にパッケージ化することもできます。PyPI のライブラリなどの依存関係を再パッケージ化するときは、それらのパッケージ内の隠しファイルやメタデータディレクトリを確認してください。

#### Warning

OS の動作によっては、これらのパッケージ化手順に正しく従うことが難しくなります。

- OSX では、\_\_MACOSX などの隠しファイルが zip ファイルの最上位に追加されることがあります。
- Windows では、ファイルが zip 内のフォルダに自動的に追加され、意図せずにネストされたフォルダが作成されることがあります。

以下の手順は、Info-Zip zip および zipinfo ユーティリティのディストリビューションを提供する Amazon Linux 2 または類似の OS でファイル进行操作していることを前提としています。予期しない動作を防ぐため、これらのツールを使用することをお勧めします。

Ray で使用する Python ファイルをパッケージ化するには

1. パッケージ名で一時ディレクトリを作成し、作業ディレクトリが親ディレクトリであることを確認します。これを行うには、次のコマンドを使用します。

```
cd parent_directory
mkdir temp_dir
```

2. ファイルを一時ディレクトリにコピーし、ディレクトリ構造を確認します。このディレクトリの内容は、Python モジュールとして直接アクセスされます。これを行うには、次のコマンドを使用します。

```
ls -AR temp_dir
my_file_1.py
my_file_2.py
```

3. zip を使用して一時フォルダを圧縮します。これを行うには、次のコマンドを使用します。

```
zip -r zip_file.zip temp_dir
```

4. ファイルが適切にパッケージ化されていることを確認します。*zip\_file.zip* は作業ディレクトリにあるはずですが、次のコマンドで検査できます。

```
zipinfo -l zip_file.zip
temp_dir/
temp_dir/my_file_1.py
temp_dir/my_file_2.py
```

Python パッケージを Ray で使用できるように再パッケージ化するには。

1. パッケージ名で一時ディレクトリを作成し、作業ディレクトリが親ディレクトリであることを確認します。これを行うには、次のコマンドを使用します。

```
cd parent_directory
mkdir temp_dir
```

2. パッケージを解凍し、その内容を一時ディレクトリにコピーします。以前のパッケージング標準に関連するファイルを削除し、モジュールの内容のみを残します。以下のコマンドを実行して、ファイル構造が正しいことを確認します。

```
ls -AR temp_dir
my_module
my_module/__init__.py
my_module/my_file_1.py
my_module/my_submodule/__init__.py
my_module/my_submodule/my_file_2.py
my_module/my_submodule/my_file_3.py
```

3. `zip` を使用して一時フォルダを圧縮します。これを行うには、次のコマンドを使用します。

```
zip -r zip_file.zip temp_dir
```

4. ファイルが適切にパッケージ化されていることを確認します。*zip\_file.zip* は作業ディレクトリにあるはずですが、次のコマンドで検査できます。

```
zipinfo -l zip_file.zip
temp_dir/my_module/
temp_dir/my_module/__init__.py
temp_dir/my_module/my_file_1.py
temp_dir/my_module/my_submodule/
temp_dir/my_module/my_submodule/__init__.py
temp_dir/my_module/my_submodule/my_file_2.py
temp_dir/my_module/my_submodule/my_file_3.py
```

## Ray ジョブのデータに接続する

AWS Glue Ray ジョブでは、データをすばやく統合するために設計されたさまざまな Python パッケージを使用できます。環境が混乱した状態にならないよう、最低限の依存関係のみが提供されています。デフォルトで含まれる内容の詳細については、「[the section called “Ray ジョブで提供されるモジュール”](#)」を参照してください。

### Note

AWS Glue 抽出、変換、ロード (ETL) は、DynamicFrame の抽象化を提供し、データセット内の行間におけるスキーマの違いを解決する ETL ワークフローを効率化します。AWS GlueETL には、ジョブのブックマークや入力ファイルのグループ化など、追加機能があります。Ray ジョブには現在、対応する機能は提供されていません。

AWS Glue for Spark では、特定のデータ形式、ソース、シンクへの接続を直接サポートしています。Ray では、AWS SDK for pandas と現在のサードパーティーライブラリがそれらのニーズを実質的にカバーしています。どのような機能を利用できるか理解するには、これらのライブラリを調べる必要があります。

AWS Glue for Ray の Amazon VPC との統合は、現在ご利用いただけません。Amazon VPC 内のリソースには、パブリックルートがないとアクセスできません。Amazon VPC での AWS Glue の使用に関する詳細は、「[the section called “VPC エンドポイント \(AWS PrivateLink\)”](#)」を参照してください。

## Ray でデータを使用するための共通ライブラリ

Ray Data — Ray Data は、一般的なデータ形式、ソース、シンクを処理する方法を提供します。Ray Data でサポートされている形式とソースの詳細については、Ray Data のドキュメントの「[Input/](#)

[Output](#)」を参照してください。Ray Data は、データセットを扱うための汎用のライブラリではなく特化型のライブラリです。

Ray では、Ray Data がユーザーの作業に最適なソリューションとなるユースケースについて、特定のガイダンスを提供します。詳細については、Ray のドキュメントの「[Ray の使用例](#)」を参照してください。

AWS SDK for pandas (aws wrangler) – AWS SDK for pandas は、pandas の DataFrames を使用して、変換によりデータを管理する際に、AWS サービスでの読み取りと書き込みを行うためのクリーンなテスト済みのソリューションを提供する AWS の製品です。AWS SDK for pandas でサポートされている形式とソースの詳細については、AWS SDK for pandas のドキュメントの「[API Reference](#)」を参照してください。

AWS SDK for pandas を使用してデータを読み書きする方法の例については、AWS SDK for pandas のドキュメントの「[Quick Start](#)」を参照してください。AWS SDK for pandas では、データの変換は行いません。ソースからの読み取りと書き込みのみサポートしています。

Modin — Modin は、一般的な pandas のオペレーションを配布可能な方法で実装している Python ライブラリです。Modin の詳細については、[Modin のドキュメント](#)を参照してください。Modin 自体は、ソースからの読み取りと書き込みをサポートしていません。一般的な変換の分散実装を行います。Modin は AWS SDK for pandas でサポートされています。

Modin を実行し、AWS SDK for pandas と共に Ray 環境で使用すると、高いパフォーマンスで一般的な ETL タスクを実行できます。Modin を AWS SDK for pandas と併用する方法の詳細については、AWS SDK for pandas のドキュメントの「[At scale](#)」を参照してください。

その他のフレームワーク — Ray がサポートしているフレームワークの詳細については、Ray のドキュメントの「[Ray エコシステム](#)」を参照してください。AWS Glue for Ray はその他のフレームワークをサポートしていません。

## データカタログを介したデータへの接続

Ray ジョブと組み合わせたデータカタログによるデータ管理は、AWS SDK for pandas でサポートされています。詳細については、AWS SDK for pandas ウェブサイトの「[Glue Catalog](#)」(Glue カタログ)を参照してください。

## AWS SDK でこのサービスを使用する

AWS Software Development Kit (SDKs)は、多くの一般的なプログラミング言語で使用できます。各 SDK には、デベロッパーが好みの言語でアプリケーションを簡単に構築できるようにする API、コード例、およびドキュメントが提供されています。

SDK ドキュメント	コード例
<a href="#">AWS SDK for C++</a>	<a href="#">AWS SDK for C++ コード例</a>
<a href="#">AWS CLI</a>	<a href="#">AWS CLI コード例</a>
<a href="#">AWS SDK for Go</a>	<a href="#">AWS SDK for Go コード例</a>
<a href="#">AWS SDK for Java</a>	<a href="#">AWS SDK for Java コード例</a>
<a href="#">AWS SDK for JavaScript</a>	<a href="#">AWS SDK for JavaScript コード例</a>
<a href="#">AWS SDK for Kotlin</a>	<a href="#">AWS SDK for Kotlin コード例</a>
<a href="#">AWS SDK for .NET</a>	<a href="#">AWS SDK for .NET コード例</a>
<a href="#">AWS SDK for PHP</a>	<a href="#">AWS SDK for PHP コード例</a>
<a href="#">AWS Tools for PowerShell</a>	<a href="#">PowerShell コード例のツール</a>
<a href="#">AWS SDK for Python (Boto3)</a>	<a href="#">AWS SDK for Python (Boto3) コード例</a>
<a href="#">AWS SDK for Ruby</a>	<a href="#">AWS SDK for Ruby コード例</a>
<a href="#">AWS SDK for Rust</a>	<a href="#">AWS SDK for Rust コード例</a>
<a href="#">AWS SDK for SAP ABAP</a>	<a href="#">AWS SDK for SAP ABAP コード例</a>
<a href="#">AWS SDK for Swift</a>	<a href="#">AWS SDK for Swift コード例</a>

このサービスに固有の例については、「[AWS GlueAWS SDKs を使用した API コード例](#)」を参照してください。

**i** 可用性の例

必要なものが見つからなかった場合。このページの下側にある [Provide feedback (フィードバックを送信)] リンクから、コードの例をリクエストしてください。

# AWS Glue API

このセクションでは、AWS Glue SDKs とツールで使用されるデータ型とプリミティブについて説明します。の外部で AWS Glue プログラムで を操作するには、次の 3 つの一般的な方法があり AWS Management Console、それぞれに独自のドキュメントがあります。

- 言語 SDK ライブラリを使用すれば、一般的なプログラミング言語で AWS リソースにアクセスできます。詳細については、「[AWSでの構築ツール](#)」をご覧ください。
- AWS CLI では、コマンドラインから AWS リソースにアクセスできます。詳細については、「[AWS CLI コマンドリファレンス](#)」を参照してください。
- AWS CloudFormation では、一貫してプロビジョニングする一連の AWS リソースを定義できます。詳細については、[AWS CloudFormation「:resource type reference AWS Glue」](#)を参照してください。

このセクションでは、これらの SDK およびツールから独立した共有プリミティブについて説明します。ツールでは、[AWS Glue ウェブ API リファレンス](#)を使用して と通信します AWS。

## 目次

- [AWS Glue のセキュリティ API](#)
  - [データ型](#)
  - [DataCatalogEncryptionSettings 構造](#)
  - [EncryptionAtRest 構造](#)
  - [ConnectionPasswordEncryption 構造](#)
  - [EncryptionConfiguration 構造](#)
  - [S3Encryption 構造](#)
  - [CloudWatchEncryption 構造](#)
  - [JobBookmarksEncryption 構造](#)
  - [SecurityConfiguration 構造](#)
  - [GluePolicy の構造](#)
  - [操作](#)
  - [GetDataCatalogEncryptionSettings アクション \(Python: `get\_data\_catalog\_encryption\_settings`\)](#)
  - [PutDataCatalogEncryptionSettings アクション \(Python: `put\_data\_catalog\_encryption\_settings`\)](#)
  - [PutResourcePolicy アクション \(Python: `put\_resource\_policy`\)](#)

- [GetResourcePolicy アクション \(Python: get\\_resource\\_policy\)](#)
- [DeleteResourcePolicy アクション \(Python: delete\\_resource\\_policy\)](#)
- [CreateSecurityConfiguration アクション \(Python: create\\_security\\_configuration\)](#)
- [DeleteSecurityConfiguration アクション \(Python: delete\\_security\\_configuration\)](#)
- [GetSecurityConfiguration アクション \(Python: get\\_security\\_configuration\)](#)
- [GetSecurityConfigurations アクション \(Python: get\\_security\\_configurations\)](#)
- [GetResourcePolicies アクション \(Python: get\\_resource\\_policies\)](#)
- [Catalog API](#)
  - [データベース API](#)
    - [データ型](#)
    - [データベース構造](#)
    - [DatabaseInput 構造](#)
    - [PrincipalPermissions 構造](#)
    - [DataLakePrincipal 構造](#)
    - [DatabaseIdentifier 構造](#)
    - [FederatedDatabase 構造](#)
    - [操作](#)
    - [CreateDatabase アクション \(Python: create\\_database\)](#)
    - [UpdateDatabase アクション \(Python: update\\_database\)](#)
    - [DeleteDatabase アクション \(Python: delete\\_database\)](#)
    - [GetDatabase アクション \(Python: get\\_database\)](#)
    - [GetDatabases アクション \(Python: get\\_databases\)](#)
  - [Table API](#)
    - [データ型](#)
    - [Table 構造](#)
    - [TableInput 構造](#)
    - [FederatedTable 構造](#)
    - [列の構造](#)
    - [StorageDescriptor 構造](#)
    - [SchemaReference 構造](#)

- [SerDeInfo 構造](#)
- [Order 構造](#)
- [SkewedInfo 構造](#)
- [TableVersion 構造](#)
- [TableError 構造](#)
- [TableVersionError 構造](#)
- [SortCriterion 構造](#)
- [TableIdentifier 構造](#)
- [KeySchemaElement 構造](#)
- [PartitionIndex 構造](#)
- [PartitionIndexDescriptor 構造](#)
- [BackfillError 構造](#)
- [IcebergInput 構造](#)
- [OpenTableFormatInput 構造](#)
- [ViewDefinition 構造](#)
- [ViewDefinitionInput 構造](#)
- [ViewRepresentation 構造](#)
- [ViewRepresentationInput 構造](#)
- [操作](#)
- [CreateTable アクション \(Python: create\\_table\)](#)
- [UpdateTable アクション \(Python: update\\_table\)](#)
- [DeleteTable アクション \(Python: delete\\_table\)](#)
- [BatchDeleteTable アクション \(Python: batch\\_delete\\_table\)](#)
- [GetTable アクション \(Python: get\\_table\)](#)
- [GetTables アクション \(Python: get\\_tables\)](#)
- [GetTableVersion アクション \(Python: get\\_table\\_version\)](#)
- [GetTableVersions アクション \(Python: get\\_table\\_versions\)](#)
- [DeleteTableVersion アクション \(Python: delete\\_table\\_version\)](#)
- [BatchDeleteTableVersion アクション \(Python: batch\\_delete\\_table\\_version\)](#)
- [SearchTables アクション \(Python: search\\_tables\)](#)

- [GetPartitionIndexes アクション \(Python: get\\_partition\\_indexes\)](#)
- [CreatePartitionIndex アクション \(Python: create\\_partition\\_index\)](#)
- [DeletePartitionIndex アクション \(Python: delete\\_partition\\_index\)](#)
- [GetColumnStatisticsForTable アクション \(Python: get\\_column\\_statistics\\_for\\_table\)](#)
- [UpdateColumnStatisticsForTable アクション \(Python: update\\_column\\_statistics\\_for\\_table\)](#)
- [DeleteColumnStatisticsForTable アクション \(Python: delete\\_column\\_statistics\\_for\\_table\)](#)
- [パーティション API](#)
  - [データ型](#)
  - [パーティションの構造](#)
  - [PartitionInput の構造](#)
  - [PartitionSpecWithSharedStorageDescriptor 構造](#)
  - [PartitionListComposingSpec 構造](#)
  - [PartitionSpecProxy 構造](#)
  - [PartitionValueList 構造](#)
  - [セグメント構造](#)
  - [PartitionError 構造](#)
  - [BatchUpdatePartitionFailureEntry 構造](#)
  - [BatchUpdatePartitionRequestEntry 構造](#)
  - [StorageDescriptor 構造](#)
  - [SchemaReference 構造](#)
  - [SerDelInfo 構造](#)
  - [SkewedInfo 構造](#)
  - [操作](#)
  - [CreatePartition アクション \(Python: create\\_partition\)](#)
  - [BatchCreatePartition アクション \(Python: batch\\_create\\_partition\)](#)
  - [UpdatePartition アクション \(Python: update\\_partition\)](#)
  - [DeletePartition アクション \(Python: delete\\_partition\)](#)
  - [BatchDeletePartition アクション \(Python: batch\\_delete\\_partition\)](#)
  - [GetPartition アクション \(Python: get\\_partition\)](#)
  - [GetPartitions アクション \(Python: get\\_partitions\)](#)

- [BatchGetPartition アクション \(Python: batch\\_get\\_partition\)](#)
- [BatchUpdatePartition アクション \(Python: batch\\_update\\_partition\)](#)
- [GetColumnStatisticsForPartition アクション \(Python: get\\_column\\_statistics\\_for\\_partition\)](#)
- [UpdateColumnStatisticsForPartition アクション \(Python: update\\_column\\_statistics\\_for\\_partition\)](#)
- [DeleteColumnStatisticsForPartition アクション \(Python: delete\\_column\\_statistics\\_for\\_partition\)](#)
- [接続 API](#)
  - [データ型](#)
  - [Connection 構造](#)
  - [ConnectionInput 構造](#)
  - [PhysicalConnectionRequirements 構造](#)
  - [GetConnectionsFilter 構造](#)
  - [操作](#)
  - [CreateConnection アクション \(Python: create\\_connection\)](#)
  - [DeleteConnection アクション \(Python: delete\\_connection\)](#)
  - [GetConnection アクション \(Python: get\\_connection\)](#)
  - [GetConnections アクション \(Python: get\\_connections\)](#)
  - [UpdateConnection アクション \(Python: update\\_connection\)](#)
  - [BatchDeleteConnection アクション \(Python: batch\\_delete\\_connection\)](#)
  - [認証の設定](#)
  - [AuthenticationConfiguration 構造](#)
  - [AuthenticationConfigurationInput 構造](#)
  - [OAuth2Properties 構造](#)
  - [OAuth2PropertiesInput 構造](#)
  - [OAuth2ClientApplication 構造](#)
  - [AuthorizationCodeProperties 構造](#)
- [ユーザー定義関数 API](#)
  - [データ型](#)
  - [UserDefinedFunction 構造](#)
  - [UserDefinedFunctionInput 構造](#)

- [操作](#)
- [CreateUserDefinedFunction アクション \(Python: create\\_user\\_defined\\_function\)](#)
- [UpdateUserDefinedFunction アクション \(Python: update\\_user\\_defined\\_function\)](#)
- [DeleteUserDefinedFunction アクション \(Python: delete\\_user\\_defined\\_function\)](#)
- [GetUserDefinedFunction アクション \(Python: get\\_user\\_defined\\_function\)](#)
- [GetUserDefinedFunctions アクション \(Python: get\\_user\\_defined\\_functions\)](#)
- [Athena カタログを AWS Glue にインポートする](#)
  - [データ型](#)
  - [CatalogImportStatus 構造](#)
  - [操作](#)
  - [ImportCatalogToGlue アクション \(Python: import\\_catalog\\_to\\_glue\)](#)
  - [GetCatalogImportStatus アクション \(Python: get\\_catalog\\_import\\_status\)](#)
- [テーブルオプティマイザー API](#)
  - [データ型](#)
  - [TableOptimizer 構造](#)
  - [TableOptimizerConfiguration 構造](#)
  - [TableOptimizerRun 構造](#)
  - [RunMetrics 構造](#)
  - [BatchGetTableOptimizerEntry 構造](#)
  - [BatchTableOptimizer 構造](#)
  - [BatchGetTableOptimizerError 構造](#)
  - [操作](#)
  - [GetTableOptimizer アクション \(Python: get\\_table\\_optimizer\)](#)
  - [BatchGetTableOptimizer アクション \(Python: batch\\_get\\_table\\_optimizer\)](#)
  - [ListTableOptimizerRuns アクション \(Python: list\\_table\\_optimizer\\_runs\)](#)
  - [CreateTableOptimizer アクション \(Python: create\\_table\\_optimizer\)](#)
  - [DeleteTableOptimizer アクション \(Python: delete\\_table\\_optimizer\)](#)
  - [UpdateTableOptimizer アクション \(Python: update\\_table\\_optimizer\)](#)
- [クローラーおよび分類子 API](#)
  - [分類子 API](#)

- [データ型](#)
- [分類子の構造](#)
- [GrokClassifier の構造](#)
- [XMLClassifier の構造](#)
- [JsonClassifier の構造](#)
- [CsvClassifier の構造](#)
- [CreateGrokClassifierRequest の構造](#)
- [UpdateGrokClassifierRequest の構造](#)
- [CreateXMLClassifierRequest の構造](#)
- [UpdateXMLClassifierRequest の構造](#)
- [CreateJsonClassifierRequest の構造](#)
- [UpdateJsonClassifierRequest の構造](#)
- [CreateCsvClassifierRequest の構造](#)
- [UpdateCsvClassifierRequest の構造](#)
- [操作](#)
- [CreateClassifier アクション \(Python: create\\_classifier\)](#)
- [DeleteClassifier アクション \(Python: delete\\_classifier\)](#)
- [GetClassifier アクション \(Python: get\\_classifier\)](#)
- [GetClassifiers アクション \(Python: get\\_classifiers\)](#)
- [UpdateClassifier アクション \(Python: update\\_classifier\)](#)
- [クローラー API](#)
  - [データ型](#)
  - [Crawler 構造](#)
  - [Schedule 構造](#)
  - [CrawlerTargets 構造](#)
  - [S3Target 構造](#)
  - [S3DeltaCatalogTarget 構造](#)
  - [S3DeltaDirectTarget 構造](#)
  - [JdbcTarget 構造](#)
  - [MongoDBTarget 構造](#)

- [DynamoDBTarget 構造](#)
- [DeltaTarget 構造](#)
- [IcebergTarget 構造](#)
- [HudiTarget 構造](#)
- [CatalogTarget 構造](#)
- [CrawlerMetrics 構造](#)
- [CrawlerHistory 構造](#)
- [CrawlsFilter 構造](#)
- [SchemaChangePolicy 構造](#)
- [LastCrawlInfo 構造](#)
- [RecrawlPolicy 構造](#)
- [LineageConfiguration 構造](#)
- [LakeFormationConfiguration 構造](#)
- [操作](#)
- [CreateCrawler アクション \(Python: create\\_crawler\)](#)
- [DeleteCrawler アクション \(Python: delete\\_crawler\)](#)
- [GetCrawler アクション \(Python: get\\_crawler\)](#)
- [GetCrawlers アクション \(Python: get\\_crawlers\)](#)
- [GetCrawlerMetrics アクション \(Python: get\\_crawler\\_metrics\)](#)
- [UpdateCrawler アクション \(Python: update\\_crawler\)](#)
- [StartCrawler アクション \(Python: start\\_crawler\)](#)
- [StopCrawler アクション \(Python: stop\\_crawler\)](#)
- [BatchGetCrawlers アクション \(Python: batch\\_get\\_crawlers\)](#)
- [ListCrawlers アクション \(Python: list\\_crawlers\)](#)
- [ListCrawls アクション \(Python: list\\_crawls\)](#)
- [列統計 API](#)
  - [データ型](#)
  - [ColumnStatisticsTaskRun の構造](#)
  - [ColumnStatisticsTaskRunningException の構造](#)
  - [ColumnStatisticsTaskNotRunningException の構造](#)

- [ColumnStatisticsTaskStoppingException の構造](#)
- [操作](#)
- [StartColumnStatisticsTaskRun アクション \(Python: start\\_column\\_statistics\\_task\\_run\)](#)
- [GetColumnStatisticsTaskRun アクション \(Python: get\\_column\\_statistics\\_task\\_run\)](#)
- [GetColumnStatisticsTaskRuns アクション \(Python: get\\_column\\_statistics\\_task\\_runs\)](#)
- [ListColumnStatisticsTaskRuns アクション \(Python: list\\_column\\_statistics\\_task\\_runs\)](#)
- [StopColumnStatisticsTaskRun アクション \(Python: stop\\_column\\_statistics\\_task\\_run\)](#)
- [クローラスケジューラ API](#)
  - [データ型](#)
  - [Schedule 構造](#)
  - [操作](#)
  - [UpdateCrawlerSchedule アクション \(Python: update\\_crawler\\_schedule\)](#)
  - [StartCrawlerSchedule アクション \(Python: start\\_crawler\\_schedule\)](#)
  - [StopCrawlerSchedule アクション \(Python: stop\\_crawler\\_schedule\)](#)
- [ETL スクリプト API の自動生成](#)
  - [データ型](#)
  - [CodeGenNode 構造](#)
  - [CodeGenNodeArg 構造](#)
  - [CodeGenEdge 構造](#)
  - [場所の構造](#)
  - [CatalogEntry 構造](#)
  - [MappingEntry 構造](#)
  - [操作](#)
  - [CreateScript アクション \(Python: create\\_script\)](#)
  - [GetDataflowGraph アクション \(Python: get\\_dataflow\\_graph\)](#)
  - [GetMapping アクション \(Python: get\\_mapping\)](#)
  - [GetPlan アクション \(Python: get\\_plan\)](#)
- [ビジュアルジョブ API](#)
  - [データ型](#)
  - [CodeGenConfigurationNode 構造](#)

- [JDBC ConnectorOptions 構造](#)
- [StreamingDataPreviewOptions 構造](#)
- [AthenaConnectorSource 構造](#)
- [JDBC ConnectorSource 構造](#)
- [SparkConnectorSource 構造](#)
- [CatalogSource 構造](#)
- [MySQLCatalogSource 構造](#)
- [PostgreSQLCatalogSource 構造](#)
- [OracleSQLCatalogSource 構造](#)
- [MicrosoftSQLServerCatalogSource 構造](#)
- [CatalogKinesisSource 構造](#)
- [DirectKinesisSource 構造](#)
- [KinesisStreamingSourceOptions 構造](#)
- [CatalogKafkaSource 構造](#)
- [DirectKafkaSource 構造](#)
- [KafkaStreamingSourceOptions 構造](#)
- [RedshiftSource 構造](#)
- [AmazonRedshiftSource 構造](#)
- [AmazonRedshiftNodeData 構造](#)
- [AmazonRedshiftAdvancedOption 構造](#)
- [Option 構造](#)
- [S3CatalogSource 構造](#)
- [S3SourceAdditionalOptions 構造](#)
- [S3CsvSource 構造](#)
- [DirectJDBCSource 構造](#)
- [S3DirectSourceAdditionalOptions 構造](#)
- [S3JsonSource 構造](#)
- [S3ParquetSource 構造](#)
- [S3DeltaSource 構造](#)
- [S3CatalogDeltaSource 構造](#)

- [CatalogDeltaSource 構造](#)
- [S3HudiSource 構造](#)
- [S3CatalogHudiSource 構造](#)
- [CatalogHudiSource 構造](#)
- [DynamoDBCatalogSource 構造](#)
- [RelationalCatalogSource 構造](#)
- [JDBC ConnectorTarget 構造](#)
- [SparkConnectorTarget 構造](#)
- [BasicCatalogTarget 構造](#)
- [MySQLCatalogTarget 構造](#)
- [PostgreSQLCatalogTarget 構造](#)
- [OracleSQLCatalogTarget 構造](#)
- [MicrosoftSQLServerCatalogTarget 構造](#)
- [RedshiftTarget 構造](#)
- [AmazonRedshiftTarget 構造](#)
- [UpsertRedshiftTargetOptions 構造](#)
- [S3CatalogTarget 構造](#)
- [S3GlueParquetTarget 構造](#)
- [CatalogSchemaChangePolicy 構造](#)
- [S3DirectTarget 構造](#)
- [S3HudiCatalogTarget 構造](#)
- [S3HudiDirectTarget 構造](#)
- [S3DeltaCatalogTarget 構造](#)
- [S3DeltaDirectTarget 構造](#)
- [DirectSchemaChangePolicy 構造](#)
- [ApplyMapping 構造](#)
- [Mapping 構造](#)
- [SelectFields 構造](#)
- [DropFields 構造](#)
- [RenameField 構造](#)

- [スピゴット構造](#)
- [Join 構造](#)
- [JoinColumn 構造](#)
- [SplitFields 構造](#)
- [SelectFromCollection 構造](#)
- [FillMissingValues 構造](#)
- [Filter 構造](#)
- [FilterExpression 構造](#)
- [FilterValue 構造](#)
- [CustomCode 構造](#)
- [SparkSQL 構造](#)
- [SqlAlias 構造](#)
- [DropNullFields 構造](#)
- [NullCheckBoxList 構造](#)
- [NullValueField 構造](#)
- [Datatype 構造](#)
- [Merge 構造](#)
- [Union 構造](#)
- [PIIDetection 構造](#)
- [Aggregate 構造](#)
- [DropDuplicates 構造](#)
- [GovernedCatalogTarget 構造](#)
- [GovernedCatalogSource 構造](#)
- [AggregateOperation 構造](#)
- [GlueSchema 構造](#)
- [GlueStudioSchemaColumn 構造](#)
- [GlueStudioColumn 構造](#)
- [DynamicTransform 構造](#)
- [TransformConfigParameter 構造](#)
- [EvaluateDataQuality 構造](#)

- [DQ ResultsPublishingOptions 構造](#)
- [DQ StopJobOnFailureOptions 構造](#)
- [EvaluateDataQualityMultiFrame 構造](#)
- [Recipe 構造](#)
- [RecipeReference 構造](#)
- [SnowflakeNodeData 構造](#)
- [SnowflakeSource 構造](#)
- [SnowflakeTarget 構造](#)
- [ConnectorDataSource 構造](#)
- [ConnectorDataTarget 構造](#)
- [ジョブ API](#)
  - [ジョブ](#)
    - [データ型](#)
    - [Job 構造](#)
    - [ExecutionProperty 構造](#)
    - [NotificationProperty 構造](#)
    - [JobCommand 構造](#)
    - [ConnectionsList 構造](#)
    - [JobUpdate 構造](#)
    - [SourceControlDetails 構造](#)
    - [操作](#)
    - [CreateJob アクション \(Python: create\\_job\)](#)
    - [UpdateJob アクション \(Python: update\\_job\)](#)
    - [GetJob アクション \(Python: get\\_job\)](#)
    - [GetJobs アクション \(Python: get\\_jobs\)](#)
    - [DeleteJob アクション \(Python: delete\\_job\)](#)
    - [ListJobs アクション \(Python: list\\_jobs\)](#)
    - [BatchGetJobs アクション \(Python: batch\\_get\\_jobs\)](#)
- [ジョブ実行](#)
  - [データ型](#)

- [JobRun 構造](#)
- [Predecessor 構造](#)
- [JobBookmarkEntry 構造](#)
- [BatchStopJobRunSuccessfulSubmission 構造](#)
- [BatchStopJobRunError 構造](#)
- [NotificationProperty 構造](#)
- [操作](#)
- [StartJobRun アクション \(Python: start\\_job\\_run\)](#)
- [BatchStopJobRun アクション \(Python: batch\\_stop\\_job\\_run\)](#)
- [GetJobRun アクション \(Python: get\\_job\\_run\)](#)
- [GetJobRuns アクション \(Python: get\\_job\\_runs\)](#)
- [GetJobBookmark アクション \(Python: get\\_job\\_bookmark\)](#)
- [GetJobBookmarks アクション \(Python: get\\_job\\_bookmarks\)](#)
- [ResetJobBookmark アクション \(Python: reset\\_job\\_bookmark\)](#)
- [トリガー](#)
  - [データ型](#)
  - [トリガー構造](#)
  - [TriggerUpdate 構造](#)
  - [述語構造](#)
  - [条件の構造](#)
  - [アクション構造](#)
  - [EventBatchingCondition 構造](#)
  - [操作](#)
  - [CreateTrigger アクション \(Python: create\\_trigger\)](#)
  - [StartTrigger アクション \(Python: start\\_trigger\)](#)
  - [GetTrigger アクション \(Python: get\\_trigger\)](#)
  - [GetTriggers アクション \(Python: get\\_triggers\)](#)
  - [UpdateTrigger アクション \(Python: update\\_trigger\)](#)
  - [StopTrigger アクション \(Python: stop\\_trigger\)](#)
  - [DeleteTrigger アクション \(Python: delete\\_trigger\)](#)

- [ListTriggers アクション \(Python: list\\_triggers\)](#)
- [BatchGetTriggers アクション \(Python: batch\\_get\\_triggers\)](#)
- [インタラクティブセッション API](#)
  - [データ型](#)
  - [セッション構造](#)
  - [SessionCommand 構造](#)
  - [Statement 構造](#)
  - [StatementOutput 構造](#)
  - [StatementOutputData 構造](#)
  - [ConnectionsList 構造](#)
  - [操作](#)
  - [CreateSession アクション \(Python: create\\_session\)](#)
  - [StopSession アクション \(Python: stop\\_session\)](#)
  - [DeleteSession アクション \(Python: delete\\_session\)](#)
  - [GetSession アクション \(Python: get\\_session\)](#)
  - [ListSessions アクション \(Python: list\\_sessions\)](#)
  - [RunStatement アクション \(Python: run\\_statement\)](#)
  - [CancelStatement アクション \(Python: cancel\\_statement\)](#)
  - [GetStatement アクション \(Python: get\\_statement\)](#)
  - [ListStatements アクション \(Python: list\\_statements\)](#)
- [開発エンドポイント API](#)
  - [データ型](#)
  - [DevEndpoint 構造](#)
  - [DevEndpointCustomLibraries 構造](#)
  - [操作](#)
  - [CreateDevEndpoint アクション \(Python: create\\_dev\\_endpoint\)](#)
  - [UpdateDevEndpoint アクション \(Python: update\\_dev\\_endpoint\)](#)
  - [DeleteDevEndpoint アクション \(Python: delete\\_dev\\_endpoint\)](#)
  - [GetDevEndpoint アクション \(Python: get\\_dev\\_endpoint\)](#)
  - [GetDevEndpoints アクション \(Python: get\\_dev\\_endpoints\)](#)

- [BatchGetDevEndpoints アクション \(Python: batch\\_get\\_dev\\_endpoints\)](#)
- [ListDevEndpoints アクション \(Python: list\\_dev\\_endpoints\)](#)
- [スキーマレジストリ](#)
  - [データ型](#)
  - [RegistryId 構造](#)
  - [RegistryListItem 構造](#)
  - [MetadataInfo 構造](#)
  - [OtherMetadataValueListItem 構造](#)
  - [SchemaListItem 構造](#)
  - [SchemaVersionListItem 構造](#)
  - [MetadataKeyValuePair 構造](#)
  - [SchemaVersionErrorItem 構造](#)
  - [ErrorDetails 構造](#)
  - [SchemaVersionNumber 構造](#)
  - [Schemald 構造](#)
  - [操作](#)
  - [CreateRegistry アクション \(Python: create\\_registry\)](#)
  - [CreateSchema アクション \(Python: create\\_schema\)](#)
  - [GetSchema アクション \(Python: get\\_schema\)](#)
  - [ListSchemaVersions アクション \(Python: list\\_schema\\_versions\)](#)
  - [GetSchemaVersion アクション \(Python: get\\_schema\\_version\)](#)
  - [GetSchemaVersionsDiff アクション \(Python: get\\_schema\\_versions\\_diff\)](#)
  - [ListRegistries アクション \(Python: list\\_registries\)](#)
  - [ListSchemas アクション \(Python: list\\_schemas\)](#)
  - [RegisterSchemaVersion アクション \(Python: register\\_schema\\_version\)](#)
  - [UpdateSchema アクション \(Python: update\\_schema\)](#)
  - [CheckSchemaVersionValidity アクション \(Python: check\\_schema\\_version\\_validity\)](#)
  - [UpdateRegistry アクション \(Python: update\\_registry\)](#)
  - [GetSchemaByDefinition アクション \(Python: get\\_schema\\_by\\_definition\)](#)
- [GetRegistry アクション \(Python: get\\_registry\)](#)

- [PutSchemaVersionMetadata アクション \(Python: put\\_schema\\_version\\_metadata\)](#)
- [QuerySchemaVersionMetadata アクション \(Python: query\\_schema\\_version\\_metadata\)](#)
- [RemoveSchemaVersionMetadata アクション \(Python: remove\\_schema\\_version\\_metadata\)](#)
- [DeleteRegistry アクション \(Python: delete\\_registry\)](#)
- [DeleteSchema アクション \(Python: delete\\_schema\)](#)
- [DeleteSchemaVersions アクション \(Python: delete\\_schema\\_versions\)](#)
- [ワークフロー](#)
  - [データ型](#)
  - [JobNodeDetails の構造](#)
  - [CrawlerNodeDetails 構造](#)
  - [TriggerNodeDetails 構造](#)
  - [Crawl 構造](#)
  - [Node 構造](#)
  - [Edge 構造](#)
  - [Workflow 構造](#)
  - [WorkflowGraph 構造](#)
  - [WorkflowRun 構造](#)
  - [WorkflowRunStatistics 構造](#)
  - [StartingEventBatchCondition 構造](#)
  - [Blueprint 構造](#)
  - [BlueprintDetails 構造](#)
  - [LastActiveDefinition 構造](#)
  - [BlueprintRun 構造](#)
  - [操作](#)
  - [CreateWorkflow アクション \(Python: create\\_workflow\)](#)
  - [UpdateWorkflow アクション \(Python: update\\_workflow\)](#)
  - [DeleteWorkflow アクション \(Python: delete\\_workflow\)](#)
  - [GetWorkflow アクション \(Python: get\\_workflow\)](#)
  - [ListWorkflows アクション \(Python: list\\_workflows\)](#)
- [BatchGetWorkflows アクション \(Python: batch\\_get\\_workflows\)](#)

- [GetWorkflowRun アクション \(Python: get\\_workflow\\_run\)](#)
- [GetWorkflowRuns アクション \(Python: get\\_workflow\\_runs\)](#)
- [GetWorkflowRunProperties アクション \(Python: get\\_workflow\\_run\\_properties\)](#)
- [PutWorkflowRunProperties アクション \(Python: put\\_workflow\\_run\\_properties\)](#)
- [CreateBlueprint アクション \(Python: create\\_blueprint\)](#)
- [UpdateBlueprint アクション \(Python: update\\_blueprint\)](#)
- [DeleteBlueprint アクション \(Python: delete\\_blueprint\)](#)
- [ListBlueprints アクション \(Python: list\\_blueprints\)](#)
- [BatchGetBlueprints アクション \(Python: batch\\_get\\_blueprints\)](#)
- [StartBlueprintRun アクション \(Python: start\\_blueprint\\_run\)](#)
- [GetBlueprintRun アクション \(Python: get\\_blueprint\\_run\)](#)
- [GetBlueprintRuns アクション \(Python: get\\_blueprint\\_runs\)](#)
- [StartWorkflowRun アクション \(Python: start\\_workflow\\_run\)](#)
- [StopWorkflowRun アクション \(Python: stop\\_workflow\\_run\)](#)
- [ResumeWorkflowRun アクション \(Python: resume\\_workflow\\_run\)](#)
- [使用プロファイル](#)
  - [データ型](#)
  - [ProfileConfiguration 構造](#)
  - [ConfigurationObject 構造](#)
  - [UsageProfileDefinition 構造](#)
  - [操作](#)
  - [CreateUsageProfile アクション \(Python: create\\_usage\\_profile\)](#)
  - [GetUsageProfile アクション \(Python: get\\_usage\\_profile\)](#)
  - [UpdateUsageProfile アクション \(Python: update\\_usage\\_profile\)](#)
  - [DeleteUsageProfile アクション \(Python: delete\\_usage\\_profile\)](#)
  - [ListUsageProfiles アクション \(Python: list\\_usage\\_profiles\)](#)
- [機械学習 API](#)
  - [データ型](#)
  - [TransformParameters 構造](#)
  - [EvaluationMetrics 構造](#)

- [MLTransform 構造](#)
- [FindMatchesParameters 構造](#)
- [FindMatchesMetrics 構造](#)
- [ConfusionMatrix 構造](#)
- [GlueTable 構造](#)
- [TaskRun 構造](#)
- [TransformFilterCriteria 構造](#)
- [TransformSortCriteria 構造](#)
- [TaskRunFilterCriteria 構造](#)
- [TaskRunSortCriteria 構造](#)
- [TaskRunProperties 構造](#)
- [FindMatchesTaskRunProperties 構造](#)
- [ImportLabelsTaskRunProperties 構造](#)
- [ExportLabelsTaskRunProperties 構造](#)
- [LabelingSetGenerationTaskRunProperties 構造](#)
- [SchemaColumn 構造](#)
- [TransformEncryption 構造](#)
- [MLUserDataEncryption 構造](#)
- [ColumnImportance 構造](#)
- [操作](#)
- [CreateMLTransform アクション \(Python: create\\_ml\\_transform\)](#)
- [UpdateMLTransform アクション \(Python: update\\_ml\\_transform\)](#)
- [DeleteMLTransform アクション \(Python: delete\\_ml\\_transform\)](#)
- [GetMLTransform アクション \(Python: get\\_ml\\_transform\)](#)
- [GetMLTransforms アクション \(Python: get\\_ml\\_transforms\)](#)
- [ListMLTransforms アクション \(Python: list\\_ml\\_transforms\)](#)
- [StartMLEvaluationTaskRun アクション \(Python: start\\_ml\\_evaluation\\_task\\_run\)](#)
- [StartMLLabelingSetGenerationTaskRun アクション \(Python: start\\_ml\\_labeling\\_set\\_generation\\_task\\_run\)](#)
- [GetMLTaskRun アクション \(Python: get\\_ml\\_task\\_run\)](#)

- [GetMLTaskRuns アクション \(Python: get\\_ml\\_task\\_runs\)](#)
- [CancelMLTaskRun アクション \(Python: cancel\\_ml\\_task\\_run\)](#)
- [StartExportLabelsTaskRun アクション \(Python: start\\_export\\_labels\\_task\\_run\)](#)
- [StartImportLabelsTaskRun アクション \(Python: start\\_import\\_labels\\_task\\_run\)](#)
- [Data Quality API](#)
  - [データ型](#)
  - [DataSource 構造](#)
  - [DataQualityRulesetListDetails 構造](#)
  - [DataQualityTargetTable 構造](#)
  - [DataQualityRulesetEvaluationRunDescription 構造](#)
  - [DataQualityRulesetEvaluationRunFilter 構造](#)
  - [DataQualityEvaluationRunAdditionalRunOptions 構造](#)
  - [DataQualityRuleRecommendationRunDescription 構造](#)
  - [DataQualityRuleRecommendationRunFilter 構造](#)
  - [DataQualityResult 構造](#)
  - [DataQualityAnalyzerResult 構造](#)
  - [DataQualityObservation 構造](#)
  - [MetricBasedObservation 構造](#)
  - [DataQualityMetricValues 構造](#)
  - [DataQualityRuleResult 構造](#)
  - [DataQualityResultDescription 構造](#)
  - [DataQualityResultFilterCriteria 構造](#)
  - [DataQualityRulesetFilterCriteria 構造](#)
  - [操作](#)
  - [StartDataQualityRulesetEvaluationRun アクション \(Python: start\\_data\\_quality\\_ruleset\\_evaluation\\_run \)](#)
  - [CancelDataQualityRulesetEvaluationRun アクション \(Python: cancel\\_data\\_quality\\_ruleset\\_evaluation\\_run \)](#)
  - [GetDataQualityRulesetEvaluationRun アクション \(Python: get\\_data\\_quality\\_ruleset\\_evaluation\\_run\)](#)

- [ListDataQualityRulesetEvaluationRuns アクション \(Python: list\\_data\\_quality\\_ruleset\\_evaluation\\_runs\)](#)
- [StartDataQualityRuleRecommendationRun アクション \(Python: start\\_data\\_quality\\_rule\\_recommendation\\_run\)](#)
- [CancelDataQualityRuleRecommendationRun アクション \(Python: cancel\\_data\\_quality\\_rule\\_recommendation\\_run\)](#)
- [GetDataQualityRuleRecommendationRun アクション \(Python: get\\_data\\_quality\\_rule\\_recommendation\\_run\)](#)
- [ListDataQualityRuleRecommendationRuns アクション \(Python: list\\_data\\_quality\\_rule\\_recommendation\\_runs\)](#)
- [GetDataQualityResult アクション \(Python: get\\_data\\_quality\\_result\)](#)
- [BatchGetDataQualityResult アクション \(Python: batch\\_get\\_data\\_quality\\_result\)](#)
- [ListDataQualityResults アクション \(Python: list\\_data\\_quality\\_results\)](#)
- [CreateDataQualityRuleset アクション \(Python: create\\_data\\_quality\\_ruleset\)](#)
- [DeleteDataQualityRuleset アクション \(Python: delete\\_data\\_quality\\_ruleset\)](#)
- [GetDataQualityRuleset アクション \(Python: get\\_data\\_quality\\_ruleset\)](#)
- [ListDataQualityRulesets アクション \(Python: list\\_data\\_quality\\_rulesets\)](#)
- [UpdateDataQualityRuleset アクション \(Python: update\\_data\\_quality\\_ruleset\)](#)
- [機密データの検出 API](#)
  - [データ型](#)
  - [CustomentityType 構造](#)
  - [操作](#)
  - [CreateCustomEntityType アクション \(Python: create\\_custom\\_entity\\_type\)](#)
  - [DeleteCustomEntityType アクション \(Python: delete\\_custom\\_entity\\_type\)](#)
  - [GetCustomEntityType アクション \(Python: get\\_custom\\_entity\\_type\)](#)
  - [BatchGetCustomEntityTypes アクション \(Python: batch\\_get\\_custom\\_entity\\_types\)](#)
  - [ListCustomEntityTypes アクション \(Python: list\\_custom\\_entity\\_types\)](#)
- [AWS Glue の API のタグ付け](#)
  - [データ型](#)
  - [Tag 構造](#)
  - [操作](#)

- [TagResource アクション \(Python: tag\\_resource\)](#)
- [UntagResource アクション \(Python: untag\\_resource\)](#)
- [GetTags アクション \(Python: get\\_tags\)](#)
- [一般的なデータ型](#)
  - [Tag 構造](#)
  - [DecimalNumber 構造](#)
  - [ErrorDetail 構造](#)
  - [PropertyPredicate 構造](#)
  - [ResourceUri 構造](#)
  - [ColumnStatistics 構造](#)
  - [ColumnStatisticsError 構造](#)
  - [ColumnError 構造](#)
  - [ColumnStatisticsData 構造](#)
  - [BooleanColumnStatisticsData 構造](#)
  - [DateColumnStatisticsData 構造](#)
  - [DecimalColumnStatisticsData 構造](#)
  - [DoubleColumnStatisticsData 構造](#)
  - [LongColumnStatisticsData 構造](#)
  - [StringColumnStatisticsData 構造](#)
  - [BinaryColumnStatisticsData 構造](#)
  - [文字列パターン](#)
- [例外](#)
  - [AccessDeniedException 構造](#)
  - [AlreadyExistsException 構造](#)
  - [ConcurrentModificationException 構造](#)
  - [ConcurrentRunsExceededException 構造](#)
  - [CrawlerNotRunningException 構造](#)
  - [CrawlerRunningException 構造](#)
  - [CrawlerStoppingException 構造](#)
  - [EntityNotFoundException 構造](#)

- [FederationSourceException](#) 構造
- [FederationSourceRetryableException](#) 構造
- [GlueEncryptionException](#) 構造
- [IdempotentParameterMismatchException](#) 構造
- [IllegalWorkflowStateException](#) 構造
- [InternalServiceException](#) 構造
- [InvalidExecutionEngineException](#) 構造
- [InvalidInputException](#) 構造
- [InvalidStateException](#) 構造
- [InvalidTaskStatusTransitionException](#) 構造
- [JobDefinitionErrorException](#) 構造
- [JobRunInTerminalStateException](#) 構造
- [JobRunInvalidStateTransitionException](#) 構造
- [JobRunNotInTerminalStateException](#) 構造
- [LateRunnerException](#) 構造
- [NoScheduleException](#) 構造
- [OperationTimeoutException](#) 構造
- [ResourceNotReadyException](#) 構造
- [ResourceNumberLimitExceededException](#) 構造
- [SchedulerNotRunningException](#) 構造
- [SchedulerRunningException](#) 構造
- [SchedulerTransitioningException](#) 構造
- [UnrecognizedRunnerException](#) 構造
- [ValidationException](#) 構造
- [VersionMismatchException](#) 構造

## AWS Glue のセキュリティ API

Security API では、セキュリティのデータ型と、AWS Glue のセキュリティに関する API について説明しています。

## データ型

- [DataCatalogEncryptionSettings](#) 構造
- [EncryptionAtRest](#) 構造
- [ConnectionPasswordEncryption](#) 構造
- [EncryptionConfiguration](#) 構造
- [S3Encryption](#) 構造
- [CloudWatchEncryption](#) 構造
- [JobBookmarksEncryption](#) 構造
- [SecurityConfiguration](#) 構造
- [GluePolicy](#) の構造

### DataCatalogEncryptionSettings 構造

データカタログのセキュリティを維持するための構成情報が含まれています。

フィールド

- `EncryptionAtRest` – [EncryptionAtRest](#) オブジェクト。

データカタログの保管時の暗号化の構成を指定します。

- `ConnectionPasswordEncryption` – [ConnectionPasswordEncryption](#) オブジェクト。

接続パスワードが有効になっている場合、データカタログでは、お客様が用意したキーを使用して、`CreateConnection` または `UpdateConnection` の一部としてパスワードを暗号化するか、接続プロパティの `ENCRYPTED_PASSWORD` フィールドに保存します。カタログの暗号化、またはパスワードの暗号化のみ有効にすることができます。

### EncryptionAtRest 構造

データカタログの保管時の暗号化の構成を指定します。

フィールド

- `CatalogEncryptionMode` – 必須: UTF-8 文字列 (有効な値: `DISABLED` | `SSE-KMS="SSEKMS"` | `SSE-KMS-WITH-SERVICE-ROLE="SSEKMSWITHSERVICEROLE"`)。

データカタログを暗号化するための保管時の暗号化モード。

- `SseAwsKmsKeyId` – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

保管時の暗号化に使用する、AWS KMS キーの ID。

- `CatalogEncryptionServiceRole` – UTF-8 文字列、「[Custom string pattern #24](#)」に一致。

呼び出し元に代わり、AWS Glue が Data Catalog オブジェクトを暗号化および復号するために引き受けるロール。

## ConnectionPasswordEncryption 構造

`CreateConnection` または `UpdateConnection` の一部としてパスワードを暗号化し、接続プロパティの `ENCRYPTED_PASSWORD` フィールドに保存するために、データカタログで使用されるデータ構造。カタログの暗号化、またはパスワードの暗号化のみ有効にすることができます。

パスワードを含む `CreationConnection` リクエストが到着すると、Data Catalog ではまず AWS KMS キーを使用してパスワードを暗号化します。カタログの暗号化が有効になっている場合は、続いて全体の接続オブジェクトも暗号化します。

この暗号化では、セキュリティ要件に従ってパスワードキーへのアクセスを有効にしたり制限したりするように AWS KMS キーのアクセス許可を設定する必要があります。たとえば、管理者にのみ、パスワードキーの復号のアクセス許可を付与する場合などです。

### フィールド

- `ReturnConnectionPasswordEncrypted` – 必須: Boolean。

`ReturnConnectionPasswordEncrypted` フラグが「true」に設定されている場合は、`GetConnection` および `GetConnections` の応答でパスワードは暗号化された状態のままになります。この暗号化は、カタログの暗号化とは関係なく、有効になります。

- `AwsKmsKeyId` – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

接続パスワードの暗号化に使用する AWS KMS キー。

接続パスワードの保護が有効になっている場合、`CreateConnection` および `UpdateConnection` の呼び出し元には少なくとも、指定された AWS KMS キーの `kms:Encrypt` アクセス許可が必要です。このアクセス許可は、Data Catalog に保存する前にパスワードを暗号化するために必要です。

復号のアクセス許可を設定することで、セキュリティ要件に従ってパスワードキーへのアクセスを有効または制限することができます。

## EncryptionConfiguration 構造

暗号化構成を指定します。

フィールド

- S3Encryption – [S3Encryption](#) オブジェクトの配列。  
Amazon Simple Storage Service (Amazon S3) データの暗号化設定。
- CloudWatchEncryption – [CloudWatchEncryption](#) オブジェクト。  
Amazon CloudWatch の暗号化構成。
- JobBookmarksEncryption – [JobBookmarksEncryption](#) オブジェクト。  
ジョブブックマークの暗号化構成。

## S3Encryption 構造

Amazon Simple Storage Service (Amazon S3) データを暗号化する方法を指定します。

フィールド

- S3EncryptionMode – UTF-8 文字列 (有効な値: DISABLED | SSE-KMS="SSEKMS" | SSE-S3="SSES3")。  
Simple Storage Service (Amazon S3) データに使用する暗号化モード。
- KmsKeyArn – UTF-8 文字列、「[Custom string pattern #25](#)」に一致。  
データの暗号化に使用する KMS キーの Amazon リソースネーム (ARN)。

## CloudWatchEncryption 構造

Amazon CloudWatch データを暗号化する方法を指定します。

## フィールド

- `CloudWatchEncryptionMode` – UTF-8 文字列 (有効な値: `DISABLED` | `SSE-KMS="SSEKMS"`)。  
CloudWatch データに使用する暗号化モード。
- `KmsKeyArn` – UTF-8 文字列、[「Custom string pattern #25」](#) に一致。  
データの暗号化に使用する KMS キーの Amazon リソースネーム (ARN)。

## JobBookmarksEncryption 構造

ジョブブックマークデータを暗号化する方法を指定します。

### フィールド

- `JobBookmarksEncryptionMode` – UTF-8 文字列 (有効な値: `DISABLED` | `CSE-KMS="CSEKMS"`)。  
ジョブブックマークデータに使用する暗号化モード。
- `KmsKeyArn` – UTF-8 文字列、[「Custom string pattern #25」](#) に一致。  
データの暗号化に使用する KMS キーの Amazon リソースネーム (ARN)。

## SecurityConfiguration 構造

セキュリティ設定を指定します。

### フィールド

- `Name` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
セキュリティ構成の名前。
- `CreatedTimeStamp` – タイムスタンプ。  
このセキュリティ設定が作成された時刻。
- `EncryptionConfiguration` – [EncryptionConfiguration](#) オブジェクト。  
このセキュリティ構成に関連する暗号化構成。

## GluePolicy の構造

リソースポリシーを返すための構造。

### フィールド

- PolicyInJson – UTF-8 文字列、少なくとも 2 バイト長。

リクエストされたポリシードキュメントを JSON 形式で含んでいます。

- PolicyHash – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このポリシーに関連付けられたハッシュ値を含んでいます。

- CreateTime – タイムスタンプ。

ポリシーの作成日時。

- UpdateTime – タイムスタンプ。

ポリシーの前回更新日時。

## 操作

- [GetDataCatalogEncryptionSettings アクション](#) (Python: [get\\_data\\_catalog\\_encryption\\_settings](#))
- [PutDataCatalogEncryptionSettings アクション](#) (Python: [put\\_data\\_catalog\\_encryption\\_settings](#))
- [PutResourcePolicy アクション](#) (Python: [put\\_resource\\_policy](#))
- [GetResourcePolicy アクション](#) (Python: [get\\_resource\\_policy](#))
- [DeleteResourcePolicy アクション](#) (Python: [delete\\_resource\\_policy](#))
- [CreateSecurityConfiguration アクション](#) (Python: [create\\_security\\_configuration](#))
- [DeleteSecurityConfiguration アクション](#) (Python: [delete\\_security\\_configuration](#))
- [GetSecurityConfiguration アクション](#) (Python: [get\\_security\\_configuration](#))
- [GetSecurityConfigurations アクション](#) (Python: [get\\_security\\_configurations](#))
- [GetResourcePolicies アクション](#) (Python: [get\\_resource\\_policies](#))

## GetDataCatalogEncryptionSettings アクション (Python: `get_data_catalog_encryption_settings`)

指定されたカタログのセキュリティ設定を取得します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

セキュリティ設定を取得するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

### レスポンス

- `DataCatalogEncryptionSettings` – [DataCatalogEncryptionSettings](#) オブジェクト。

リクエストされたセキュリティ設定。

### エラー

- `InternalServiceException`
- `InvalidInputException`
- `OperationTimeoutException`

## PutDataCatalogEncryptionSettings アクション (Python: `put_data_catalog_encryption_settings`)

指定したカタログのセキュリティ構成を設定します。構成が完了すると、以後、指定した暗号化がすべてのカタログ書き込みに適用されます。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

セキュリティ設定を行うデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- `DataCatalogEncryptionSettings` – 必須: [DataCatalogEncryptionSettings](#) オブジェクト。

設定するセキュリティ設定。

## レスポンス

- 応答パラメータはありません。

## エラー

- `InternalServiceException`
- `InvalidInputException`
- `OperationTimeoutException`

## PutResourcePolicy アクション (Python: `put_resource_policy`)

アクセスコントロールのためにデータカタログのリソースポリシーを設定します。

### リクエスト

- `PolicyInJson` – 必須: UTF-8 文字列、少なくとも 2 バイト長。

設定する、JSON 形式のポリシードキュメントを含んでいます。

- `ResourceArn` – UTF-8 文字列。1~10240 バイト長。 [Custom string pattern #22](#) に一致。

使用しません。内部使用のみ。

- `PolicyHashCondition` – UTF-8 文字列、1~255 バイト長、 [Single-line string pattern](#) に一致。

前のポリシーが `PutResourcePolicy` を使用して設定された時に返るハッシュ値。その目的は、ポリシーの同時変更を防ぐことです。以前に設定されたポリシーがない場合は、このパラメータは使用しないでください。

- `PolicyExistsCondition` – UTF-8 文字列 (有効な値: `MUST_EXIST` | `NOT_EXIST` | `NONE`)。

ポリシーの値を更新するために値 `MUST_EXIST` が使用されます。新しいポリシーを作成するために値 `NOT_EXIST` が使用されます。値 `NONE` または `null` 値を使用すると、呼び出しはポリシーが存在するかどうか依存しなくなります。

- `EnableHybrid` – UTF-8 文字列 (有効な値: `TRUE` | `FALSE`)。

'TRUE' の場合、Data Catalog のリソースへのクロスアカウントアクセスを許可するために次の両方の方法を使用していることを示しています。

- リソースポリシーを PutResourcePolicy で直接更新
- AWS Management Console で Grant permissions コマンドを使用

マネジメントコンソールを使用してクロスアカウントアクセスを既に許可している場合、'TRUE' に設定する必要があります。そうしないと、呼び出しは失敗します。デフォルトは FALSE です。

## レスポンス

- PolicyHash – UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

設定されたばかりのポリシーのハッシュです。このポリシーを上書きまたは更新するこれ以降の呼び出しには、これを含める必要があります。

## エラー

- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException
- InvalidInputException
- ConditionCheckFailureException

## GetResourcePolicy アクション (Python: get\_resource\_policy)

指定されたリソースポリシーを取得します。

### リクエスト

- ResourceArn – UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。

リソースポリシーを取得する対象の AWS Glue リソースの ARN です。指定しない場合は、Data Catalog のリソースポリシーが返されます。GetResourcePolicies を使用して、既存のリソースポリシーをすべて表示します。詳細については、「[AWS Glue リソース ARN を特定する](#)」を参照してください。

## レスポンス

- PolicyInJson – UTF-8 文字列、少なくとも 2 バイト長。

リクエストされたポリシードキュメントを JSON 形式で含んでいます。

- PolicyHash – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このポリシーに関連付けられたハッシュ値を含んでいます。

- CreateTime – タイムスタンプ。

ポリシーの作成日時。

- UpdateTime – タイムスタンプ。

ポリシーの前回更新日時。

## エラー

- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException
- InvalidInputException

## DeleteResourcePolicy アクション (Python: delete\_resource\_policy)

指定されたポリシーを削除します。

### リクエスト

- PolicyHashCondition – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

このポリシーが設定されたときに返されたハッシュ値です。

- ResourceArn – UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。

削除されるリソースポリシーの AWS Glue リソースの ARN です。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException
- InvalidInputException
- ConditionCheckFailureException

## CreateSecurityConfiguration アクション (Python: create\_security\_configuration)

新しいセキュリティ構成を作成します。セキュリティ構成は、AWS Glue が使用できるセキュリティプロパティのセットです。セキュリティ構成を使用して、保管時にデータを暗号化できます。AWS Glue のセキュリティ設定を使用する方法については、「[クローラ、ジョブ、および開発エンドポイントによって書き込まれたデータの暗号化](#)」を参照してください。

## リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

新しいセキュリティ設定の名前。

- EncryptionConfiguration – 必須: [EncryptionConfiguration](#) オブジェクト。

新しいセキュリティ設定に関連する暗号化設定。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

新しいセキュリティ設定に割り当てられた名前。

- CreatedTimestamp – タイムスタンプ。

新しいセキュリティ設定が作成された時刻。

## エラー

- `AlreadyExistsException`
- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`
- `ResourceNumberLimitExceededException`

## DeleteSecurityConfiguration アクション (Python: `delete_security_configuration`)

指定したセキュリティ設定を削除します。

### リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
削除するセキュリティ設定の名前。

### レスポンス

- 応答パラメータはありません。

## エラー

- `EntityNotFoundException`
- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`

## GetSecurityConfiguration アクション (Python: `get_security_configuration`)

指定したセキュリティ設定を取得します。

### リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

取得するセキュリティ設定の名前。

## レスポンス

- SecurityConfiguration – [SecurityConfiguration](#) オブジェクト。

リクエストされたセキュリティ設定。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## GetSecurityConfigurations アクション (Python: `get_security_configurations`)

すべてのセキュリティ設定のリストを取得します。

### リクエスト

- MaxResults – 1 ~ 1000 の数値 (整数)。

返される結果の最大数。

- NextToken – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

### 応答

- SecurityConfigurations – [SecurityConfiguration](#) オブジェクトの配列。

セキュリティ設定のリスト。

- NextToken – UTF-8 文字列。

返されるセキュリティ設定がさらにある場合は、継続トークンを返します。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## GetResourcePolicies アクション (Python: get\_resource\_policies)

クロスアカウントのアクセス許可の付与中に AWS Resource Access Manager によって個々のリソースに設定されたリソースポリシーを取得します。Data Catalog のリソースポリシーも取得しません。

Data Catalog 設定でメタデータの暗号化を有効にしている、AWS KMS キーへのアクセス許可がない場合、オペレーションで Data Catalog のリソースポリシーを返すことはできません。

## リクエスト

- NextToken – UTF-8 文字列。  
継続トークン (これが継続リクエストの場合)。
- MaxResults – 1 ~ 1000 の数値 (整数)。  
返されるリストの最大サイズ。

## 応答

- GetResourcePoliciesResponseList – [GluePolicy](#) オブジェクトの配列。  
個々のリソースポリシーとアカウントレベルのリソースポリシーのリスト。
- NextToken – UTF-8 文字列。  
継続トークン (戻されたリストに使用可能な最後のリソースポリシーが含まれていない場合)。

## エラー

- InternalServiceException
- OperationTimeoutException

- `InvalidInputException`
- `GlueEncryptionException`

## Catalog API

Catalog API では、AWS Glue でのカタログの操作に関連するデータ型と API について説明します。

### トピック

- [データベース API](#)
- [Table API](#)
- [パーティション API](#)
- [接続 API](#)
- [ユーザー定義関数 API](#)
- [Athena カタログを AWS Glue にインポートする](#)

## データベース API

Database API では、データベースのデータ型について説明し、これにはデータベースを作成、削除、検索、更新、および一覧表示するための API が含まれます。

### データ型

- [データベース構造](#)
- [DatabaseInput 構造](#)
- [PrincipalPermissions 構造](#)
- [DataLakePrincipal 構造](#)
- [DatabaseIdentifier 構造](#)
- [FederatedDatabase 構造](#)

### データベース構造

Database オブジェクトは、Hive メタストアまたは RDBMS に存在する可能性のあるテーブルの論理グループを表します。

## フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データベースの名前。Hive との互換性を保つため、これは保存時には小文字で折りた畳まれます。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

データベースの説明。

- LocationUri – Uniform resource identifier (uri)、1~1024 バイト長、[URI address multi-line string pattern](#) に一致。

データベースの場所 (たとえば HDFS パスなど)。

- Parameters – キーバリューペアのマップ配列。

各キーはキー文字列。1~255 バイト長。[Single-line string pattern](#) に一致。

各値は UTF-8 文字列で、512,000 バイト長以下です。

これらのキー/値ペアはデータベースのパラメータとプロパティを定義します。

- CreateTime – タイムスタンプ。

メタデータデータベースがカタログに作成された時刻。

- CreateTableDefaultPermissions – [PrincipalPermissions](#) オブジェクトの配列。

プリンシパル用のデフォルトのアクセス権限セットをテーブルに作成します。AWS Lake Formation で使用されます。AWS Glue オペレーションの通常の過程では使用されません。

- TargetDatabase – [DatabaseIdentifier](#) オブジェクト。

リソースリンク用のターゲットデータベースを記述する DatabaseIdentifier 構造。

- CatalogId – カatalog ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データベースが存在するデータカタログの ID。

- FederatedDatabase – [FederatedDatabase](#) オブジェクト。

AWS Glue Data Catalog の外部のエンティティを参照する FederatedDatabase 構造。

## Databaselnput 構造

データベースの作成または更新に使用される構造体。

### フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
データベースの名前。Hive との互換性を保つため、これは保存時には小文字で折りた畳まれます。
- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
データベースの説明。
- LocationUri – Uniform resource identifier (uri)、1~1024 バイト長、[URI address multi-line string pattern](#) に一致。  
データベースの場所 (たとえば HDFS パスなど)。
- Parameters – キーバリュールペアのマッピング配列。  
各キーはキー文字列。1~255 バイト長。[Single-line string pattern](#) に一致。  
各値は UTF-8 文字列で、512,000 バイト長以下です。  
これらのキー/値ペアはデータベースのパラメータとプロパティを定義します。  
これらのキー/値ペアはデータベースのパラメータとプロパティを定義します。
- CreateTableDefaultPermissions – [PrincipalPermissions](#) オブジェクトの配列。  
プリンシパル用のデフォルトのアクセス権限セットをテーブルに作成します。AWS Lake Formation で使用されます。AWS Glue オペレーションの通常の過程では使用されません。
- TargetDatabase – [DatabaseIdentifier](#) オブジェクト。  
リソースリンク用のターゲットデータベースを記述する DatabaseIdentifier 構造。
- FederatedDatabase – [FederatedDatabase](#) オブジェクト。  
AWS Glue Data Catalog の外部のエンティティを参照する FederatedDatabase 構造。

## PrincipalPermissions 構造

プリンシパルに付与される許可です。

## フィールド

- **Principal** – [DataLakePrincipal](#) オブジェクト。

許可を付与されたプリンシパル。

- **Permissions** – UTF-8 文字列の配列。

プリンシパルに付与される許可。

## DataLakePrincipal 構造

AWS Lake Formation プリンシパルです。

### フィールド

- **DataLakePrincipalIdentifier** – UTF-8 文字列、1~255 バイト長。

AWS Lake Formation プリンシパルの識別子。

## DatabasIdentifier 構造

リソースリンク用のターゲットデータベースを記述する 構造。

### フィールド

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データベースが存在するデータカタログの ID。

- **DatabaseName** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

カタログデータベースの名前。

- **Region** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ターゲットデータベースのリージョン。

## FederatedDatabase 構造

AWS Glue Data Catalog の外部のエンティティを指すデータベース。

## フィールド

- Identifier - UTF-8 文字列。1 ~ 512 バイト長。 [Single-line string pattern](#) に一致。  
フェデレーションデータベースの一意の識別子。
- ConnectionName - UTF-8 文字列、1 ~ 255 バイト長、 [Single-line string pattern](#) に一致。  
外部メタストアへの接続の名前。

## 操作

- [CreateDatabase アクション \(Python: create\\_database\)](#)
- [UpdateDatabase アクション \(Python: update\\_database\)](#)
- [DeleteDatabase アクション \(Python: delete\\_database\)](#)
- [GetDatabase アクション \(Python: get\\_database\)](#)
- [GetDatabases アクション \(Python: get\\_databases\)](#)

## CreateDatabase アクション (Python: create\_database)

データカタログに新しいデータベースを作成します。

### リクエスト

- CatalogId - カタログ ID 文字列、1 ~ 255 バイト長、 [Single-line string pattern](#) に一致。  
データベースを作成するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。
- DatabaseInput - 必須: [DatabaseInput](#) オブジェクト。  
データベースのメタデータ。
- Tags - キーと値のペアのマップ配列。50 ペア以下。  
各キーは UTF-8 文字列で、1 ~ 128 バイト長です。  
各値は UTF-8 文字列で、256 バイト長以下です。  
データベースに割り当てるタグ。

## レスポンス

- 応答パラメータはありません。

## エラー

- `InvalidInputException`
- `AlreadyExistsException`
- `ResourceNumberLimitExceededException`
- `InternalServiceException`
- `OperationTimeoutException`
- `GlueEncryptionException`
- `ConcurrentModificationException`
- `FederatedResourceAlreadyExistsException`

## UpdateDatabase アクション (Python: `update_database`)

データカタログの既存のデータベース定義を更新します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

メタデータのデータベースが存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

カタログで更新するデータベースの名前。Hive 互換性のために、これは小文字で表記されます。

- `DatabaseInput` – 必須: [DatabaseInput](#) オブジェクト。

カタログ内のメタデータデータベースの新しい定義を指定する `DatabaseInput` オブジェクト。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException
- ConcurrentModificationException

## DeleteDatabase アクション (Python: delete\_database)

指定されたデータベースをデータカタログから削除します。

### Note

このオペレーションが完了すると、削除されたデータベース内のテーブル (それらのテーブルに属する可能性のあるすべてのテーブルのバージョンとパーティションも含む) とユーザー定義関数にはアクセスできなくなります。AWS Glue では、このサービスの裁量により、これらの「孤立した」リソースを適時に非同期的に削除します。

関連するすべてのリソースを確実にすぐに削除するには、DeleteDatabase を呼び出す前に、DeleteTableVersion/BatchDeleteTableVersion、DeletePartition/BatchDeletePartition を使用して、データベースに属するすべてのリソースを削除します。

## リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データベースが存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除するデータベースの名前。Hive との互換性を保つため、名前はすべて小文字にする必要があります。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- ConcurrentModificationException

## GetDatabase アクション (Python: get\_database)

指定されたデータベースの定義を取得します。

## リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データベースが存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

取得するデータベースの名前。Hive 互換性のために、これはすべて小文字にする必要があります。

## レスポンス

- Database – [データベース](#) オブジェクト。

データカタログ内の指定されたデータベースの定義。

## エラー

- InvalidInputException
- EntityNotFoundException

- `InternalServiceException`
- `OperationTimeoutException`
- `GlueEncryptionException`
- `FederationSourceException`

## GetDatabases アクション (Python: `get_databases`)

指定されたデータカタログで定義されているすべてのデータベースを取得します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

`Databases` を取得するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- `NextToken` – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

- `MaxResults` – 数値 (integer)。1~100。

1 回の応答で返すデータベースの最大数。

- `ResourceShareType` – UTF-8 文字列 (有効な値: FOREIGN | ALL | FEDERATED)。

アカウントと共有しているデータベースを一覧表示するように指定できます。指定できる値は、FEDERATED、FOREIGN または ALL です。

- FEDERATED に設定した場合は、アカウントと共有されているフェデレーションデータベース (外部エンティティを参照) が一覧表示されます。
- FOREIGN に設定した場合は、アカウントと共有されているデータベースが一覧表示されます。
- ALL に設定した場合は、アカウントと共有されているデータベースと、ローカルアカウントのデータベースが一覧表示されます。

### レスポンス

- `DatabaseList` – 必須: [データベース](#) オブジェクトの配列。

指定されたカタログの Database オブジェクトのリスト。

- `NextToken` – UTF-8 文字列。

返されたトークンのリストをページ分割するための継続トークン。リストの現在のセグメントが最後のセグメントではない場合に返されます。

## エラー

- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`
- `GlueEncryptionException`

## Table API

テーブル API は、データ型とテーブルに関連する操作について説明しています。

## データ型

- [Table 構造](#)
- [TableInput 構造](#)
- [FederatedTable 構造](#)
- [列の構造](#)
- [StorageDescriptor 構造](#)
- [SchemaReference 構造](#)
- [SerDeInfo 構造](#)
- [Order 構造](#)
- [SkewedInfo 構造](#)
- [TableVersion 構造](#)
- [TableError 構造](#)
- [TableVersionError 構造](#)
- [SortCriterion 構造](#)
- [TableIdentifier 構造](#)
- [KeySchemaElement 構造](#)

- [PartitionIndex 構造](#)
- [PartitionIndexDescriptor 構造](#)
- [BackfillError 構造](#)
- [IcebergInput 構造](#)
- [OpenTableFormatInput 構造](#)
- [ViewDefinition 構造](#)
- [ViewDefinitionInput 構造](#)
- [ViewRepresentation 構造](#)
- [ViewRepresentationInput 構造](#)

## Table 構造

列と行で構成されている関連データのコレクションを表します。

### フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの名前。Hive 互換性のために、これはすべて小文字であることが必要です。

- DatabaseName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルメタデータが存在するデータベースの名前。Hive との互換性を保つため、名前はすべて小文字にする必要があります。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

テーブルの説明。

- Owner – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの所有者。

- CreateTime – タイムスタンプ。

データカタログでテーブル定義が作成された時刻。

- UpdateTime – タイムスタンプ。

テーブルが最後に更新された時刻。

- LastAccessTime – タイムスタンプ。

テーブルが最後にアクセスされた時刻。これは通常 HDFS から取得されるため、信頼性が低い場合があります。

- LastAnalyzedTime – タイムスタンプ。

このテーブルの列統計が最後に計算された時刻。

- Retention – 数値 (整数)、None 以下。

このテーブルの保持時間。

- StorageDescriptor – [StorageDescriptor](#) オブジェクト。

このテーブルの物理ストレージに関する情報を含むストレージ記述子。

- PartitionKeys – [Column](#) オブジェクトの配列。

テーブルがパーティション分割される列のリスト。パーティションキーとしてプリミティブ型のみがサポートされています。

Amazon Athena によって使用されるテーブルを作成しており、partitionKeys を指定していない場合は必ず、partitionKeys の値を空のリストに設定する必要があります。例:

```
"PartitionKeys": []
```

- ViewOriginalText - UTF-8 文字列。409,600 バイト長以下。

Apache Hive との互換性のために含まれています。通常の AWS Glue オペレーションでは使用されません。テーブルが の場合 VIRTUAL\_VIEW、base64 でエンコードされた特定の Athena 設定。

- ViewExpandedText - UTF-8 文字列。409,600 バイト長以下。

Apache Hive との互換性のために含まれています。通常の AWS Glue オペレーションでは使用されません。

- TableType - UTF-8 文字列。255 バイト長以下。

このテーブルのタイプ。AWS Glue は、EXTERNAL\_TABLE タイプのテーブルを作成します。などの他のサービスでは Athena、追加のテーブルタイプを持つテーブルが作成される場合があります。

AWS Glue 関連テーブルタイプ :

EXTERNAL\_TABLE

Hive 互換属性 - Hive 以外の管理テーブルを示します。

## GOVERNED

によって使用されます AWS Lake Formation。AWS Glue Data Catalog は を理解し  
ます GOVERNED。

- Parameters – キーバリューペアのマップ配列。

各キーはキー文字列。1~255 バイト長。 [Single-line string pattern](#) に一致。

各値は UTF-8 文字列で、512,000 バイト長以下です。

これらのキー/値ペアは、テーブルに関連付けられたプロパティを定義します。

- CreatedBy – UTF-8 文字列、1~255 バイト長、 [Single-line string pattern](#) に一致。

テーブルを作成した個人または団体。

- IsRegisteredWithLakeFormation – ブール。

テーブルが に登録されているかどうかを示します AWS Lake Formation。

- TargetTable – [TableIdentifier](#) オブジェクト。

リソースリンク用のターゲットテーブルを記述する TableIdentifier 構造。

- CatalogId – カタログ ID 文字列、1~255 バイト長、 [Single-line string pattern](#) に一致。

テーブルデータベースがあるデータカタログの ID。

- VersionId – UTF-8 文字列、1~255 バイト長、 [Single-line string pattern](#) に一致。

テーブルバージョンの ID。

- FederatedTable – [FederatedTable](#) オブジェクト。

AWS Glue Data Catalogの外部のエンティティを参照する FederatedTable 構造。

- ViewDefinition – [ViewDefinition](#) オブジェクト。

ビューのダイアレクト、クエリなど、ビューを定義するすべての情報を含む構造。

- IsMultiDialectView – ブール。

ビューがさまざまなクエリエンジンの SQL ダイアレクトをサポートし、それらのエンジンで読み  
取れるようにするかどうかを指定します。

## TableInput 構造

テーブルの定義に使用される構造。

### フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの名前。Hive との互換性を保つため、これは保存時には小文字で折りた畳まれます。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

テーブルの説明。

- Owner – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの所有者。Apache Hive との互換性のために含まれています。通常の AWS Glue オペレーションでは使用されません。

- LastAccessTime – タイムスタンプ。

テーブルが最後にアクセスされた時刻。

- LastAnalyzedTime – タイムスタンプ。

このテーブルの列統計が最後に計算された時刻。

- Retention – 数値 (整数)、None 以下。

このテーブルの保持時間。

- StorageDescriptor – [StorageDescriptor](#) オブジェクト。

このテーブルの物理ストレージに関する情報を含むストレージ記述子。

- PartitionKeys – [Column](#) オブジェクトの配列。

テーブルがパーティション分割される列のリスト。パーティションキーとしてプリミティブ型のみがサポートされています。

Amazon Athena によって使用されるテーブルを作成しており、partitionKeys を指定していない場合は必ず、partitionKeys の値を空のリストに設定する必要があります。例:

```
"PartitionKeys": []
```

- ViewOriginalText – UTF-8 文字列。409,600 バイト長以下。

Apache Hive との互換性のために含まれています。通常の AWS Glue オペレーションでは使用されません。テーブルが の場合VIRTUAL\_VIEW、base64 でエンコードされた特定の Athena 設定。

- ViewExpandedText - UTF-8 文字列。409,600 バイト長以下。

Apache Hive との互換性のために含まれています。通常の AWS Glue オペレーションでは使用されません。

- TableType - UTF-8 文字列。255 バイト長以下。

このテーブルのタイプ。AWS Glue は、EXTERNAL\_TABLEタイプのテーブルを作成します。などの他のサービスでは Athena、追加のテーブルタイプを持つテーブルが作成される場合があります。

AWS Glue 関連テーブルタイプ :

EXTERNAL\_TABLE

Hive 互換属性 - Hive 以外の管理テーブルを示します。

GOVERNED

によって使用されます AWS Lake Formation。AWS Glue Data Catalog は を理解しません GOVERNED。

- Parameters – キーバリューペアのマップ配列。

各キーはキー文字列。1~255 バイト長。 [Single-line string pattern](#) に一致。

各値は UTF-8 文字列で、512,000 バイト長以下です。

これらのキー/値ペアは、テーブルに関連付けられたプロパティを定義します。

- TargetTable – [TableIdentifier](#) オブジェクト。

リソースリンク用のターゲットテーブルを記述する TableIdentifier 構造。

- ViewDefinition – [ViewDefinitionInput](#) オブジェクト。

ビューのダイレクト、クエリなど、ビューを定義するすべての情報を含む構造。

## FederatedTable 構造

AWS Glue Data Catalogの外部のエンティティを指すテーブル。

## フィールド

- Identifier - UTF-8 文字列。1 ~ 512 バイト長。 [Single-line string pattern](#) に一致。  
フェデレーションテーブルの一意の識別子。
- DatabaseIdentifier - UTF-8 文字列。1 ~ 512 バイト長。 [Single-line string pattern](#) に一致。  
フェデレーションデータベースの一意の識別子。
- ConnectionName - UTF-8 文字列、1 ~ 255 バイト長、 [Single-line string pattern](#) に一致。  
外部メタストアへの接続の名前。

## 列の構造

Table 内の列。

### フィールド

- Name - 必須: UTF-8 文字列、1 ~ 255 バイト長、 [Single-line string pattern](#) に一致。  
Column の名前。
- Type - UTF-8 文字列。131,072 バイト長以下。 [Single-line string pattern](#) に一致。  
Column のデータ型。
- Comment - コメント文字列。255 バイト長以下。 [Single-line string pattern](#) に一致。  
自由形式のテキストコメント。
- Parameters - キーバリューペアのマップ配列。  
各キーはキー文字列。1 ~ 255 バイト長。 [Single-line string pattern](#) に一致。  
各値は UTF-8 文字列で、512,000 バイト長以下です。  
これらのキーバリューペアは、列に関連付けられたプロパティを定義します。

## StorageDescriptor 構造

テーブルデータの物理ストレージについて説明します。

## フィールド

- Columns – [Column](#) オブジェクトの配列。

テーブル内の Columns のリストです。

- Location - 場所文字列。2,056 バイト長以下。 [URI address multi-line string pattern](#) に一致。

テーブルの物理的な場所。デフォルトでは、ウェアハウスの場所、その後にウェアハウス内のデータベースの場所、その後にテーブル名が続く形式になります。

- AdditionalLocations – UTF-8 文字列の配列。

Delta テーブルへのパスを指す場所のリスト。

- InputFormat - 書式設定文字列。128 バイト長以下。 [Single-line string pattern](#) に一致。

入力形式: SequenceFileInputFormat (バイナリ)、または TextInputFormat、もしくはカスタム形式。

- OutputFormat - 書式設定文字列。128 バイト長以下。 [Single-line string pattern](#) に一致。

出力形式: SequenceFileOutputFormat (バイナリ)、または IgnoreKeyTextOutputFormat、もしくはカスタム形式。

- Compressed – ブール。

テーブルのデータが圧縮されている場合は True、圧縮されていない場合は False。

- NumberOfBuckets – 数値 (整数)。

テーブルにディメンション列が含まれている場合、指定する必要があります。

- SerdeInfo – [SerDeInfo](#) オブジェクト。

シリアル化/逆シリアル化 (SerDe) 情報。

- BucketColumns – UTF-8 文字列の配列。

テーブルのリデューサーグループ化列、クラスター列、およびバケット列のリスト。

- SortColumns – [Order](#) オブジェクトの配列。

テーブル内の各バケットのソート順を指定するリスト。

- Parameters – キーバリューペアのマップ配列。

各キーはキー文字列。1 ~ 255 バイト長。 [Single-line string pattern](#) に一致。

各値は UTF-8 文字列で、512,000 バイト長以下です。

ユーザーが指定したキーバリュー形式のプロパティ。

- SkewedInfo – [SkewedInfo](#) オブジェクト。

列に高い頻度で表示される値 (歪んだ値) に関する情報。

- StoredAsSubDirectories – ブール。

テーブルデータがサブディレクトリに保管されている場合は True、保管されていない場合は False。

- SchemaReference – [SchemaReference](#) オブジェクト。

Schema Registry に保存されている AWS Glue スキーマを参照するオブジェクト。

テーブルを作成するときにスキーマに空の列リストを渡し、代わりにスキーマリファレンスを使用します。

## SchemaReference 構造

Schema Registry に保存されている AWS Glue スキーマを参照するオブジェクト。

フィールド

- SchemaId – [Schemald](#) オブジェクト。

スキーマアイデンティティフィールドを含む構造。こちらか SchemaVersionId のいずれかを指定する必要があります。

- SchemaVersionId – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。

スキーマのバージョンに割り当てられた一意の ID。こちらか SchemaId のいずれかを指定する必要があります。

- SchemaVersionNumber – 数値 (long)。1 ~ 100000。

スキーマのバージョン番号。

## SerDeInfo 構造

エクストラクタとローダーとして機能するシリアル化/逆シリアル化プログラム (SerDe) に関する情報。

### フィールド

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

の名前 SerDe。

- SerializationLibrary – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

通常、を実装するクラス SerDe。例は

`org.apache.hadoop.hive.serde2.columnar.ColumnarSerDe` です。

- Parameters – キーバリューペアのマップ配列。

各キーはキー文字列。1~255 バイト長。[Single-line string pattern](#) に一致。

各値は UTF-8 文字列で、512,000 バイト長以下です。

これらのキーと値のペアは、の初期化パラメータを定義します SerDe。

## Order 構造

ソートされた列のソート順を指定します。

### フィールド

- Column – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

列の名前。

- SortOrder – 必須: 数値 (integer)。1 以下。

列が昇順 (== 1) または降順 (==0) でソートされていることを指定します。

## SkewedInfo 構造

テーブルで歪んだ値を指定します。「歪んだ値」とは、非常に高い頻度で発生する値です。

## フィールド

- `SkewedColumnNames` – UTF-8 文字列の配列。  
歪んだ値を含む列名のリスト。
- `SkewedColumnValues` – UTF-8 文字列の配列。  
頻繁に出現するため、歪んだとみなされる値のリスト。
- `SkewedColumnValueLocationMaps` – キーバリューペアのマップ配列。  
各キーは UTF-8 文字列。  
各値は UTF-8 文字列。  
歪んだ値が含まれている列へのマッピング。

## TableVersion 構造

テーブルのバージョンを指定します。

### フィールド

- `Table` – [Table](#) オブジェクト。  
該当するテーブル。
- `VersionId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
このテーブルのバージョンを特定する ID 値。VersionId は整数の文字列表現です。各バージョンは 1 ずつ増加します。

## TableError 構造

テーブル操作に対するエラーレコード。

### フィールド

- `TableName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルの名前。Hive 互換性のために、これはすべて小文字であることが必要です。
- `ErrorDetail` – [ErrorDetail](#) オブジェクト。

エラーに関する詳細。

## TableVersionError 構造

テーブルバージョン操作に対するエラーレコード。

フィールド

- **TableName** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当するテーブルの名前。

- **VersionId** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当するバージョンの ID 値。VersionID は整数の文字列表現です。各バージョンは 1 ずつ増加します。

- **ErrorDetail** – [ErrorDetail](#) オブジェクト。

エラーに関する詳細。

## SortCriterion 構造

ソート基準となるフィールドとソート順を指定します。

フィールド

- **FieldName** – Value 文字列、1024 バイト長以下。

ソートするフィールドの名前。

- **Sort** – UTF-8 文字列 (有効な値: ASC="ASCENDING" | DESC="DESCENDING")。

昇順または降順のソート。

## TableIdentifier 構造

リソースリンク用のターゲットテーブルを記述する構造。

フィールド

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルデータベースがあるデータカタログの ID。

- DatabaseName – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

ターゲットテーブルを含むカタログデータベースの名前。

- Name – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

ターゲットテーブルの名前。

- Region – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

ターゲットテーブルのリージョン。

## KeySchemaElement 構造

名前とタイプで構成されるパーティションキーペア。

フィールド

- Name – 必須: UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

パーティションキーの名前。

- Type – 必須: UTF-8 文字列。131072 バイト長以下。[Single-line string pattern](#) に一致。

パーティションキーのタイプ。

## PartitionIndex 構造

パーティションインデックスの構造。

フィールド

- Keys – 必須: UTF-8 文字列の配列。1 個の以上の文字列。

パーティションインデックスのキー。

- IndexName – 必須: UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

パーティションインデックスの名前。

## PartitionIndexDescriptor 構造

テーブル内のパーティションインデックスの記述子。

フィールド

- **IndexName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションインデックスの名前。

- **Keys** – 必須: [KeySchemaElement](#) オブジェクトの配列、少なくとも 1 つの構造体。

[KeySchemaElement](#) 構造体としてのパーティションインデックスの 1 つ以上のキーのリスト。

- **IndexStatus** – 必須: UTF-8 文字列 (有効な値: CREATING | ACTIVE | DELETING | FAILED)。

インデックスパーティションのステータス。

次のようなステータスがあります。

- **CREATING**: インデックスは作成中です。インデックスが CREATING 状態の場合、インデックスまたはそのテーブルは削除できません。
  - **ACTIVE**: インデックスの作成に成功しました。
  - **FAILED**: インデックスの作成に失敗しました。
  - **DELETING**: インデックスは、インデックスのリストから削除されます。
- **BackfillErrors** – [BackfillError](#) オブジェクトの配列。

既存のテーブルのパーティションインデックスを登録するときに発生する可能性があるエラーのリスト。

## BackfillError 構造

既存のテーブルのパーティションインデックスを登録するときに発生する可能性があるエラーのリスト。

これらのエラーでは、インデックスの登録が失敗した理由の詳細が表示され、応答として限られた数のパーティションが示されます。これにより、障害が発生したパーティションを修正し、インデックスの登録を再試行できます。発生する可能性のある最も一般的なエラーのセットは、次のように分類されます。

- **EncryptedPartitionError**: パーティションは暗号化されます。

- `InvalidPartitionTypeError`: パーティション値が、そのパーティション列のデータ型と一致しません。
- `MissingPartitionValueError`: パーティションは暗号化されます。
- `UnsupportedPartitionCharacterError`: パーティション値内の文字はサポートされていません。例: `U+0000`、`U+0001`、`U+0002`。
- `InternalError`: 他のエラーコードに属さないエラー。

## フィールド

- `Code` – UTF-8 文字列 (有効な値: `ENCRYPTED_PARTITION_ERROR` | `INTERNAL_ERROR` | `INVALID_PARTITION_TYPE_DATA_ERROR` | `MISSING_PARTITION_VALUE_ERROR` | `UNSUPPORTED_PARTITION_CHARACTER_ERROR`)。

既存のテーブルのパーティションインデックスを登録するときに発生したエラーのエラーコード。

- `Partitions` – [PartitionValueList](#) オブジェクトの配列。

応答内の限られた数のパーティションのリスト。

## IcebergInput 構造

カタログ内に作成する Apache Iceberg メタデータテーブルを定義する構造。

### フィールド

- `MetadataOperation` – 必須: UTF-8 文字列 (有効な値: `CREATE`)。  
必須のメタデータオペレーション。これは `CREATE` にのみ設定できます。
- `Version` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

Iceberg テーブルのテーブルバージョン。デフォルトは 2 です。

## OpenTableFormatInput 構造

オープン形式のテーブルを表す構造。

### フィールド

- `IcebergInput` – [IcebergInput](#) オブジェクト。

Apache Iceberg メタデータテーブルを定義する IcebergInput 構造を指定します。

## ViewDefinition 構造

表現の詳細を含む構造。

フィールド

- IsProtected - ブール。

このフラグを true に設定すると、クエリ計画中にユーザー提供のオペレーションをビューの論理プランにプッシュしないようにエンジンに指示できます。ただし、このフラグを設定しても、エンジンが従うとは限りません。エンジンのドキュメントを参照して、提供される保証を理解してください。

- Definer - UTF-8 文字列。1 ~ 512 バイト長。[Single-line string pattern](#) に一致。

SQL のビューの定義者。

- SubObjects - UTF-8 文字列の配列。文字列 10 個以下。

テーブル Amazon リソースネーム (ARN) のリスト。

- Representations - [ViewRepresentation](#) オブジェクトの配列。1 ~ 1000 個の構造体。

表現のリスト。

## ViewDefinitionInput 構造

AWS Glue ビューを作成または更新するための詳細を含む構造。

フィールド

- IsProtected - ブール。

このフラグを true に設定すると、クエリ計画中にユーザー提供のオペレーションをビューの論理プランにプッシュしないようにエンジンに指示できます。ただし、このフラグを設定しても、エンジンが従うとは限りません。エンジンのドキュメントを参照して、提供される保証を理解してください。

- Definer - UTF-8 文字列。1 ~ 512 バイト長。[Single-line string pattern](#) に一致。

SQL のビューの定義者。

- Representations – [ViewRepresentationInput](#) オブジェクトの配列、1～10 個の構造。

ビューのダイレクトとビューを定義するクエリを含む構造のリスト。

- SubObjects - UTF-8 文字列の配列。文字列 10 個以下。

ビューを構成するベーステーブルの ARN のリスト。

## ViewRepresentation 構造

ビューのダイレクトとビューを定義するクエリを含む構造。

フィールド

- Dialect – UTF-8 文字列 (有効な値: REDSHIFT | ATHENA | SPARK)。

クエリエンジンのダイレクト。

- DialectVersion – UTF-8 文字列、1～255 バイト長。

クエリエンジンのダイレクトのバージョン。3.0.0 などです。

- ViewOriginalText - UTF-8 文字列。409,600 バイト長以下。

CREATE VIEW DDL 実行中に顧客が提供した SELECT クエリ。この SQL は、ビューのクエリ中には使用されません (ViewExpandedText が代わりに使用されます)。ViewOriginalText は、ユーザーがビューを作成した元の DDL コマンドを表示したい SHOW CREATE VIEW のような場合に使用されます。

- ViewExpandedText - UTF-8 文字列。409,600 バイト長以下。

ビューの拡張された SQL。この SQL は、ビューでクエリを処理するときにエンジンによって使用されます。エンジンは、ビューの作成中にオペレーションを実行して ViewOriginalText を ViewExpandedText に変換する場合があります。例:

- 完全修飾識別子: SELECT \* from table1 -> SELECT \* from db1.table1
- ValidationConnection – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

ビューの特定の表現の検証に使用される接続の名前。

- IsStale – ブール。

古いとマークされたダイレクトは無効になり、それぞれのクエリエンジンでクエリを実行する前に更新する必要があります。

## ViewRepresentationInput 構造

Lake Formation ビューを更新または作成するための表現の詳細を含む構造。

### フィールド

- `Dialect` – UTF-8 文字列 (有効な値: REDSHIFT | ATHENA | SPARK)。

特定の表現のエンジンタイプを指定するパラメータ。

- `DialectVersion` – UTF-8 文字列、1~255 バイト長。

特定の表現のエンジンのバージョンを指定するパラメータ。

- `ViewOriginalText` - UTF-8 文字列。409,600 バイト長以下。

ビューを記述する元の SQL クエリを表す文字列。

- `ValidationConnection` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ビューの特定の表現の検証に使用される接続の名前。

- `ViewExpandedText` - UTF-8 文字列。409,600 バイト長以下。

ビューを記述する SQL クエリおよび拡張されたリソースの ARN を表す文字列

### 操作

- [CreateTable アクション \(Python: create\\_table\)](#)
- [UpdateTable アクション \(Python: update\\_table\)](#)
- [DeleteTable アクション \(Python: delete\\_table\)](#)
- [BatchDeleteTable アクション \(Python: batch\\_delete\\_table\)](#)
- [GetTable アクション \(Python: get\\_table\)](#)
- [GetTables アクション \(Python: get\\_tables\)](#)
- [GetTableVersion アクション \(Python: get\\_table\\_version\)](#)
- [GetTableVersions アクション \(Python: get\\_table\\_versions\)](#)

- [DeleteTableVersion アクション \(Python: delete\\_table\\_version\)](#)
- [BatchDeleteTableVersion アクション \(Python: batch\\_delete\\_table\\_version\)](#)
- [SearchTables アクション \(Python: search\\_tables\)](#)
- [GetPartitionIndexes アクション \(Python: get\\_partition\\_indexes\)](#)
- [CreatePartitionIndex アクション \(Python: create\\_partition\\_index\)](#)
- [DeletePartitionIndex アクション \(Python: delete\\_partition\\_index\)](#)
- [GetColumnStatisticsForTable アクション \(Python: get\\_column\\_statistics\\_for\\_table\)](#)
- [UpdateColumnStatisticsForTable アクション \(Python: update\\_column\\_statistics\\_for\\_table\)](#)
- [DeleteColumnStatisticsForTable アクション \(Python: delete\\_column\\_statistics\\_for\\_table\)](#)

## CreateTable アクション (Python: create\_table)

データカタログで新しいテーブル定義を作成します。

### リクエスト

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

Table を作成するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- **DatabaseName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

新しいテーブルを作成するカタログデータベース。Hive 互換性のために、この名前はすべて小文字であることが必要です。

- **TableInput** – 必須: [TableInput](#) オブジェクト。

カタログで作成するメタデータテーブルを定義する TableInput オブジェクト。

- **PartitionIndexes** – [PartitionIndex](#) オブジェクトの配列。構造 3 個以下。

テーブル内に作成されるパーティションインデックスのリストである PartitionIndex 構造。

- **TransactionId** – UTF-8 文字列、1~255 バイト長、「[Custom string pattern #16](#)」に一致。

トランザクションの ID。

- **OpenTableFormatInput** – [OpenTableFormatInput](#) オブジェクト。

オープン形式のテーブルを作成するときに OpenTableFormatInput 構造を指定します。

## レスポンス

- 応答パラメータはありません。

## エラー

- AlreadyExistsException
- InvalidInputException
- EntityNotFoundException
- ResourceNumberLimitExceededException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException
- ConcurrentModificationException
- ResourceNotReadyException

## UpdateTable アクション (Python: update\_table)

データカタログのメタデータテーブルを更新します。

### リクエスト

- CatalogId – カタログ ID 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログデータベースの名前。Hive 互換性のために、この名前はすべて小文字であることが必要です。

- TableInput – 必須: [TableInput](#) オブジェクト。

カタログのメタデータテーブルを定義するための、更新された TableInput オブジェクト。

- SkipArchive – ブール。

UpdateTable は、デフォルトでは、更新する前に常にテーブルのアーカイブされたバージョンを作成します。ただし、skipArchive が true に設定されている場合、UpdateTable ではアーカイブされたバージョンは作成されません。

- TransactionId – UTF-8 文字列、1~255 バイト長、[Custom string pattern #16](#) に一致。

テーブルの内容を更新するトランザクション ID。

- VersionId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの内容を更新するバージョン ID。

- ViewUpdateAction – UTF-8 文字列 (有効な値: ADD | REPLACE | ADD\_OR\_REPLACE | DROP)。

ビューの更新時に実行されるオペレーション。

- Force – ブール。

一致するストレージ記述子とサブオブジェクト一致要件を無視するように true に設定できるフラグ。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- ConcurrentModificationException
- ResourceNumberLimitExceededException
- GlueEncryptionException
- ResourceNotReadyException

## DeleteTable アクション (Python: delete\_table)

データカタログからテーブル定義を削除します。

**Note**

このオペレーションが完了すると、削除されたテーブルに属するテーブルのバージョンとパーティションにはアクセスできなくなります。AWS Glue では、このサービスの裁量により、これらの「孤立した」リソースを適時に非同期的に削除します。

すべての関連リソースをすぐに削除するには、DeleteTableVersion または BatchDeleteTableVersion を使用し、DeletePartition または BatchDeletePartition を使用して、テーブルに属するすべてのリソースを削除してから、DeleteTable を呼び出します。

**リクエスト**

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログデータベースの名前。Hive 互換性のために、この名前はすべて小文字であることが必要です。

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除するテーブルの名前。Hive 互換性のために、この名前はすべて小文字であることが必要です。

- TransactionId – UTF-8 文字列、1~255 バイト長、[Custom string pattern #16](#) に一致。

テーブルの内容を更新するトランザクション ID。

**レスポンス**

- 応答パラメータはありません。

**エラー**

- EntityNotFoundException
- InvalidInputException

- `InternalServiceException`
- `OperationTimeoutException`
- `ConcurrentModificationException`
- `ResourceNotReadyException`

## BatchDeleteTable アクション (Python: `batch_delete_table`)

一度に複数のテーブルを削除します。

### Note

このオペレーションが完了すると、削除されたテーブルに属するテーブルのバージョンとパーティションにはアクセスできなくなります。AWS Glue では、このサービスの裁量により、これらの「孤立した」リソースを適時に非同期的に削除します。

すべての関連リソースをすぐに削除するには、`DeleteTableVersion` または `BatchDeleteTableVersion` を使用し、`DeletePartition` または `BatchDeletePartition` を使用して、テーブルに属するすべてのリソースを削除してから、`BatchDeleteTable` を呼び出します。

## リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除するテーブルが存在するカタログデータベースの名前。Hive 互換性のために、この名前はすべて小文字であることが必要です。

- `TablesToDelete` – 必須: UTF-8 文字列の配列。文字列 100 個以下。

削除するテーブルのリスト。

- `TransactionId` – UTF-8 文字列、1~255 バイト長、[Custom string pattern #16](#) に一致。

テーブルの内容を更新するトランザクション ID。

## 応答

- Errors – [TableError](#) オブジェクトの配列。

指定されたテーブルの削除試行中に発生したエラーのリスト。

## エラー

- InvalidInputException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException
- ResourceNotReadyException

## GetTable アクション (Python: get\_table)

指定されたテーブルのデータカタログ内にある Table 定義を取得します。

### リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログのデータベースの名前。Hive 互換性のために、この名前はすべて小文字であることが必要です。

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

定義を取得するテーブルの名前です。Hive 互換性のために、この名前はすべて小文字であることが必要です。

- TransactionId – UTF-8 文字列、1~255 バイト長、[Custom string pattern #16](#) に一致。

テーブルの内容を更新するトランザクション ID。

- QueryAsOfTime – タイムスタンプ。

テーブルの内容を読み取る時点の時間。設定されていない場合は、最新のトランザクションコミット時間が使用されます。TransactionIdと一緒に指定することはできません。

## レスポンス

- Table – [Table](#) オブジェクト。

指定したテーブルを定義する Table オブジェクト。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException
- ResourceNotReadyException
- FederationSourceException
- FederationSourceRetryableException

## GetTables アクション (Python: get\_tables)

特定の Database 内の、一部またはすべてのテーブル定義を取得します。

### リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルをリスト表示するカタログ内のデータベース。Hive 互換性のために、この名前はすべて小文字であることが必要です。

- Expression – UTF-8 文字列。2,048 バイト長以下。[Single-line string pattern](#) に一致。

正規表現パターン。存在する場合、パターンに名前が一致するテーブルのみが返されます。

- NextToken – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

- MaxResults – 数値 (integer)。1 ~ 100。

1 回の応答で返されるテーブルの最大数。

- TransactionId – UTF-8 文字列、1 ~ 255 バイト長、[Custom string pattern #16](#) に一致。

テーブルの内容を更新するトランザクション ID。

- QueryAsOfTime – タイムスタンプ。

テーブルの内容を読み取る時点の時間。設定されていない場合は、最新のトランザクションコミット時間が使用されます。TransactionId と一緒に指定することはできません。

## 応答

- TableList – [Table](#) オブジェクトの配列。

リクエストされた Table オブジェクトのリスト。

- NextToken – UTF-8 文字列。

継続トークン (現在のリストセグメントが最後ではない場合に存在)。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException
- GlueEncryptionException
- FederationSourceException
- FederationSourceRetryableException

## GetTableVersion アクション (Python: get\_table\_version)

テーブルの指定されたバージョンを取得します。

### リクエスト

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- **DatabaseName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログのデータベース。Hive 互換性のために、この名前はすべて小文字である必要があります。

- **TableName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの名前。Hive 互換性のために、この名前はすべて小文字である必要があります。

- **VersionId** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

取得するテーブルバージョンの ID 値。VersionID は整数の文字列表現です。各バージョンは 1 ずつ増加します。

### レスポンス

- **TableVersion** – [TableVersion](#) オブジェクト。

リクエストされたテーブルのバージョン。

### エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException

## GetTableVersions アクション (Python: get\_table\_versions)

指定したテーブルで使用可能なバージョンを指定する文字列のリストを取得します。

### リクエスト

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- **DatabaseName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログのデータベース。Hive 互換性のために、この名前はすべて小文字である必要があります。

- **TableName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの名前。Hive 互換性のために、この名前はすべて小文字である必要があります。

- **NextToken** – UTF-8 文字列。

継続トークン (これが最初の呼び出しでない場合)。

- **MaxResults** – 数値 (integer)。1~100。

1 回の応答で返すテーブルバージョンの最大数。

### 応答

- **TableVersions** – [TableVersion](#) オブジェクトの配列。

指定したテーブルの使用可能なバージョンを指定する文字列のリスト。

- **NextToken** – UTF-8 文字列。

継続トークン (使用可能なバージョンのリストに最後のバージョンが含まれていない場合)。

### エラー

- `EntityNotFoundException`
- `InvalidInputException`
- `InternalServiceException`

- `OperationTimeoutException`
- `GlueEncryptionException`

## DeleteTableVersion アクション (Python: `delete_table_version`)

テーブルの指定されたバージョンを削除します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログのデータベース。Hive 互換性のために、この名前はすべて小文字であることが必要です。

- `TableName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの名前。Hive 互換性のために、この名前はすべて小文字であることが必要です。

- `VersionId` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除するテーブルバージョンの ID。VersionID は整数の文字列表現です。各バージョンは 1 ずつ増加します。

### レスポンス

- 応答パラメータはありません。

### エラー

- `EntityNotFoundException`
- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`

## BatchDeleteTableVersion アクション (Python: batch\_delete\_table\_version)

テーブルの指定されたバージョンのバッチを削除します。

### リクエスト

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- **DatabaseName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログのデータベース。Hive 互換性のために、この名前はすべて小文字である必要があります。

- **TableName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの名前。Hive 互換性のために、この名前はすべて小文字である必要があります。

- **VersionIds** – 必須: UTF-8 文字列の配列。文字列 100 個以下。

削除するバージョンの ID リスト。VersionId は整数の文字列表現です。各バージョンは 1 ずつ増加します。

### 応答

- **Errors** – [TableVersionError](#) オブジェクトの配列。

指定されたテーブルバージョンを削除しようとしているときに発生したエラーのリスト。

### エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## SearchTables アクション (Python: search\_tables)

テーブルメタデータのプロパティと親データベースに基づいて一連のテーブルを検索します。テキストまたはフィルタ条件で検索できます。

Lake Formation に定義されているセキュリティポリシーに基づいて、アクセスできるテーブルのみを取得できます。テーブルが返されるようにするには、少なくともテーブルへの読み取り専用アクセスが必要です。テーブル内の一部の列にアクセスできない場合、テーブルのリストが返されるときに、これらの列は検索されません。列にはアクセスできても、これら列内のデータにアクセスできない場合は、これらの列とこれらの列に関連付けられたメタデータが検索に含まれます。

### リクエスト

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

一意の識別子 (*account\_id* で構成)。

- **NextToken** – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

- **Filters** – [PropertyPredicate](#) オブジェクトの配列。

キーバリューペアのリスト、および検索結果のフィルタリングに使用する比較演算子。述語に一致するすべてのエンティティを返します。

[PropertyPredicate](#) 構造体の `Comparator` メンバーは時間フィールドにのみ使用され、他のフィールドタイプでは省略できます。また、`Key=Name` など文字列値を比較する場合、ファジーマッチアルゴリズムが使用されます。`Key` フィールド (例えば、`Name` フィールドの値) は、特定の句読点文字 (`-`、`:`、`#` など) でトークンに分割されます。次に、各トークンは [PropertyPredicate](#) の `Value` メンバーとエグザクトマッチで比較されます。例えば、`Key=Name` および `Value=link` の場合、`customer-link` や `xx-link-yy` という名前のテーブルは返されますが、`xxlinkyy` の場合は返されません。

- **SearchText** – `Value` 文字列、1024 バイト長以下。

テキスト検索に使用する文字列。

値を引用符で囲んで指定すると、値との完全一致に基づいてフィルタリングされます。

- **SortCriteria** – [SortCriterion](#) オブジェクトの配列。構造 1 個以下。

結果をフィールド名で昇順または降順でソートするための条件のリスト。

- `MaxResults` – 1 ~ 1000 の数値 (整数)。

1 回の応答で返されるテーブルの最大数。

- `ResourceShareType` – UTF-8 文字列 (有効な値: FOREIGN | ALL | FEDERATED)。

アカウントと共有されているテーブルを検索するように指定できます。指定できる値は、FOREIGN または ALL です。

- FOREIGN に設定されている場合、アカウントと共有されているテーブルが検索されます。
- ALL に設定されている場合、アカウントと共有されているテーブルとローカルアカウントのテーブルが検索されます。

## レスポンス

- `NextToken` – UTF-8 文字列。

継続トークン (現在のリストセグメントが最後ではない場合に存在)。

- `TableList` – [Table](#) オブジェクトの配列。

リクエストされた Table オブジェクトのリスト。SearchTables レスポンスは、ユーザーがアクセスできるテーブルのみを返します。

## エラー

- `InternalServiceException`
- `InvalidInputException`
- `OperationTimeoutException`

## GetPartitionIndexes アクション (Python: `get_partition_indexes`)

テーブルに関連付けられたパーティションインデックスを取得します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログの ID。

- `DatabaseName` – 必須: UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。

パーティションインデックスを取得するデータベースの名前を指定します。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションインデックスを取得するテーブルの名前を指定します。

- NextToken – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

## 応答

- PartitionIndexDescriptorList – [PartitionIndexDescriptor](#) オブジェクトの配列。

インデックス記述子のリスト。

- NextToken – UTF-8 文字列。

継続トークン (現在のリストセグメントが最後ではない場合に存在)。

## エラー

- InternalServiceException
- OperationTimeoutException
- InvalidInputException
- EntityNotFoundException
- ConflictException

## CreatePartitionIndex アクション (Python: create\_partition\_index)

既存のテーブルに指定されたパーティションインデックスを作成します。

### リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログの ID。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションインデックスを作成するデータベースの名前を指定します。

- `TableName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションインデックスを作成するテーブルの名前を指定します。

- `PartitionIndex` – 必須: [PartitionIndex](#) オブジェクト。

`PartitionIndex` 構造を指定して、既存のテーブルにパーティションインデックスを作成します。

## レスポンス

- 応答パラメータはありません。

## エラー

- `AlreadyExistsException`
- `InvalidInputException`
- `EntityNotFoundException`
- `ResourceNumberLimitExceededException`
- `InternalServiceException`
- `OperationTimeoutException`
- `GlueEncryptionException`

## DeletePartitionIndex アクション (Python: `delete_partition_index`)

既存のテーブルから指定されたパーティションインデックスを削除します。

## リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログの ID。

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションインデックスを削除するデータベースの名前を指定します。

- `TableName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションインデックスを削除するテーブルの名前を指定します。

- IndexName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
削除するパーティションインデックスの名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- InternalServiceException
- OperationTimeoutException
- InvalidInputException
- EntityNotFoundException
- ConflictException
- GlueEncryptionException

## GetColumnStatisticsForTable アクション (Python: `get_column_statistics_for_table`)

列のテーブル統計を取得します。

このオペレーションに必要な Identity and Access Management (IAM) アクセス許可は GetTable です。

## リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
該当するパーティションが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。
- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
パーティションが存在するカタログデータベースの名前。
- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
パーティションのテーブルの名前。
- ColumnNames – 必須: UTF-8 文字列の配列。文字列 100 個以下。

列名のリスト。

応答

- ColumnStatisticsList – [ColumnStatistics](#) オブジェクトの配列。

のリスト ColumnStatistics。

- Errors – [ColumnError](#) オブジェクトの配列。

取得に ColumnStatistics 失敗した のリスト。

エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException

UpdateColumnStatisticsForTable アクション (Python:  
update\_column\_statistics\_for\_table)

列のテーブル統計を作成または更新します。

このオペレーションに必要な Identity and Access Management (IAM) アクセス許可は UpdateTable です。

リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当するパーティションが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが存在するカタログデータベースの名前。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションのテーブルの名前。

- ColumnStatisticsList – 必須: [ColumnStatistics](#) オブジェクトの配列。構造 25 個以下。

列の統計のリスト。

## 応答

- Errors – [ColumnStatisticsError](#) オブジェクトの配列。

のリスト ColumnStatisticsErrors。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException

## DeleteColumnStatisticsForTable アクション (Python: delete\_column\_statistics\_for\_table)

列のテーブル統計を取得します。

このオペレーションに必要な Identity and Access Management (IAM) アクセス許可は DeleteTable です。

## リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当するパーティションが存在するデータカタログの ID。何も指定しない場合、AWS アカウント ID はデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが存在するカタログデータベースの名前。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションのテーブルの名前。

- ColumnName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

列の名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException

## パーティション API

Partition API では、パーティションを使用するためのデータ型とオペレーションについて説明します。

### データ型

- [パーティションの構造](#)
- [PartitionInput の構造](#)
- [PartitionSpecWithSharedStorageDescriptor 構造](#)
- [PartitionListComposingSpec 構造](#)
- [PartitionSpecProxy 構造](#)
- [PartitionValueList 構造](#)
- [セグメント構造](#)
- [PartitionError 構造](#)
- [BatchUpdatePartitionFailureEntry 構造](#)
- [BatchUpdatePartitionRequestEntry 構造](#)

- [StorageDescriptor](#) 構造
- [SchemaReference](#) 構造
- [SerDeInfo](#) 構造
- [SkewedInfo](#) 構造

## パーティションの構造

テーブルデータのスライスを表します。

フィールド

- Values – UTF-8 文字列の配列。  
パーティションの値。
- DatabaseName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
パーティションを作成するカタログデータベースの名前。
- TableName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
パーティションを作成するデータベーステーブルの名前。
- CreationTime – タイムスタンプ。  
パーティションが作成された時刻。
- LastAccessTime – タイムスタンプ。  
パーティションが最後にアクセスされた時刻。
- StorageDescriptor – [StorageDescriptor](#) オブジェクト。  
パーティションが格納されている物理的な場所に関する情報を提供します。
- Parameters – キーバリュールールのマップ配列。  
各キーはキー文字列。1~255 バイト長。[Single-line string pattern](#) に一致。  
各値は UTF-8 文字列で、512,000 バイト長以下です。  
これらのキー/値ペアはパーティションパラメータを定義します。
- LastAnalyzedTime – タイムスタンプ。  
このパーティションの列統計が最後に計算された時刻。

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションがある Data Catalog の ID。

## PartitionInput の構造

パーティションの作成と更新に使用される構造。

### フィールド

- Values – UTF-8 文字列の配列。

パーティションの値。SDK ではこのパラメータは必要ありませんが、有効な入力としてこのパラメータを指定する必要があります。

新しいパーティションのキーの値は、文字列オブジェクトの配列として渡す必要があります。この配列は、Amazon S3 プレフィックスに表示されるパーティションキーと同じ順序で順序付けする必要があります。そうしないと、AWS Glue では、間違ったキーに値が追加されます。

- LastAccessTime – タイムスタンプ。

パーティションが最後にアクセスされた時刻。

- StorageDescriptor – [StorageDescriptor](#) オブジェクト。

パーティションが格納されている物理的な場所に関する情報を提供します。

- Parameters – キーバリューペアのマップ配列。

各キーはキー文字列。1~255 バイト長。[Single-line string pattern](#) に一致。

各値は UTF-8 文字列で、512,000 バイト長以下です。

これらのキー/値ペアはパーティションパラメータを定義します。

- LastAnalyzedTime – タイムスタンプ。

このパーティションの列統計が最後に計算された時刻。

## PartitionSpecWithSharedStorageDescriptor 構造

物理的な場所を共有するパーティションのパーティション仕様。

## フィールド

- StorageDescriptor – [StorageDescriptor](#) オブジェクト。  
共有物理ストレージ情報。
- Partitions – [パーティション](#) オブジェクトの配列。  
この物理的な場所を共有するパーティションのリスト。

## PartitionListComposingSpec 構造

関連するパーティションを一覧表示します。

### フィールド

- Partitions – [パーティション](#) オブジェクトの配列。  
構成仕様のパーティションのリスト。

## PartitionSpecProxy 構造

指定されたパーティションへのルートパスを提供します。

### フィールド

- DatabaseName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
パーティションが存在するカタログデータベース。
- TableName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
パーティションを含むテーブルの名前。
- RootPath – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
パーティションをアドレス指定するためのプロキシのルートパス。
- PartitionSpecWithSharedSD – [PartitionSpecWithSharedStorageDescriptor](#) オブジェクト。  
同じ物理ストレージの場所を共有するパーティションの仕様。
- PartitionListComposingSpec – [PartitionListComposingSpec](#) オブジェクト。  
パーティションのリストを指定します。

## PartitionValueList 構造

パーティションを定義する値のリストが含まれています。

フィールド

- `Values` – 必須: UTF-8 文字列の配列。

値のリスト。

## セグメント構造

テーブルのパーティションの重複しないリージョンを定義し、複数のリクエストを並行して実行できるようにします。

フィールド

- `SegmentNumber` – 必須: 数値 (integer)。None 以下。

セグメントのゼロベースのインデックス番号。たとえば、セグメントの合計数が 4 の場合、`SegmentNumber` 値の範囲は 0 ~ 3 です。

- `TotalSegments` – 必須: 数値 (整数)。1 ~ 10。

セグメントの合計数。

## PartitionError 構造

パーティションのエラーに関する情報が含まれます。

フィールド

- `PartitionValues` – UTF-8 文字列の配列。

パーティションを定義する値。

- `ErrorDetail` – [ErrorDetail](#) オブジェクト。

パーティションエラーの詳細。

## BatchUpdatePartitionFailureEntry 構造

バッチ更新パーティションエラーに関する情報が含まれます。

フィールド

- PartitionValueList - UTF-8 文字列の配列、文字列 100 個以下。

パーティションを定義する値のリスト。

- ErrorDetail - [ErrorDetail](#) オブジェクト。

バッチ更新パーティションエラーの詳細。

## BatchUpdatePartitionRequestEntry 構造

パーティションの更新に使用される値と構造体を含む構造体。

フィールド

- PartitionValueList - 必須: UTF-8 文字列の配列。文字列 100 個以下。

パーティションを定義する値のリスト。

- PartitionInput - 必須: [PartitionInput](#) オブジェクト。

パーティションの更新に使用される構造体。

## StorageDescriptor 構造

テーブルデータの物理ストレージについて説明します。

フィールド

- Columns - [Column](#) オブジェクトの配列。

テーブル内の Columns のリストです。

- Location - 場所文字列。2,056 バイト長以下。 [URI address multi-line string pattern](#) に一致。

テーブルの物理的な場所。デフォルトでは、ウェアハウスの場所、その後にウェアハウス内のデータベースの場所、その後にテーブル名が続く形式になります。

- AdditionalLocations - UTF-8 文字列の配列。

Delta テーブルへのパスを指す場所のリスト。

- InputFormat - 書式設定文字列。128 バイト長以下。 [Single-line string pattern](#) に一致。

入力形式: SequenceFileInputFormat (バイナリ)、または TextInputFormat、もしくはカスタム形式。

- OutputFormat - 書式設定文字列。128 バイト長以下。 [Single-line string pattern](#) に一致。

出力形式: SequenceFileOutputFormat (バイナリ)、または IgnoreKeyTextOutputFormat、もしくはカスタム形式。

- Compressed - ブール。

テーブルのデータが圧縮されている場合は True、圧縮されていない場合は False。

- NumberOfBuckets - 数値 (整数)。

テーブルにディメンション列が含まれている場合、指定する必要があります。

- SerdeInfo - [SerDeInfo](#) オブジェクト。

シリアライズ/デシリアライズ (SerDe) 情報。

- BucketColumns - UTF-8 文字列の配列。

テーブルのリデューサーグループ化列、クラスター列、およびバケット列のリスト。

- SortColumns - [Order](#) オブジェクトの配列。

テーブル内の各バケットのソート順を指定するリスト。

- Parameters - キーバリューペアのマップ配列。

各キーはキー文字列。1~255 バイト長。 [Single-line string pattern](#) に一致。

各値は UTF-8 文字列で、512,000 バイト長以下です。

ユーザーが指定したキーバリュー形式のプロパティ。

- SkewedInfo - [SkewedInfo](#) オブジェクト。

列に高い頻度で表示される値 (歪んだ値) に関する情報。

- StoredAsSubDirectories - ブール。

テーブルデータがサブディレクトリに保管されている場合は True、保管されていない場合は False。

- SchemaReference – [SchemaReference](#) オブジェクト。

AWS Glue スキーマレジストリに保存されているスキーマを参照するオブジェクト。

テーブルを作成するときにスキーマに空の列リストを渡し、代わりにスキーマリファレンスを使用します。

## SchemaReference 構造

AWS Glue スキーマレジストリに保存されているスキーマを参照するオブジェクト。

### フィールド

- SchemaId – [Schemald](#) オブジェクト。

スキーマアイデンティティフィールドを含む構造。こちらか SchemaVersionId のいずれかを指定する必要があります。

- SchemaVersionId – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。

スキーマのバージョンに割り当てられた一意の ID。こちらか SchemaId のいずれかを指定する必要があります。

- SchemaVersionNumber – 数値 (long)。1 ~ 100000。

スキーマのバージョン番号。

## SerDeInfo 構造

エクストラクターおよびローダーとして機能するシリアライズ/デシリアライズプログラム (SerDe) に関する情報。

### フィールド

- Name – UTF-8 文字列、1 ~ 255 バイト長、 [Single-line string pattern](#) に一致。

SerDe 名。

- SerializationLibrary – UTF-8 文字列、1 ~ 255 バイト長、 [Single-line string pattern](#) に一致。

通常、SerDe を実装するクラス。例は

`org.apache.hadoop.hive.serde2.columnar.ColumnarSerDe` です。

- **Parameters** – キーバリューペアのマップ配列。

各キーはキー文字列。1~255 バイト長。 [Single-line string pattern](#) に一致。

各値は UTF-8 文字列で、512,000 バイト長以下です。

これらのキーバリューペアは SerDe の初期化パラメータを定義します。

## SkewedInfo 構造

テーブルで歪んだ値を指定します。「歪んだ値」とは、非常に高い頻度で発生する値です。

### フィールド

- **SkewedColumnNames** – UTF-8 文字列の配列。

歪んだ値を含む列名のリスト。

- **SkewedColumnValues** – UTF-8 文字列の配列。

頻繁に出現するため、歪んだとみなされる値のリスト。

- **SkewedColumnValueLocationMaps** – キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

歪んだ値が含まれている列へのマッピング。

## 操作

- [CreatePartition アクション \(Python: create\\_partition\)](#)
- [BatchCreatePartition アクション \(Python: batch\\_create\\_partition\)](#)
- [UpdatePartition アクション \(Python: update\\_partition\)](#)
- [DeletePartition アクション \(Python: delete\\_partition\)](#)
- [BatchDeletePartition アクション \(Python: batch\\_delete\\_partition\)](#)
- [GetPartition アクション \(Python: get\\_partition\)](#)
- [GetPartitions アクション \(Python: get\\_partitions\)](#)
- [BatchGetPartition アクション \(Python: batch\\_get\\_partition\)](#)

- [BatchUpdatePartition アクション \(Python: batch\\_update\\_partition\)](#)
- [GetColumnStatisticsForPartition アクション \(Python: get\\_column\\_statistics\\_for\\_partition\)](#)
- [UpdateColumnStatisticsForPartition アクション \(Python: update\\_column\\_statistics\\_for\\_partition\)](#)
- [DeleteColumnStatisticsForPartition アクション \(Python: delete\\_column\\_statistics\\_for\\_partition\)](#)

## CreatePartition アクション (Python: create\_partition)

新しいパーティションを作成します。

### リクエスト

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
パーティションを作成するカタログの AWS アカウント ID。
- **DatabaseName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
パーティションが作成されるメタデータデータベースの名前。
- **TableName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
パーティションが作成されるメタデータテーブルの名前。
- **PartitionInput** – 必須: [PartitionInput](#) オブジェクト。  
作成されるパーティションを定義する PartitionInput 構造。

### レスポンス

- 応答パラメータはありません。

### エラー

- `InvalidInputException`
- `AlreadyExistsException`
- `ResourceNumberLimitExceededException`
- `InternalServiceException`
- `EntityNotFoundException`
- `OperationTimeoutException`

- `GlueEncryptionException`

## BatchCreatePartition アクション (Python: `batch_create_partition`)

バッチオペレーションで1つ以上のパーティションを作成します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションを作成するカタログの ID。現在、これは AWS アカウント ID である必要があります。

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが作成されるメタデータデータベースの名前。

- `TableName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが作成されるメタデータテーブルの名前。

- `PartitionInputList` – 必須: [PartitionInput](#) オブジェクトの配列。構造 100 個以下。

作成されるパーティションを定義する `PartitionInput` 構造のリスト。

### 応答

- `Errors` – [PartitionError](#) オブジェクトの配列。

リクエストされたパーティションを作成しようとしたときにエラーが発生しました。

### エラー

- `InvalidInputException`
- `AlreadyExistsException`
- `ResourceNumberLimitExceededException`
- `InternalServiceException`
- `EntityNotFoundException`
- `OperationTimeoutException`
- `GlueEncryptionException`

## UpdatePartition アクション (Python: update\_partition)

パーティションを更新します。

### リクエスト

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

更新されるパーティションが存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- **DatabaseName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当する表が存在するカタログデータベースの名前。

- **TableName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

更新するパーティションが存在するテーブルの名前。

- **PartitionValueList** – 必須: UTF-8 文字列の配列。文字列 100 個以下。

更新するパーティションを定義するパーティションキー値のリスト。

- **PartitionInput** – 必須: [PartitionInput](#) オブジェクト。

パーティションを更新する新しいパーティションオブジェクト。

Values プロパティは変更できません。パーティションのパーティションキー値を変更する場合は、パーティションを削除して再作成します。

### レスポンス

- 応答パラメータはありません。

### エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException

## DeletePartition アクション (Python: delete\_partition)

指定したパーティションを削除します。

### リクエスト

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除されるパーティションが存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- **DatabaseName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当する表が存在するカタログデータベースの名前。

- **TableName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除するパーティションを含むテーブルの名前。

- **PartitionValues** – 必須: UTF-8 文字列の配列。

パーティションを定義する値。

### レスポンス

- 応答パラメータはありません。

### エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## BatchDeletePartition アクション (Python: batch\_delete\_partition)

バッチオペレーションで1つ以上のパーティションを削除します。

### リクエスト

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除されるパーティションが存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当する表が存在するカタログデータベースの名前。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除するパーティションを含むテーブルの名前。

- PartitionsToDelete – 必須: [PartitionValueList](#) オブジェクトの配列。構造 25 個以下。

削除されるパーティションを定義する PartitionInput 構造のリスト。

## 応答

- Errors – [PartitionError](#) オブジェクトの配列。

リクエストされたパーティションを削除しようとしたときにエラーが発生しました。

## エラー

- InvalidInputException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException

## GetPartition アクション (Python: get\_partition)

指定したパーティションに関する情報を取得します。

### リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

対象のパーティションが存在する Data Catalog の ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが存在するカタログデータベースの名前。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションのテーブルの名前。

- PartitionValues – 必須: UTF-8 文字列の配列。

パーティションを定義する値。

## レスポンス

- Partition – [パーティション](#) オブジェクト。

Partition オブジェクトの形式で、リクエストされた情報。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException
- FederationSourceException
- FederationSourceRetryableException

## GetPartitions アクション (Python: get\_partitions)

テーブルのパーティションについての情報を取得します。

### リクエスト

- CatalogId – カatalog ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当するパーティションが存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが存在するカタログデータベースの名前。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションのテーブルの名前。

- Expression – 述語文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

返されるパーティションをフィルタリングする式。

この式では、SQL WHERE フィルタ句と似た SQL 構文を使用します。SQL ステートメントパーサー [JSQLParser](#) が式を解析します。

演算子: Expression API 呼び出しで使用できる演算子を以下に示します。

=

2つのオペランドの値が等しいかどうかを確認します。「はい」の場合、条件は true です。

例: 「変数 a」が 10、「変数 b」が 20 とします。

(a = b) は true ではありません。

<>

2つのオペランドの値が等しいかどうかを確認します。値が等しくない場合、条件は true です。

例: (a <> b) は true です。

>

左のオペランドの値が右のオペランドの値よりも大きいかどうかを確認します。「はい」の場合、条件は true です。

例: (a > b) は true ではありません。

<

左のオペランドの値が右のオペランドの値よりも小さいかどうかを確認します。「はい」の場合、条件は true です。

例: (a < b) は true です。

>=

左のオペランドの値が右のオペランドの値以上かどうかを確認します。「はい」の場合、条件は true です。

例: (a >= b) は true ではありません。

<=

左のオペランドの値が右のオペランドの値以下かどうかを確認します。「はい」の場合、条件は true です。

例: (a <= b) は true です。

AND、OR、IN、BETWEEN、LIKE、NOT、IS NULL

論理演算子。

サポートされているパーティションキーの型: サポートされているパーティションキーを以下に示します。

- string
- date
- timestamp
- int
- bigint
- long
- tinyint
- smallint
- decimal

無効なタイプが存在した場合は、例外が発生します。

各タイプの有効な演算子のリストを次に示します。クローラを定義する場合、partitionKey タイプはカタログパーティションとの互換性を確保するため STRING として作成されます。

サンプル API コール:

Example

```
year = 2015
 year = 2016
 year = 2017
```

### Example

パーティション `year` を 2015 と等しくする

```
aws glue get-partitions --database-name dbname --table-name twitter_partition
 --expression "year=*'2015'"
```

### Example

パーティション `year` を 2016-2018 (これらの値を含まない) にする

```
aws glue get-partitions --database-name dbname --table-name twitter_partition
 --expression "year>'2016' AND year<'2018'"
```

### Example

パーティション `year` を 2015-2018 (これらの値を含む) にする 次の API コールは、相互に等しくなります。

```
aws glue get-partitions --database-name dbname --table-name twitter_partition
 --expression "year>='2015' AND year<='2018'"

aws glue get-partitions --database-name dbname --table-name
twitter_partition
 --expression "year BETWEEN 2015 AND 2018"

aws glue get-partitions --database-name dbname --table-name
twitter_partition
 --expression "year IN (2015,2016,2017,2018)"
```

### Example

ワイルドカードパーティションフィルタ。次のコール出力がパーティション `year=2017` になります。正規表現は LIKE ではサポートされません。

```
aws glue get-partitions --database-name dbname --table-name twitter_partition
--expression "year LIKE '%7'"
```

- NextToken – UTF-8 文字列。

これらのパーティションを取得する最初の呼び出しでない場合は、継続トークン。

- Segment – [Segment](#) オブジェクト。

このリクエストでスキャンするテーブルのパーティションのセグメント。

- MaxResults – 1 ~ 1000 の数値 (整数)。

1 回の応答で返されるパーティションの最大数。

- ExcludeColumnSchema – ブール。

true の場合、パーティション列スキーマを返さないことを指定します。パーティション値や場所など、他のパーティション属性にのみ関心がある場合に便利です。この方法では、重複データを返さないことで、大きな応答の問題を回避できます。

- TransactionId – UTF-8 文字列、1 ~ 255 バイト長、[Custom string pattern #16](#) に一致。

テーブルの内容を更新するトランザクション ID。

- QueryAsOfTime – タイムスタンプ。

パーティションの内容を読み取る時点の時間。設定されていない場合は、最新のトランザクションコミット時間が使用されます。TransactionId と一緒に指定することはできません。

## 応答

- Partitions – [パーティション](#) オブジェクトの配列。

リクエストされたパーティションのリスト。

- NextToken – UTF-8 文字列。

継続トークン (戻されたパーティションのリストに最後のパーティションが含まれていない場合)。

## エラー

- EntityNotFoundException

- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`
- `GlueEncryptionException`
- `InvalidStateException`
- `ResourceNotReadyException`
- `FederationSourceException`
- `FederationSourceRetryableException`

## BatchGetPartition アクション (Python: `batch_get_partition`)

バッチリクエストのパーティションを取得します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当するパーティションが存在するデータカタログの ID。提供されない場合は、AWS アカウント ID がデフォルトで使用されます。

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが存在するカタログデータベースの名前。

- `TableName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションのテーブルの名前。

- `PartitionsToGet` – 必須: [PartitionValueList](#) オブジェクトの配列。構造 1000 個以下。

取得するパーティションを識別するパーティション値のリスト。

### 応答

- `Partitions` – [パーティション](#) オブジェクトの配列。

リクエストされたパーティションのリスト。

- `UnprocessedKeys` – [PartitionValueList](#) オブジェクトの配列。構造 1000 個以下。

パーティションが返されなかったリクエスト内のパーティション値のリスト。

## エラー

- `InvalidInputException`
- `EntityNotFoundException`
- `OperationTimeoutException`
- `InternalServiceException`
- `GlueEncryptionException`
- `InvalidStateException`
- `FederationSourceException`
- `FederationSourceRetryableException`

## BatchUpdatePartition アクション (Python: `batch_update_partition`)

バッチ操作で 1 つ以上のパーティションを更新します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションを更新するカタログの ID。現在、これは AWS アカウント ID である必要がありません。

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが更新されるメタデータデータベースの名前。

- `TableName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが更新されるメタデータテーブルの名前。

- `Entries` – 必須: [BatchUpdatePartitionRequestEntry](#) オブジェクトの配列。1~100 個の構造。

更新する最大 100 個の `BatchUpdatePartitionRequestEntry` オブジェクトのリスト。

### 応答

- `Errors` – [BatchUpdatePartitionFailureEntry](#) オブジェクトの配列。

リクエストされたパーティションを更新しようとしたときにエラーが発生しました。 `BatchUpdatePartitionFailureEntry` オブジェクトのリスト。

## エラー

- `InvalidInputException`
- `EntityNotFoundException`
- `OperationTimeoutException`
- `InternalServiceException`
- `GlueEncryptionException`

## GetColumnStatisticsForPartition アクション (Python: `get_column_statistics_for_partition`)

列のパーティション統計を取得します。

このオペレーションに必要な Identity and Access Management (IAM) アクセス許可は `GetPartition` です。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当するパーティションが存在するデータカタログの ID。提供されない場合は、AWS アカウント ID がデフォルトで使用されます。

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが存在するカタログデータベースの名前。

- `TableName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションのテーブルの名前。

- `PartitionValues` – 必須: UTF-8 文字列の配列。

パーティションを識別するパーティション値のリスト。

- `ColumnNames` – 必須: UTF-8 文字列の配列。文字列 100 個以下。

列名のリスト。

### 応答

- `ColumnStatisticsList` – [ColumnStatistics](#) オブジェクトの配列。

取得に失敗した ColumnStatistics のリスト。

- Errors – [ColumnError](#) オブジェクトの配列。

列の統計データの取得中にエラーが発生しました。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException

## UpdateColumnStatisticsForPartition アクション (Python: update\_column\_statistics\_for\_partition)

カラムのパーティション統計を作成または更新します。

このオペレーションに必要な Identity and Access Management (IAM) アクセス許可は UpdatePartition です。

## リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当するパーティションが存在するデータカタログの ID。提供されない場合は、AWS アカウント ID がデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが存在するカタログデータベースの名前。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションのテーブルの名前。

- PartitionValues – 必須: UTF-8 文字列の配列。

パーティションを識別するパーティション値のリスト。

- ColumnStatisticsList – 必須: [ColumnStatistics](#) オブジェクトの配列。構造 25 個以下。

列の統計のリスト。

応答

- Errors – [ColumnStatisticsError](#) オブジェクトの配列。

列の統計データの更新中にエラーが発生しました。

エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException

DeleteColumnStatisticsForPartition アクション (Python: delete\_column\_statistics\_for\_partition)

列のパーティション列統計を削除します。

このオペレーションに必要な Identity and Access Management (IAM) アクセス許可は DeletePartition です。

リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

該当するパーティションが存在するデータカタログの ID。提供されない場合は、AWS アカウント ID がデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションが存在するカタログデータベースの名前。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

パーティションのテーブルの名前。

- PartitionValues – 必須: UTF-8 文字列の配列。

パーティションを識別するパーティション値のリスト。

- ColumnName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

列の名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- GlueEncryptionException

## 接続 API

接続 API では、AWS Glue 接続データ型と、接続を作成、削除、更新、および一覧表示するための API について説明します。

## データ型

- [Connection 構造](#)
- [ConnectionInput 構造](#)
- [PhysicalConnectionRequirements 構造](#)
- [GetConnectionsFilter 構造](#)

## Connection 構造

データソースへの接続を定義します。

### フィールド

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

接続定義の名前。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

接続の説明。

- ConnectionType – UTF-8 文字列 (有効な値: JDBC | SFTP | MONGODB | KAFKA | NETWORK | MARKETPLACE | CUSTOM | SALESFORCE)。

接続のタイプ。現在、SFTP はサポートされていません。

- MatchCriteria - UTF-8 文字列の配列。文字列 10 個以下。

この接続を選択する際に使用可能な条件のリスト。

- ConnectionProperties – キーと値のペアのマップ配列。100 ペア以下。

各キーは UTF-8 文字列です (有効な値: HOST | PORT | USERNAME="USER\_NAME" | PASSWORD | ENCRYPTED\_PASSWORD | JDBC\_DRIVER\_JAR\_URI | JDBC\_DRIVER\_CLASS\_NAME | JDBC\_ENGINE | JDBC\_ENGINE\_VERSION | CONFIG\_FILES | INSTANCE\_ID | JDBC\_CONNECTION\_URL | JDBC\_ENFORCE\_SSL | CUSTOM\_JDBC\_CERT | SKIP\_CUSTOM\_JDBC\_CERT\_VALIDATION | CUSTOM\_JDBC\_CERT\_STRING | CONNECTION\_URL | KAFKA\_BOOTSTRAP\_SERVERS | KAFKA\_SSL\_ENABLED | KAFKA\_CUSTOM\_CERT | KAFKA\_SKIP\_CUSTOM\_CERT\_VALIDATION | KAFKA\_CLIENT\_KEYSTORE | KAFKA\_CLIENT\_KEYSTORE\_PASSWORD | KAFKA\_CLIENT\_KEY\_PASSWORD | ENCRYPTED\_KAFKA\_CLIENT\_KEYSTORE\_PASSWORD | ENCRYPTED\_KAFKA\_CLIENT\_KEY\_PASSWORD | SECRET\_ID | CONNECTOR\_URL | CONNECTOR\_TYPE | CONNECTOR\_CLASS\_NAME | KAFKA\_SASL\_MECHANISM | KAFKA\_SASL\_PLAIN\_USERNAME | KAFKA\_SASL\_PLAIN\_PASSWORD | ENCRYPTED\_KAFKA\_SASL\_PLAIN\_PASSWORD | KAFKA\_SASL\_SCRAM\_USERNAME | KAFKA\_SASL\_SCRAM\_PASSWORD | KAFKA\_SASL\_SCRAM\_SECRETS\_ARN | ENCRYPTED\_KAFKA\_SASL\_SCRAM\_PASSWORD | KAFKA\_SASL\_GSSAPI\_KEYTAB | KAFKA\_SASL\_GSSAPI\_KRB5\_CONF | KAFKA\_SASL\_GSSAPI\_SERVICE | KAFKA\_SASL\_GSSAPI\_PRINCIPAL | ROLE\_ARN)。

各値は Value 文字列で、1,024 バイト長以下です。

以下のキー/値ペアは接続のパラメータを定義します。

- HOST - ホスト URI: データベースホストの完全修飾ドメイン名 (FQDN) または IPv4 アドレス。
- PORT - データベースホストがデータベース接続をリッスンするポート番号 (1024 ~ 65535)。

- **USER\_NAME** - データベースにログインするための名前。USER\_NAME の値文字列は「USERNAME」です。
- **PASSWORD** - ユーザー名に対応するパスワード (使用されている場合)。
- **ENCRYPTED\_PASSWORD** - データカタログの暗号化設定の設定 (ConnectionPasswordEncryption) で接続パスワードの保護が有効になっている場合、このフィールドには、暗号化されたパスワードが保存されます。
- **JDBC\_DRIVER\_JAR\_URI** - 使用する JDBC ドライバーを含む JAR ファイルの Amazon Simple Storage Service (Amazon S3) パス。
- **JDBC\_DRIVER\_CLASS\_NAME** - 使用する JDBC ドライバーのクラス名。
- **JDBC\_ENGINE** - 使用する JDBC エンジンの名前。
- **JDBC\_ENGINE\_VERSION** - 使用する JDBC エンジンのバージョン。
- **CONFIG\_FILES** - (将来の利用のために予約されています。)
- **INSTANCE\_ID** - 使用するインスタンス ID。
- **JDBC\_CONNECTION\_URL** - JDBC データソースに接続するための URL。
- **JDBC\_ENFORCE\_SSL** - ホスト名のマッチによる Secure Sockets Layer (SSL) がクライアント上の JDBC 接続に適用されるかどうかを指定するブール文字列 (true、false)。デフォルトは False です。
- **CUSTOM\_JDBC\_CERT** - お客様のルート証明書を指定する Amazon S3 の場所。AWS Glue では、このルート証明書を使用して、お客様データベースに接続するときにお客様の証明書を検証します。AWS Glue で扱うのは、X.509 証明書のみです。提供する証明書は、DER エンコードし、Base64 エンコード PEM 形式で指定する必要があります。
- **SKIP\_CUSTOM\_JDBC\_CERT\_VALIDATION** - デフォルトでは、false です。AWS Glue は、お客様の証明書の署名アルゴリズムとサブジェクトパブリックキーアルゴリズムを検証します。署名アルゴリズムで許可されるアルゴリズムは、SHA256withRSA、SHA384withRSA、または SHA512withRSA のみです。サブジェクトパブリックキーアルゴリズムでは、キーの長さは 2048 以上である必要があります。このプロパティの値を true に設定して、AWS Glue によるお客様の証明書の検証をスキップできます。
- **CUSTOM\_JDBC\_CERT\_STRING** - 中間者 (man-in-the-middle) 攻撃を防ぐためにドメインの一致または識別名の一致に使用するカスタム JDBC 証明書文字列。Oracle データベースでは、これは SSL\_SERVER\_CERT\_DN として使用され、Microsoft SQL Server では hostNameInCertificate として使用されます。
- **CONNECTION\_URL** - 一般的な (JDBC 以外の) データソースに接続するための URL。
- **SECRET\_ID** - 認証情報のシークレットマネージャーに使用されるシークレット ID。

- CONNECTOR\_URL – MARKETPLACE または CUSTOM 接続のコネクタ URL。
- CONNECTOR\_TYPE – MARKETPLACE または CUSTOM 接続のコネクタの種類。
- CONNECTOR\_CLASS\_NAME – MARKETPLACE または CUSTOM 接続のコネクタクラス名。
- KAFKA\_BOOTSTRAP\_SERVERS - Kafkaクライアントが接続してそれ自身をブートストラップする先である Kafkaクラスター内の Apache Kafka ブローカーのアドレスを示すホストとポートのペアのコンマ区切りリスト。
- KAFKA\_SSL\_ENABLED – Apache Kafka 接続で SSL を有効にするか無効にするか。デフォルト値は「true」です。
- KAFKA\_CUSTOM\_CERT – プライベート CA 証明書ファイル (.pem 形式) の Amazon S3 URL。デフォルトは空の文字列です。
- KAFKA\_SKIP\_CUSTOM\_CERT\_VALIDATION – CA 証明書ファイルの検証をスキップするかどうか。AWS Glue では、SHA256withRSA、SHA384withRSA、SHA512withRSA の 3 つのアルゴリズムについて検証されます。デフォルト値は「false」です。
- KAFKA\_CLIENT\_KEYSTORE – Kafka クライアント側認証用のクライアントキーストアファイルの Amazon S3 の場所 (オプション)。
- KAFKA\_CLIENT\_KEYSTORE\_PASSWORD – 指定されたキーストアにアクセスするためのパスワード (オプション)。
- KAFKA\_CLIENT\_KEY\_PASSWORD – キーストアは、複数のキーで構成することができるので、これは Kafka サーバー側のキーで使用するクライアントキーにアクセスするためのパスワードです (オプション)。
- ENCRYPTED\_KAFKA\_CLIENT\_KEYSTORE\_PASSWORD – Kafka クライアントキーストアのパスワードの暗号化されたバージョン (ユーザーが AWS Glue 暗号化パスワードの設定を選択している場合)。
- ENCRYPTED\_KAFKA\_CLIENT\_KEY\_PASSWORD – Kafka クライアントキーパスワードの暗号化されたバージョン (ユーザーが AWS Glue 暗号化パスワードの設定を選択している場合)。
- KAFKA\_SASL\_MECHANISM – "SCRAM-SHA-512"、"GSSAPI"、"AWS\_MSK\_IAM"、"PLAIN"。これらはサポートされている [SASL メカニズム](#)です。
- KAFKA\_SASL\_PLAIN\_USERNAME – 「PLAIN」メカニズムで認証するとき使用されるプレーンテキストのユーザーネーム。
- KAFKA\_SASL\_PLAIN\_PASSWORD – 「PLAIN」メカニズムで認証するとき使用されるプレーンテキストのパスワード。
- ENCRYPTED\_KAFKA\_SASL\_PLAIN\_PASSWORD – Kafka SASL PLAIN パスワードの暗号化バージョン (ユーザーが AWS Glue 暗号化パスワードの設定を選択した場合)。

- KAFKA\_SASL\_SCRAM\_USERNAME - 「SCRAM-SHA-512」メカニズムでの認証に使用されるプレーンテキストのユーザーネーム。
- KAFKA\_SASL\_SCRAM\_PASSWORD - 「SCRAM-SHA-512」メカニズムでの認証に使用されるプレーンテキストのパスワード。
- ENCRYPTED\_KAFKA\_SASL\_SCRAM\_PASSWORD - Kafka SASL SCRAM パスワードの暗号化されたバージョン (ユーザーが AWS Glue 暗号化パスワードの設定を選択している場合)。
- KAFKA\_SASL\_SCRAM\_SECRETS\_ARN - AWS Secrets Manager のシークレットの Amazon リソースネーム。
- KAFKA\_SASL\_GSSAPI\_KEYTAB - Kerberos keytab ファイルの S3 の場所。キータブには、1 つ以上のプリンシパルの長期キーが保存されます。詳細については、[MIT Kerberos ドキュメント: キータブ](#)を参照してください。
- KAFKA\_SASL\_GSSAPI\_KRB5\_CONF - Kerberos krb5.conf ファイルの S3 の場所。krb5.conf には、KDC サーバーの場所などの Kerberos 構成情報が保存されます。詳細については、[MIT Kerberos ドキュメント: krb5.conf](#)を参照してください。
- KAFKA\_SASL\_GSSAPI\_SERVICE - [Kafka 構成](#) の `sasl.kerberos.service.name` で設定した通りの Kerberos サービスの名前。
- KAFKA\_SASL\_GSSAPI\_PRINCIPAL - AWS Glue によって使用される Kerberos プリンシパルの名前。詳細については、[Kafka ドキュメント: Kafka ブローカーの構成](#)を参照してください。
- PhysicalConnectionRequirements - [PhysicalConnectionRequirements](#) オブジェクト。

この接続を正常に行うために必要な物理接続要件 (例: 仮想プライベートクラウド (VPC)、SecurityGroup)。

- CreationTime - タイムスタンプ。

この接続定義が作成された時刻のタイムスタンプ。

- LastUpdatedTime - タイムスタンプ。

この接続定義が更新された最終時刻のタイムスタンプ。

- LastUpdatedBy - UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

この接続定義を最終更新したユーザー、グループ、またはロール。

- Status - UTF-8 文字列 (有効な値: READY | IN\_PROGRESS | FAILED)。

接続のステータス。READY、IN\_PROGRESS、または FAILED のいずれか。

- StatusReason - UTF-8 文字列、1~16384 バイト長。

接続ステータスの理由。

- LastConnectionValidationTime – タイムスタンプ。

この接続が最後に検証された時刻のタイムスタンプ。

- AuthenticationConfiguration – [AuthenticationConfiguration](#) オブジェクト。

接続の認証プロパティ。

## ConnectionInput 構造

作成または更新する接続を指定するために使用される構造。

フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

コレクションの名前。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

接続の説明。

- ConnectionType – 必須: UTF-8 文字列 (有効な値: JDBC | SFTP | MONGODB | KAFKA | NETWORK | MARKETPLACE | CUSTOM | SALESFORCE)。

接続のタイプ。現在、次のタイプがサポートされています。

- JDBC - Java データベース接続 (JDBC) を介してデータベースへの接続を指定します。

JDBC 接続では、次の ConnectionParameters を使用します。

- 必須: (HOST、PORT、JDBC\_ENGINE) のすべて、または JDBC\_CONNECTION\_URL。

- 必須: (USERNAME、PASSWORD) の両方、または SECRET\_ID。

- オプション:

JDBC\_ENFORCE\_SSL、CUSTOM\_JDBC\_CERT、CUSTOM\_JDBC\_CERT\_STRING、SKIP\_CUSTOM\_JDB  
これらのパラメータは、JDBC で SSL を設定するために使用されます。

- KAFKA - Apache Kafka ストリーミングプラットフォームへの接続を指定します。

KAFKA 接続では、次の ConnectionParameters を使用します。

- 必須: KAFKA\_BOOTSTRAP\_SERVERS。

- オプション:  
KAFKA\_SSL\_ENABLED、KAFKA\_CUSTOM\_CERT、KAFKA\_SKIP\_CUSTOM\_CERT\_VALIDATION。  
これらのパラメータは KAFKA で SSL を設定するために使用されます。
- オプション:  
KAFKA\_CLIENT\_KEYSTORE、KAFKA\_CLIENT\_KEYSTORE\_PASSWORD、KAFKA\_CLIENT\_KEY\_PASS  
これらのパラメータは、KAFKA の中で SSL で TLS クライアント設定を設定するために使用  
されます。
- オプション: KAFKA\_SASL\_MECHANISM。SCRAM-SHA-512、GSSAPI、または AWS\_MSK\_IAM  
として指定できます。
- オプション:  
KAFKA\_SASL\_SCRAM\_USERNAME、KAFKA\_SASL\_SCRAM\_PASSWORD、ENCRYPTED\_KAFKA\_SASL\_S  
これらのパラメータは、KAFKA で SASL/SCRAM-SHA-512 認証を設定するために使用されま  
す。
- オプション:  
KAFKA\_SASL\_GSSAPI\_KEYTAB、KAFKA\_SASL\_GSSAPI\_KRB5\_CONF、KAFKA\_SASL\_GSSAPI\_SER  
これらのパラメータは、KAFKA で SASL/GSSAPI 認証を設定するために使用されます。
- MONGODB - MongoDB ドキュメントデータベースへの接続を指定します。

MONGODB 接続では、次の ConnectionParameters を使用します。

- 必須: CONNECTION\_URL。
- 必須: (USERNAME、PASSWORD) の両方、または SECRET\_ID。
- SALESFORCE - OAuth 認証を使用して Salesforce への接続を指定します。
  - AuthenticationConfiguration メンバーを設定する必要があります。
- NETWORK - Amazon Virtual Private Cloud 環境 (Amazon VPC) 内のデータソースへのネットワー  
ク接続を指定します。

NETWORK 接続には、ConnectionParameters は必要ありません。かわり  
に、PhysicalConnectionRequirements を用意します。

- MARKETPLACE – AWS Marketplace から購入したコネクタに含まれる設定を使用して、AWS  
Glue がネイティブにサポートしていないデータストアに対する読み取りと書き込みを行いま  
す。

MARKETPLACE 接続では、次の ConnectionParameters を使用します。

- 必須:

~~CONNECTOR\_TYPE、CONNECTOR\_URL、CONNECTOR\_CLASS\_NAME、CONNECTION\_URL。~~

- JDBC\_CONNECTOR\_TYPE 接続に必須: (USERNAME、PASSWORD) の両方、または SECRET\_ID。
- CUSTOM – カスタムコネクタに含まれる設定を使用して、AWS Glue がネイティブにサポートしていないデータストアに対する読み取りと書き込みを行います。

SFTP はサポートされていません。

オプションの ConnectionProperties を使用して AWS Glue の機能を設定する方法の詳細については、「[AWS Glue 接続プロパティ](#)」を参照してください。

オプションの ConnectionProperties を使用して AWS Glue Studio の機能を設定する方法の詳細については、「[コネクタと接続の使用](#)」を参照してください。

- MatchCriteria - UTF-8 文字列の配列。文字列 10 個以下。

この接続を選択する際に使用可能な条件のリスト。

- ConnectionProperties – 必須: キーと値のペアのマップ配列。100 ペア以下。

各キーは UTF-8 文字列です (有効な値: HOST | PORT | USERNAME="USER\_NAME" | PASSWORD | ENCRYPTED\_PASSWORD | JDBC\_DRIVER\_JAR\_URI | JDBC\_DRIVER\_CLASS\_NAME | JDBC\_ENGINE | JDBC\_ENGINE\_VERSION | CONFIG\_FILES | INSTANCE\_ID | JDBC\_CONNECTION\_URL | JDBC\_ENFORCE\_SSL | CUSTOM\_JDBC\_CERT | SKIP\_CUSTOM\_JDBC\_CERT\_VALIDATION | CUSTOM\_JDBC\_CERT\_STRING | CONNECTION\_URL | KAFKA\_BOOTSTRAP\_SERVERS | KAFKA\_SSL\_ENABLED | KAFKA\_CUSTOM\_CERT | KAFKA\_SKIP\_CUSTOM\_CERT\_VALIDATION | KAFKA\_CLIENT\_KEYSTORE | KAFKA\_CLIENT\_KEYSTORE\_PASSWORD | KAFKA\_CLIENT\_KEY\_PASSWORD | ENCRYPTED\_KAFKA\_CLIENT\_KEYSTORE\_PASSWORD | ENCRYPTED\_KAFKA\_CLIENT\_KEY\_PASSWORD | SECRET\_ID | CONNECTOR\_URL | CONNECTOR\_TYPE | CONNECTOR\_CLASS\_NAME | KAFKA\_SASL\_MECHANISM | KAFKA\_SASL\_PLAIN\_USERNAME | KAFKA\_SASL\_PLAIN\_PASSWORD | ENCRYPTED\_KAFKA\_SASL\_PLAIN\_PASSWORD | KAFKA\_SASL\_SCRAM\_USERNAME | KAFKA\_SASL\_SCRAM\_PASSWORD | KAFKA\_SASL\_SCRAM\_SECRETS\_ARN | ENCRYPTED\_KAFKA\_SASL\_SCRAM\_PASSWORD | KAFKA\_SASL\_GSSAPI\_KEYTAB | KAFKA\_SASL\_GSSAPI\_KRB5\_CONF | KAFKA\_SASL\_GSSAPI\_SERVICE | KAFKA\_SASL\_GSSAPI\_PRINCIPAL | ROLE\_ARN)。

各値は Value 文字列で、1,024 バイト長以下です。

以下のキーバリューペアは接続のパラメータを定義します。

- PhysicalConnectionRequirements – [PhysicalConnectionRequirements](#) オブジェクト。

この接続を正常に行うために必要な物理接続要件 (例: 仮想プライベートクラウド (VPC)、SecurityGroup)。

- AuthenticationConfiguration – [AuthenticationConfigurationInput](#) オブジェクト。

接続の認証プロパティ。Salesforce 接続に使用されます。

- ValidateCredentials – ブール。

接続の作成中に認証情報を検証するフラグ。Salesforce 接続に使用されます。デフォルトは true です。

## PhysicalConnectionRequirements 構造

GetConnection レスポンスの OAuth クライアントアプリ。

フィールド

- SubnetId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

接続で使用されるサブネット ID。

- SecurityGroupIdList - UTF-8 文字列の配列。文字列 50 個以下。

接続で使用されるセキュリティグループ ID のリスト。

- AvailabilityZone – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

接続のアベイラビリティゾーン。

## GetConnectionsFilter 構造

GetConnections API オペレーションによって返る接続定義をフィルタリングします。

フィールド

- MatchCriteria - UTF-8 文字列の配列。文字列 10 個以下。

接続定義を返すためにその接続定義に記録された条件と一致する必要がある条件文字列。

- ConnectionType – UTF-8 文字列 (有効な値: JDBC | SFTP | MONGODB | KAFKA | NETWORK | MARKETPLACE | CUSTOM | SALESFORCE)。

返す接続のタイプ。現在、SFTP はサポートされていません。

## 操作

- [CreateConnection アクション \(Python: create\\_connection\)](#)
- [DeleteConnection アクション \(Python: delete\\_connection\)](#)
- [GetConnection アクション \(Python: get\\_connection\)](#)
- [GetConnections アクション \(Python: get\\_connections\)](#)
- [UpdateConnection アクション \(Python: update\\_connection\)](#)
- [BatchDeleteConnection アクション \(Python: batch\\_delete\\_connection\)](#)

## CreateConnection アクション (Python: create\_connection)

データカタログで接続定義を作成します。

フェデレーションリソースの作成に使用される接続には IAM `glue:PassConnection` 権限が必要です。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

接続を作成するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- `ConnectionInput` – 必須: [ConnectionInput](#) オブジェクト。

作成する接続を定義する `ConnectionInput` オブジェクト。

- `Tags` – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

接続に割り当てるタグ。

## レスポンス

- `CreateConnectionStatus` – UTF-8 文字列 (有効な値: `READY` | `IN_PROGRESS` | `FAILED`)。

接続作成リクエストのステータス。リクエストは、VPC 経由でトークン交換を使用して OAuth 接続を作成する場合など、特定の認証タイプでは時間がかかる場合があります。

## エラー

- `AlreadyExistsException`
- `InvalidInputException`
- `OperationTimeoutException`
- `ResourceNumberLimitExceededException`
- `GlueEncryptionException`

## DeleteConnection アクション (Python: `delete_connection`)

データカタログから接続を削除します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

接続が存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- `ConnectionName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除する接続の名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- `EntityNotFoundException`
- `OperationTimeoutException`

## GetConnection アクション (Python: get\_connection)

データカタログから接続定義を取得します。

### リクエスト

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

接続が存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- **Name** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

取得する接続定義の名前。

- **HidePassword** – ブール。

パスワードを返さずに接続メタデータを取得できます。例えば、AWS Glue コンソールではこのフラグを使用して接続を取得し、パスワードは表示しません。呼び出し元に AWS KMS キーを使用してパスワードを復号するアクセス許可がない可能性があるが、それ以外の接続プロパティへのアクセス許可がある場合は、このパラメータを設定します。

### レスポンス

- **Connection** – [Connection](#) オブジェクト。

リクエストされた接続定義。

### エラー

- `EntityNotFoundException`
- `OperationTimeoutException`
- `InvalidInputException`
- `GlueEncryptionException`

## GetConnections アクション (Python: get\_connections)

データカタログから接続定義のリストを取得します。

## リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

接続が存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- `Filter` – [GetConnectionsFilter](#) オブジェクト。

返される接続をコントロールするフィルター。

- `HidePassword` – ブール。

パスワードを返さずに接続メタデータを取得できます。例えば、AWS Glue コンソールではこのフラグを使用して接続を取得し、パスワードは表示しません。呼び出し元に AWS KMS キーを使用してパスワードを復号するアクセス許可がない可能性があるが、それ以外の接続プロパティへのアクセス許可がある場合は、このパラメータを設定します。

- `NextToken` – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

- `MaxResults` – 1~1000 の数値 (整数)。

1 回の応答で返す接続の最大数。

## レスポンス

- `ConnectionList` – [Connection](#) オブジェクトの配列。

リクエストされた接続定義のリスト。

- `NextToken` – UTF-8 文字列。

返された接続のリストにフィルタリングされた接続の最後のものが含まれていない場合は、継続トークン。

## エラー

- `EntityNotFoundException`
- `OperationTimeoutException`
- `InvalidInputException`

- `GlueEncryptionException`

## UpdateConnection アクション (Python: `update_connection`)

データカタログで接続定義を更新します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

接続が存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

更新する接続定義の名前。

- `ConnectionInput` – 必須: [ConnectionInput](#) オブジェクト。

該当する接続を再定義する `ConnectionInput` オブジェクト。

### レスポンス

- 応答パラメータはありません。

### エラー

- `InvalidInputException`
- `EntityNotFoundException`
- `OperationTimeoutException`
- `InvalidInputException`
- `GlueEncryptionException`

## BatchDeleteConnection アクション (Python: `batch_delete_connection`)

データカタログから接続定義のリストを削除します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

接続が存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- `ConnectionNameList` – 必須: UTF-8 文字列の配列。文字列 25 個以下。

削除する接続の名前のリスト。

## レスポンス

- `Succeeded` – UTF-8 文字列の配列。

正常に削除された接続定義の名前のリスト。

- `Errors` – キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。

各値は [ErrorDetail](#) オブジェクトです。

エラーの詳細への正常に削除されなかった接続の名前のマップ。

## エラー

- `InternalServiceException`
- `OperationTimeoutException`

## 認証の設定

- [AuthenticationConfiguration](#) 構造
- [AuthenticationConfigurationInput](#) 構造
- [OAuth2Properties](#) 構造
- [OAuth2PropertiesInput](#) 構造
- [OAuth2ClientApplication](#) 構造
- [AuthorizationCodeProperties](#) 構造

## AuthenticationConfiguration 構造

認証の設定を含む構造。

## フィールド

- `AuthenticationType` – UTF-8 文字列 (有効な値: BASIC | OAUTH2 | CUSTOM)。

認証の設定を含む構造。

- `SecretArn` – UTF-8 文字列、「[Custom string pattern #11](#)」に一致。

認証情報を保存するシークレットマネージャー ARN。

- `OAuth2Properties` – [OAuth2Properties](#) オブジェクト。

OAuth2 認証のプロパティ。

## AuthenticationConfigurationInput 構造

CreateConnection リクエストの認証の設定を含む構造。

### フィールド

- `AuthenticationType` – UTF-8 文字列 (有効な値: BASIC | OAUTH2 | CUSTOM)。

CreateConnection リクエストの認証の設定を含む構造。

- `SecretArn` – UTF-8 文字列、「[Custom string pattern #11](#)」に一致。

CreateConnection リクエストに認証情報を保存するシークレットマネージャー ARN。

- `OAuth2Properties` – [OAuth2PropertiesInput](#) オブジェクト。

CreateConnection リクエストの OAuth2 認証のプロパティ。

## OAuth2Properties 構造

OAuth2 認証のプロパティを含む構造。

### フィールド

- `OAuth2GrantType` – UTF-8 文字列 (有効な値: AUTHORIZATION\_CODE | CLIENT\_CREDENTIALS | JWT\_BEARER)。

OAuth2 付与タイプ。例えば、AUTHORIZATION\_CODE、JWT\_BEARER、または CLIENT\_CREDENTIALS などです。

- `OAuth2ClientApplication` – [OAuth2ClientApplication](#) オブジェクト。  
クライアントアプリケーションタイプ。たとえば、`AWS_MANAGED` または `USER_MANAGED`。
- `TokenUrl` - UTF-8 文字列、256 バイト長以下、[Custom string pattern #12](#) に一致。  
認証コードをアクセストークンと交換するための、プロバイダーの認証サーバーの URL。
- `TokenUrlParametersMap` – キーバリューペアのマップ配列。  
各キーは UTF-8 文字列で、1~128 バイト長です。  
各値は UTF-8 文字列で、1~512 バイト長です。  
トークン GET リクエストに追加されるパラメータのマップ。

## OAuth2PropertiesInput 構造

`CreateConnection` リクエストの OAuth2 のプロパティを含む構造。

### フィールド

- `OAuth2GrantType` – UTF-8 文字列 (有効な値: `AUTHORIZATION_CODE` | `CLIENT_CREDENTIALS` | `JWT_BEARER`)。  
`CreateConnection` リクエストの OAuth2 付与タイプ。例えば、`AUTHORIZATION_CODE`、`JWT_BEARER`、または `CLIENT_CREDENTIALS` などです。
- `OAuth2ClientApplication` – [OAuth2ClientApplication](#) オブジェクト。  
`CreateConnection` リクエストのクライアントアプリケーションタイプ。例えば、`AWS_MANAGED`、`USER_MANAGED` などです。
- `TokenUrl` - UTF-8 文字列、256 バイト長以下、[Custom string pattern #12](#) に一致。  
認証コードをアクセストークンと交換するための、プロバイダーの認証サーバーの URL。
- `TokenUrlParametersMap` – キーバリューペアのマップ配列。  
各キーは UTF-8 文字列で、1~128 バイト長です。  
各値は UTF-8 文字列で、1~512 バイト長です。  
トークン GET リクエストに追加されるパラメータのマップ。
- `AuthorizationCodeProperties` – [AuthorizationCodeProperties](#) オブジェクト。

OAuth2 AUTHORIZATION\_CODE 付与タイプに必要なプロパティのセット。

## OAuth2ClientApplication 構造

接続に使用される OAuth2 クライアントアプリ。

フィールド

- `UserManagedClientApplicationClientId` - UTF-8 文字列。2,048 バイト長以下。 [Custom string pattern #13](#) に一致。

`ClientAppType` が `USER_MANAGED` の場合のクライアントアプリケーションの `clientId`。

- `AWSManagedClientApplicationReference` - UTF-8 文字列。2,048 バイト長以下。 [Custom string pattern #13](#) に一致。

AWS によって管理されている SaaS 側のクライアントアプリケーションへの参照。

## AuthorizationCodeProperties 構造

OAuth2 AUTHORIZATION\_CODE 付与タイプのワークフローに必要なプロパティのセット。

フィールド

- `AuthorizationCode` - UTF-8 文字列、1~4096 バイト長、 [Custom string pattern #13](#) に一致。

AUTHORIZATION\_CODE で付与されるワークフローの 3 番目のログで使用される認証コード。これは、アクセストークンと交換されると無効になるワンタイムコードであるため、この値をリクエストパラメータとして指定できます。

- `RedirectUri` - UTF-8 文字列、512 バイト長以下、 [Custom string pattern #14](#) に一致。

認証コードの発行時にユーザーが認証サーバーによってリダイレクトされるリダイレクト URI。この URI は、認証コードがアクセストークンと交換されるときに使用されます。

## ユーザー定義関数 API

ユーザー定義関数 API では、関数の処理に使用する AWS Glue データ型およびオペレーションについて説明します。

## データ型

- [UserDefinedFunction 構造](#)
- [UserDefinedFunctionInput 構造](#)

### UserDefinedFunction 構造

Hive ユーザー定義関数 (UDF) 定義と同等のものを表します。

#### フィールド

- `FunctionName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

#### 関数の名前

- `DatabaseName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数を含むカタログデータベースの名前。

- `ClassName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数コードを含む Java クラス。

- `OwnerName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数の所有者。

- `OwnerType` – UTF-8 文字列 (有効な値: USER | ROLE | GROUP)。

所有者のタイプ。

- `CreateTime` – タイムスタンプ。

関数の作成時刻。

- `ResourceUri` – [ResourceUri](#) オブジェクトの配列。構造 1000 個以下。

関数のリソース URI。

- `CatalogId` – カatalog ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数が存在する Data Catalog の ID。

## UserDefinedFunctionInput 構造

ユーザー定義関数の作成または更新に使用される構造。

### フィールド

- `FunctionName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

#### 関数の名前

- `ClassName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数コードを含む Java クラス。

- `OwnerName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数の所有者。

- `OwnerType` – UTF-8 文字列 (有効な値: USER | ROLE | GROUP)。

所有者のタイプ。

- `ResourceUris` – [ResourceUri](#) オブジェクトの配列。構造 1000 個以下。

関数のリソース URI。

### 操作

- [CreateUserDefinedFunction アクション \(Python: create\\_user\\_defined\\_function\)](#)
- [UpdateUserDefinedFunction アクション \(Python: update\\_user\\_defined\\_function\)](#)
- [DeleteUserDefinedFunction アクション \(Python: delete\\_user\\_defined\\_function\)](#)
- [GetUserDefinedFunction アクション \(Python: get\\_user\\_defined\\_function\)](#)
- [GetUserDefinedFunctions アクション \(Python: get\\_user\\_defined\\_functions\)](#)

## CreateUserDefinedFunction アクション (Python: create\_user\_defined\_function)

データカタログで新しい関数定義を作成します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数を作成するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数を作成するカタログデータベースの名前。

- FunctionInput – 必須: [UserDefinedFunctionInput](#) オブジェクト。

データカタログで作成する関数を定義する FunctionInput オブジェクト。

## レスポンス

- 応答パラメータはありません。

## エラー

- AlreadyExistsException
- InvalidInputException
- InternalServiceException
- EntityNotFoundException
- OperationTimeoutException
- ResourceNumberLimitExceededException
- GlueEncryptionException

## UpdateUserDefinedFunction アクション (Python: update\_user\_defined\_function)

データカタログで既存の関数定義を更新します。

### リクエスト

- CatalogId – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

更新する関数が存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

更新する関数が存在するカタログデータベースの名前。

- `FunctionName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数の名前

- `FunctionInput` – 必須: [UserDefinedFunctionInput](#) オブジェクト。

データカタログで関数を再定義する `FunctionInput` オブジェクト。

## レスポンス

- 応答パラメータはありません。

## エラー

- `EntityNotFoundException`
- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`
- `GlueEncryptionException`

## DeleteUserDefinedFunction アクション (Python: `delete_user_defined_function`)

データカタログから既存の関数定義を削除します。

## リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除する関数が存在するデータカタログの ID。提供されない場合は、AWS アカウント ID がデフォルトで使用されます。

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除する関数が存在するカタログデータベースの名前。

- `FunctionName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除する関数定義の名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## GetUserDefinedFunction アクション (Python: `get_user_defined_function`)

データカタログから指定された関数定義を取得します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

取得する関数が存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除する関数が存在するカタログデータベースの名前。

- `FunctionName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数の名前

## レスポンス

- `UserDefinedFunction` – [UserDefinedFunction](#) オブジェクト。

リクエストされた関数定義。

## エラー

- EntityNotFoundException
- InvalidInputException

- `InternalServiceException`
- `OperationTimeoutException`
- `GlueEncryptionException`

## GetUserDefinedFunctions アクション (Python: `get_user_defined_functions`)

データカタログから複数の関数定義を取得します。

### リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

取得する関数が存在するデータカタログの ID。設定しない場合は、AWS アカウント ID がデフォルトで使用されます。

- `DatabaseName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

関数が存在するカタログデータベースの名前。何も指定されていない場合は、カタログ全体のすべてのデータベースからの関数が返されます。

- `Pattern` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

返される関数定義をフィルタリングするオプションの `function-name` パターン文字列。

- `NextToken` – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

- `MaxResults` – 数値 (integer)。1~100。

1 回の応答で返す関数の最大数。

### 応答

- `UserDefinedFunctions` – [UserDefinedFunction](#) オブジェクトの配列。

リクエストされた関数定義のリスト。

- `NextToken` – UTF-8 文字列。

戻された関数のリストに最後のリクエストされた関数が含まれていない場合は、継続トークン。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException
- GlueEncryptionException

## Athena カタログを AWS Glue にインポートする

Migration API では、Athena データカタログの AWS Glue への移行と関係のある AWS Glue データ型とオペレーションについて説明します。

### データ型

- [CatalogImportStatus](#) 構造

### CatalogImportStatus 構造

移行ステータス情報を含む構造。

#### フィールド

- ImportCompleted – ブール。

移行が完了した場合は True、それ以外の場合は False です。

- ImportTime – タイムスタンプ。

移行を開始した時刻。

- ImportedBy – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

移行を開始したユーザーの名前。

### 操作

- [ImportCatalogToGlue アクション](#) (Python: `import_catalog_to_glue`)
- [GetCatalogImportStatus アクション](#) (Python: `get_catalog_import_status`)

## ImportCatalogToGlue アクション (Python: import\_catalog\_to\_glue)

既存の Amazon Athena Data Catalog を AWS Glue にインポートする

リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
インポートするカタログの ID。現在、これは AWS アカウント ID である必要があります。

レスポンス

- 応答パラメータはありません。

エラー

- `InternalServiceException`
- `OperationTimeoutException`

## GetCatalogImportStatus アクション (Python: get\_catalog\_import\_status)

移行操作のステータスを取得します。

リクエスト

- `CatalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
移行するカタログの ID。現在、これは AWS アカウント ID である必要があります。

レスポンス

- `ImportStatus` – [CatalogImportStatus](#) オブジェクト。  
指定したカタログ移行のステータス。

エラー

- `InternalServiceException`
- `OperationTimeoutException`

# テーブル 옵ティマイザー API

テーブル 옵ティマイ저 API は、圧縮を有効にして読み取りパフォーマンスを向上させるための AWS Glue API について説明します。

## データ型

- [TableOptimizer 構造](#)
- [TableOptimizerConfiguration 構造](#)
- [TableOptimizerRun 構造](#)
- [RunMetrics 構造](#)
- [BatchGetTableOptimizerEntry 構造](#)
- [BatchTableOptimizer 構造](#)
- [BatchGetTableOptimizerError 構造](#)

## TableOptimizer 構造

テーブルに関連付けられた 옵ティマイ저に関する詳細が含まれます。

### フィールド

- `type` – UTF-8 文字列 (有効な値: `compaction="COMPACTION"`)。  
テーブル 옵ティマイ저のタイプ。現在、有効な値は `compaction` のみです。
- `configuration` – [TableOptimizerConfiguration](#) オブジェクト。  
テーブル 옵ティマイ저を作成または更新する際に指定された `TableOptimizerConfiguration` オブジェクト。
- `lastRun` – [TableOptimizerRun](#) オブジェクト。  
テーブル 옵ティマイ저の前の実行を表す `TableOptimizerRun` オブジェクト。

## TableOptimizerConfiguration 構造

テーブル 옵ティマイ저の設定に関する詳細が含まれます。テーブル 옵ティマイ저を作成または更新する際に、この設定を渡します。

## フィールド

- `roleArn` - UTF-8 文字列。1 ~ 512 バイト長。 [Single-line string pattern](#) に一致。

呼び出し元によって渡されるロール。このロールは、呼び出し元に代わってオプティマイザーに関連付けられたリソースを更新するための許可をサービスに付与します。

- `enabled` - ブール。

テーブル最適化が有効かどうか。

## TableOptimizerRun 構造

テーブルオプティマイザーの実行の詳細が含まれます。

### フィールド

- `eventType` - UTF-8 文字列 (有効な値: `starting="STARTING"` | `completed="COMPLETED"` | `failed="FAILED"` | `in_progress="IN_PROGRESS"`)。

テーブルオプティマイザーの実行のステータスを表すイベントタイプ。

- `startTimestamp` - タイムスタンプ。

Lake Formation 内で圧縮ジョブが開始された時点のエポックタイムスタンプを表します。

- `endTimestamp` - タイムスタンプ。

圧縮ジョブが終了した時点のエポックタイムスタンプを表します。

- `metrics` - [RunMetrics](#) オブジェクト。

オプティマイザーの実行のメトリクスを含む `RunMetrics` オブジェクト。

- `error` - UTF-8 文字列。

オプティマイザーの実行中に発生したエラー。

## RunMetrics 構造

オプティマイザーの実行のメトリクス。

## フィールド

- `NumberOfBytesCompacted` – UTF-8 文字列。  
圧縮ジョブの実行によって削除されたバイト数。
- `NumberOfFilesCompacted` – UTF-8 文字列。  
圧縮ジョブの実行によって削除されたファイル数。
- `NumberOfDpus` – UTF-8 文字列。  
ジョブによって消費された DPU 時間数。
- `JobDurationInHour` – UTF-8 文字列。  
ジョブの実行時間 (単位: 時間)。

## BatchGetTableOptimizerEntry 構造

`BatchGetTableOptimizer` オペレーションで取得するテーブル 옵ティマイザーを表します。

### フィールド

- `catalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルのカタログ ID。
- `databaseName` – UTF-8 文字列、少なくとも 1 バイト長。  
テーブルが存在するカタログのデータベースの名前。
- `tableName` – UTF-8 文字列、少なくとも 1 バイト長。  
テーブルの名前。
- `type` – UTF-8 文字列 (有効な値: `compaction="COMPACTION"`)。  
テーブル 옵ティマイザーのタイプ。

## BatchTableOptimizer 構造

`BatchGetTableOptimizer` オペレーションによって返されたいずれかのテーブル 옵ティマイザーの詳細が含まれます。

## フィールド

- `catalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルのカタログ ID。
- `databaseName` – UTF-8 文字列、少なくとも 1 バイト長。  
テーブルが存在するカタログのデータベースの名前。
- `tableName` – UTF-8 文字列、少なくとも 1 バイト長。  
テーブルの名前。
- `tableOptimizer` – [TableOptimizer](#) オブジェクト。  
テーブル 옵ティマイザーの設定と前回の実行に関する詳細を含む `TableOptimizer` オブジェクト。

## BatchGetTableOptimizerError 構造

`BatchGetTableOptimizer` オペレーションによって返されたエラーリスト内のいずれかのエラーに関する詳細が含まれます。

## フィールド

- `error` – [ErrorDetail](#) オブジェクト。  
エラーに関するコードとメッセージの詳細を含む `ErrorDetail` オブジェクト。
- `catalogId` – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルのカタログ ID。
- `databaseName` – UTF-8 文字列、少なくとも 1 バイト長。  
テーブルが存在するカタログのデータベースの名前。
- `tableName` – UTF-8 文字列、少なくとも 1 バイト長。  
テーブルの名前。
- `type` – UTF-8 文字列 (有効な値: `compaction="COMPACTION"`)。  
テーブル 옵ティマイザーのタイプ。

## 操作

- [GetTableOptimizer アクション \(Python: get\\_table\\_optimizer\)](#)
- [BatchGetTableOptimizer アクション \(Python: batch\\_get\\_table\\_optimizer\)](#)
- [ListTableOptimizerRuns アクション \(Python: list\\_table\\_optimizer\\_runs\)](#)
- [CreateTableOptimizer アクション \(Python: create\\_table\\_optimizer\)](#)
- [DeleteTableOptimizer アクション \(Python: delete\\_table\\_optimizer\)](#)
- [UpdateTableOptimizer アクション \(Python: update\\_table\\_optimizer\)](#)

### GetTableOptimizer アクション (Python: get\_table\_optimizer)

指定されたテーブルに関連付けられているすべてのオプティマイザーの設定を返します。

#### リクエスト

- **CatalogId** – 必須: カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルのカタログ ID。
- **DatabaseName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルが存在するカタログのデータベースの名前。
- **TableName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルの名前。
- **Type** – 必須: UTF-8 文字列 (有効な値: compaction="COMPACTION")。  
テーブルオプティマイザーのタイプ。

#### レスポンス

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルのカタログ ID。
- **DatabaseName** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルが存在するカタログのデータベースの名前。

- `TableName` – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

テーブルの名前。

- `TableOptimizer` – [TableOptimizer](#) オブジェクト。

指定されたテーブルに関連付けられたオプティマイザー。

## エラー

- `EntityNotFoundException`
- `InvalidInputException`
- `AccessDeniedException`
- `InternalServiceException`

## BatchGetTableOptimizer アクション (Python: `batch_get_table_optimizer`)

指定されたテーブルオプティマイザーの設定を返します。

### リクエスト

- `Entries` – 必須: [BatchGetTableOptimizerEntry](#) オブジェクトの配列。

取得するテーブルオプティマイザーを指定する `BatchGetTableOptimizerEntry` オブジェクトのリスト。

### 応答

- `TableOptimizers` – [BatchTableOptimizer](#) オブジェクトの配列。

`BatchTableOptimizer` オブジェクトのリスト。

- `Failures` – [BatchGetTableOptimizerError](#) オブジェクトの配列。

オペレーションで発生したエラーのリスト。

## エラー

- `InternalServiceException`

## ListTableOptimizerRuns アクション (Python: list\_table\_optimizer\_runs)

特定のテーブルについての以前のオプティマイザーの実行の履歴をリストします。

### リクエスト

- **CatalogId** – 必須: カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルのカタログ ID。
- **DatabaseName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルが存在するカタログのデータベースの名前。
- **TableName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルの名前。
- **Type** – 必須: UTF-8 文字列 (有効な値: compaction="COMPACTION")。  
テーブルオプティマイザーのタイプ。現在、有効な値は compaction のみです。
- **MaxResults** – 数値 (整数)。  
各呼び出しで返すオプティマイザーの実行の最大数。
- **NextToken** – UTF-8 文字列。  
継続トークン (これが継続呼び出しの場合)。

### 応答

- **CatalogId** – カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルのカタログ ID。
- **DatabaseName** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルが存在するカタログのデータベースの名前。
- **TableName** – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
テーブルの名前。
- **NextToken** – UTF-8 文字列。

返されたオブティマイザーの実行のリストをページ分割するための継続トークン。リストの現在のセグメントが最後のセグメントではない場合に返されます。

- TableOptimizerRuns – [TableOptimizerRun](#) オブジェクトの配列。

テーブルに関連付けられているオブティマイザーの実行のリスト。

## エラー

- EntityNotFoundException
- AccessDeniedException
- InvalidInputException
- InternalServiceException

## CreateTableOptimizer アクション (Python: create\_table\_optimizer)

特定の関数のために新しいテーブルオブティマイザーを作成します。現在サポートされているオブティマイザーの種類は compaction のみです。

### リクエスト

- CatalogId – 必須: カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルのカタログ ID。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログのデータベースの名前。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの名前。

- Type – 必須: UTF-8 文字列 (有効な値: compaction="COMPACTION")。

テーブルオブティマイザーのタイプ。現在、有効な値は compaction のみです。

- TableOptimizerConfiguration – 必須: [TableOptimizerConfiguration](#) オブジェクト。

テーブルオブティマイザーの設定を表す TableOptimizerConfiguration オブジェクト。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- AccessDeniedException
- AlreadyExistsException
- InternalServiceException

## DeleteTableOptimizer アクション (Python: delete\_table\_optimizer)

オプティマイザーと、テーブルに関連付けられているすべてのメタデータを削除します。最適化はテーブルに対して実行されなくなります。

### リクエスト

- CatalogId – 必須: カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルのカタログ ID。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するカタログのデータベースの名前。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの名前。

- Type – 必須: UTF-8 文字列 (有効な値: compaction="COMPACTION")。

テーブルオプティマイザーのタイプ。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- AccessDeniedException
- InternalServiceException

## UpdateTableOptimizer アクション (Python: update\_table\_optimizer)

既存のテーブルオプティマイザーの設定を更新します。

### リクエスト

- CatalogId – 必須: カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルのカタログ ID。
- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルが存在するカタログのデータベースの名前。
- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルの名前。
- Type – 必須: UTF-8 文字列 (有効な値: compaction="COMPACTION")。  
テーブルオプティマイザーのタイプ。現在、有効な値は compaction のみです。
- TableOptimizerConfiguration – 必須: [TableOptimizerConfiguration](#) オブジェクト。  
テーブルオプティマイザーの設定を表す TableOptimizerConfiguration オブジェクト。

### レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException

- `AccessDeniedException`
- `InternalServiceException`

## クローラーおよび分類子 API

クローラーおよび分類子 API では、AWS Glue のクローラーおよび分類子のデータ型について説明します。これには、クローラーまたは分類子を作成、削除、更新、および一覧表示するための API が含まれます。

### トピック

- [分類子 API](#)
- [クローラー API](#)
- [列統計 API](#)
- [クローラースケジューラ API](#)

## 分類子 API

分類子 API では、AWS Glue 分類子のデータ型について説明します。これには、分類子を作成、削除、更新、および一覧表示するための API が含まれます。

### データ型

- [分類子の構造](#)
- [GrokClassifier の構造](#)
- [XMLClassifier の構造](#)
- [JsonClassifier の構造](#)
- [CsvClassifier の構造](#)
- [CreateGrokClassifierRequest の構造](#)
- [UpdateGrokClassifierRequest の構造](#)
- [CreateXMLClassifierRequest の構造](#)
- [UpdateXMLClassifierRequest の構造](#)
- [CreateJsonClassifierRequest の構造](#)
- [UpdateJsonClassifierRequest の構造](#)

- [CreateCsvClassifierRequest の構造](#)
- [UpdateCsvClassifierRequest の構造](#)

## 分類子の構造

分類子は、クローラタスクでトリガーされます。分類子は、指定されたファイルが処理可能な形式であるかどうかをチェックします。処理可能な場合、分類子は、そのデータ形式と一致する StructType オブジェクトの形式でスキーマを作成します。

AWS Glue に用意されている標準の分類子を使用することも、独自の分類子を作成して、データソースの最適な分類を行い、これらに使用する適切なスキーマを指定することもできます。分類子として、Classifier オブジェクトの各フィールドに指定されている grok 分類子、XML 分類子、JSON 分類子、または CSV 分類子のいずれかを使用できます。

### フィールド

- GrokClassifier – [GrokClassifier](#) オブジェクト。  
grok を使用する分類子。
- XMLClassifier – [XMLClassifier](#) オブジェクト。  
XML コンテンツの分類子。
- JsonClassifier – [JsonClassifier](#) オブジェクト。  
JSON コンテンツの分類子。
- CsvClassifier – [CsvClassifier](#) オブジェクト。  
カンマ区切り値 (CSV) の分類子。

## GrokClassifier の構造

grok パターンを使用する分類子。

### フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
分類子名。
- Classification – 必須: UTF-8 文字列。

Twitter、JSON、Omniure ログなど、分類子が一致するデータ形式の識別子。

- **CreationTime** – タイムスタンプ。

分類子が登録された時刻。

- **LastUpdated** – タイムスタンプ。

分類子が最後に更新された時刻。

- **Version** – 数値 (long 型)。

分類子のバージョン。

- **GrokPattern** – 必須: UTF-8 文字列。1 ~ 2048 バイト長。 [A Logstash Grok string pattern](#) に一致。

分類子によってデータストアに適用される grok パターン。詳細については、「[カスタム分類子の書き込み](#)」の組み込みパターンを参照してください。

- **CustomPatterns** - UTF-8 文字列。16,000 バイト長以下。 [URI address multi-line string pattern](#) に一致。

この分類子によって定義されたオプションのカスタム Grok パターン。詳細については、「[カスタム分類子の書き込み](#)」のカスタムパターンを参照してください。

## XMLClassifier の構造

XML コンテンツの分類子。

フィールド

- **Name** – 必須: UTF-8 文字列、1 ~ 255 バイト長、 [Single-line string pattern](#) に一致。

分類子名。

- **Classification** – 必須: UTF-8 文字列。

分類子が一致するデータ形式の識別子。

- **CreationTime** – タイムスタンプ。

分類子が登録された時刻。

- **LastUpdated** – タイムスタンプ。

分類子が最後に更新された時刻。

- Version – 数値 (long 型)。

分類子のバージョン。

- RowTag – UTF-8 文字列。

解析中の XML ドキュメントの各レコードを含む要素を指定する XML タグ。これは、自己終了要素 (`</>` で終了) を識別することはできません。属性のみを含む空の行要素は、終了タグで終わる場合は解析できます (例: `<row item_a="A" item_b="B"></row>` は解析できますが、`<row item_a="A" item_b="B" />` は解析できません)。

## JsonClassifier の構造

JSON コンテンツの分類子。

フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

分類子名。

- CreationTime – タイムスタンプ。

分類子が登録された時刻。

- LastUpdated – タイムスタンプ。

分類子が最後に更新された時刻。

- Version – 数値 (long 型)。

分類子のバージョン。

- JsonPath – 必須: UTF-8 文字列。

分類子が分類するための JSON データを定義する JsonPath 文字列。AWS Glue は、「[Writing JsonPath Custom Classifiers](#)」に記載されているように、JsonPath のサブセットをサポートしています。

## CsvClassifier の構造

カスタム CSV コンテンツの分類子。

## フィールド

- **Name** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
分類子名。
- **CreationTime** – タイムスタンプ。  
分類子が登録された時刻。
- **LastUpdated** – タイムスタンプ。  
分類子が最後に更新された時刻。
- **Version** – 数値 (long 型)。  
分類子のバージョン。
- **Delimiter** – UTF-8 文字列。1 バイト長ちょうど。[Custom string pattern #10](#) に一致。  
行内の各列エントリを区切るものを示すカスタム記号。
- **QuoteSymbol** – UTF-8 文字列。1 バイト長ちょうど。[Custom string pattern #10](#) に一致。  
コンテンツを単一の列の値に結合するものを示すカスタム記号。列の区切り文字とは異なる必要があります。
- **ContainsHeader** – UTF-8 文字列 (有効な値: UNKNOWN | PRESENT | ABSENT)。  
CSV ファイルにヘッダーが含まれているかどうかを示します。
- **Header** – UTF-8 文字列の配列。  
列名を表す文字列のリスト。
- **DisableValueTrimming** – ブール。  
列の値のタイプを識別する前に値をトリミングしないことを指定します。デフォルト値は true です。
- **AllowSingleColumn** – ブール。  
1 つの列のみを含むファイルの処理を有効にします。
- **CustomDatatypeConfigured** – ブール。  
カスタムデータ型を設定できるようにします。
- **CustomDatatypes** – UTF-8 文字列の配列。

「BINARY」、「BOOLEAN」、「DATE」、「DECIMAL」、「DOUBLE」、「FLOAT」、「INT」、「LONG」、「SHORT」、「STRING」、「TIMESTAMP」などのカスタムデータ型のリスト。

- Serde – UTF-8 文字列 (有効な値: OpenCSVSerDe | LazySimpleSerDe | None)。

分類子で CSV を処理するため、データカタログに適用される SerDe を設定します。有効な値は、OpenCSVSerDe、LazySimpleSerDe、None です。クローラーが検出を行う場合は、None 値を指定できます。

## CreateGrokClassifierRequest の構造

作成する CreateClassifier の grok 分類子を指定します。

### フィールド

- Classification – 必須: UTF-8 文字列。

Twitter、JSON、Omniure ログ、Amazon CloudWatch Logs など、分類子が一致するデータ形式の識別子。

- Name – 必須: UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。

新しい分類子の名前。

- GrokPattern – 必須: UTF-8 文字列。1 ~ 2048 バイト長。[A Logstash Grok string pattern](#) に一致。

この分類子によって使用される grok パターン。

- CustomPatterns - UTF-8 文字列。16,000 バイト長以下。[URI address multi-line string pattern](#) に一致。

この分類子によって使用されたオプションのカスタム Grok パターン。

## UpdateGrokClassifierRequest の構造

UpdateClassifier に渡すときに更新する grok 分類子を指定します。

### フィールド

- Name – 必須: UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。

GrokClassifier の名前。

- Classification – UTF-8 文字列。

Twitter、JSON、Omnicore ログ、Amazon CloudWatch Logs など、分類子が一致するデータ形式の識別子。

- GrokPattern – UTF-8 文字列。1 ~ 2,048 バイト長。[A Logstash Grok string pattern](#) に一致。

この分類子によって使用される grok パターン。

- CustomPatterns – UTF-8 文字列。16,000 バイト長以下。[URI address multi-line string pattern](#) に一致。

この分類子によって使用されたオプションのカスタム Grok パターン。

## CreateXMLClassifierRequest の構造

作成する CreateClassifier の XML 分類子を指定します。

### フィールド

- Classification – 必須: UTF-8 文字列。

分類子が一致するデータ形式の識別子。

- Name – 必須: UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。

分類子名。

- RowTag – UTF-8 文字列。

解析中の XML ドキュメントの各レコードを含む要素を指定する XML タグ。これは、自己終了要素 (`</>` で終了) を識別することはできません。属性のみを含む空の行要素は、終了タグで終わる場合は解析できます (例: `<row item_a="A" item_b="B"></row>` は解析できますが、`<row item_a="A" item_b="B" />` は解析できません)。

## UpdateXMLClassifierRequest の構造

更新する XML 分類子を指定します。

## フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

分類子名。

- Classification – UTF-8 文字列。

分類子が一致するデータ形式の識別子。

- RowTag – UTF-8 文字列。

解析中の XML ドキュメントの各レコードを含む要素を指定する XML タグ。これは、自己終了要素 (`</>` で終了) を識別することはできません。属性のみを含む空の行要素は、終了タグで終わる場合は解析できます (例: `<row item_a="A" item_b="B"></row>` は解析できますが、`<row item_a="A" item_b="B" />` は解析できません)。

## CreateJsonClassifierRequest の構造

作成する `CreateClassifier` の JSON 分類子を指定します。

### フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

分類子名。

- JsonPath – 必須: UTF-8 文字列。

分類子が分類するための JSON データを定義する JsonPath 文字列。AWS Glue は、「[Writing JsonPath Custom Classifiers](#)」に記載されているように、JsonPath のサブセットをサポートしています。

## UpdateJsonClassifierRequest の構造

更新する JSON 分類子を指定します。

### フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

分類子名。

- JsonPath – UTF-8 文字列。

分類子が分類するための JSON データを定義する JsonPath 文字列。AWS Glue は、「[Writing JsonPath Custom Classifiers](#)」に記載されているように、JsonPath のサブセットをサポートしています。

## CreateCsvClassifierRequest の構造

作成する CreateClassifier のカスタム CSV 分類子を指定します。

### フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
分類子名。
- Delimiter – UTF-8 文字列。1 バイト長ちょうど。[Custom string pattern #10](#) に一致。  
行内の各列エントリを区切るものを示すカスタム記号。
- QuoteSymbol – UTF-8 文字列。1 バイト長ちょうど。[Custom string pattern #10](#) に一致。  
コンテンツを単一の列の値に結合するものを示すカスタム記号。列の区切り文字とは異なる必要があります。
- ContainsHeader – UTF-8 文字列 (有効な値: UNKNOWN | PRESENT | ABSENT)。  
CSV ファイルにヘッダーが含まれているかどうかを示します。
- Header – UTF-8 文字列の配列。  
列名を表す文字列のリスト。
- DisableValueTrimming – ブール。  
列の値のタイプを識別する前に値をトリミングしないことを指定します。デフォルト値は True です。
- AllowSingleColumn – ブール。  
1 つの列のみを含むファイルの処理を有効にします。
- CustomDatatypeConfigured – ブール。  
カスタムデータ型の設定を有効にします。
- CustomDatatypes – UTF-8 文字列の配列。

サポートされているカスタムデータ型のリストを作成します。

- Serde – UTF-8 文字列 (有効な値: OpenCSVSerDe | LazySimpleSerDe | None)。

分類子で CSV を処理するため、データカタログに適用される SerDe を設定します。有効な値は、OpenCSVSerDe、LazySimpleSerDe、None です。クローラーが検出を行う場合は、None 値を指定できます。

## UpdateCsvClassifierRequest の構造

更新する カスタム CSV 分類子を指定します。

### フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

分類子名。

- Delimiter – UTF-8 文字列。1 バイト長ちょうど。[Custom string pattern #10](#) に一致。

行内の各列エントリを区切るものを示すカスタム記号。

- QuoteSymbol – UTF-8 文字列。1 バイト長ちょうど。[Custom string pattern #10](#) に一致。

コンテンツを単一の列の値に結合するものを示すカスタム記号。列の区切り文字とは異なる必要があります。

- ContainsHeader – UTF-8 文字列 (有効な値: UNKNOWN | PRESENT | ABSENT)。

CSV ファイルにヘッダーが含まれているかどうかを示します。

- Header – UTF-8 文字列の配列。

列名を表す文字列のリスト。

- DisableValueTrimming – ブール。

列の値のタイプを識別する前に値をトリミングしないことを指定します。デフォルト値は True です。

- AllowSingleColumn – ブール。

1 つの列のみを含むファイルの処理を有効にします。

- CustomDatatypeConfigured – ブール。

カスタムデータ型の設定を指定します。

- CustomDatatypes – UTF-8 文字列の配列。

サポートされているカスタムデータ型のリストを指定します。

- Serde – UTF-8 文字列 (有効な値: OpenCSVSerDe | LazySimpleSerDe | None)。

分類子で CSV を処理するため、データカタログに適用される SerDe を設定します。有効な値は、OpenCSVSerDe、LazySimpleSerDe、None です。クローラーが検出を行う場合は、None 値を指定できます。

## 操作

- [CreateClassifier アクション \(Python: create\\_classifier\)](#)
- [DeleteClassifier アクション \(Python: delete\\_classifier\)](#)
- [GetClassifier アクション \(Python: get\\_classifier\)](#)
- [GetClassifiers アクション \(Python: get\\_classifiers\)](#)
- [UpdateClassifier アクション \(Python: update\\_classifier\)](#)

## CreateClassifier アクション (Python: create\_classifier)

ユーザーのアカウントに分類子を作成します。これは、どのリクエストのフィールドが存在するかに応じて、GrokClassifier、XMLClassifier、JsonClassifier、または CsvClassifier である場合があります。

### リクエスト

- GrokClassifier – [CreateGrokClassifierRequest](#) オブジェクト。  
作成する分類子を指定する GrokClassifier オブジェクト。
- XMLClassifier – [CreateXMLClassifierRequest](#) オブジェクト。  
作成する分類子を指定する XMLClassifier オブジェクト。
- JsonClassifier – [CreateJsonClassifierRequest](#) オブジェクト。  
作成する分類子を指定する JsonClassifier オブジェクト。
- CsvClassifier – [CreateCsvClassifierRequest](#) オブジェクト。

作成する分類子を指定する `CsvClassifier` オブジェクト。

## レスポンス

- 応答パラメータはありません。

## エラー

- `AlreadyExistsException`
- `InvalidInputException`
- `OperationTimeoutException`

## DeleteClassifier アクション (Python: `delete_classifier`)

データカタログから分類子を削除します。

## リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

削除する分類子の名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- `EntityNotFoundException`
- `OperationTimeoutException`

## GetClassifier アクション (Python: `get_classifier`)

分類子を名前で取得します。

## リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

取得する分類子の名前。

## レスポンス

- Classifier – [分類子](#) オブジェクト。

リクエストされた分類子。

## エラー

- EntityNotFoundException
- OperationTimeoutException

## GetClassifiers アクション (Python: get\_classifiers)

データカタログの分類子オブジェクトのすべてを一覧表示します。

## リクエスト

- MaxResults – 1~1000 の数値 (整数)。

返されるリストのサイズ (オプション)。

- NextToken – UTF-8 文字列。

オプションの継続トークン。

## 応答

- Classifiers – [分類子](#) オブジェクトの配列。

分類子オブジェクトのリクエストされたリスト。

- NextToken – UTF-8 文字列。

継続トークン。

## エラー

- `OperationTimeoutException`

## UpdateClassifier アクション (Python: `update_classifier`)

既存の分類子を変更します (フィールドが存在するかどうかに応じて `GrokClassifier`、`XMLClassifier`、`JsonClassifier`、または `CsvClassifier`)。

### リクエスト

- `GrokClassifier` – [UpdateGrokClassifierRequest](#) オブジェクト。

フィールドが更新された `GrokClassifier` オブジェクト。

- `XMLClassifier` – [UpdateXMLClassifierRequest](#) オブジェクト。

フィールドが更新された `XMLClassifier` オブジェクト。

- `JsonClassifier` – [UpdateJsonClassifierRequest](#) オブジェクト。

フィールドが更新された `JsonClassifier` オブジェクト。

- `CsvClassifier` – [UpdateCsvClassifierRequest](#) オブジェクト。

フィールドが更新された `CsvClassifier` オブジェクト。

### レスポンス

- 応答パラメータはありません。

## エラー

- `InvalidInputException`
- `VersionMismatchException`
- `EntityNotFoundException`
- `OperationTimeoutException`

## クローラー API

クローラー API は、AWS Glue クローラーのデータ型と、クローラーを作成、削除、更新、および一覧表示するための API について説明します。

### データ型

- [Crawler 構造](#)
- [Schedule 構造](#)
- [CrawlerTargets 構造](#)
- [S3Target 構造](#)
- [S3DeltaCatalogTarget 構造](#)
- [S3DeltaDirectTarget 構造](#)
- [JdbcTarget 構造](#)
- [MongoDBTarget 構造](#)
- [DynamoDBTarget 構造](#)
- [DeltaTarget 構造](#)
- [IcebergTarget 構造](#)
- [HudiTarget 構造](#)
- [CatalogTarget 構造](#)
- [CrawlerMetrics 構造](#)
- [CrawlerHistory 構造](#)
- [CrawlsFilter 構造](#)
- [SchemaChangePolicy 構造](#)
- [LastCrawlInfo 構造](#)
- [RecrawlPolicy 構造](#)
- [LineageConfiguration 構造](#)
- [LakeFormationConfiguration 構造](#)

## Crawler 構造

データソースを検査し、分類子を使用してスキーマを判別しようとするクローラープログラムを指定します。成功すると、クローラーはデータソースに関するメタデータを AWS Glue Data Catalog に記録します。

### フィールド

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

クローラー名。

- Role – UTF-8 文字列。

Amazon Simple Storage Service (Amazon S3) データなどの顧客リソースへのアクセスに使用される IAM ロールの Amazon リソースネーム (ARN)。

- Targets – [CrawlerTargets](#) オブジェクト。

クロールするターゲットのコレクション。

- DatabaseName – UTF-8 文字列。

クローラーの出力が保存されている場所のデータベース名。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

クローラーの説明。

- Classifiers – UTF-8 文字列の配列。

クローラーに関連付けられたカスタム分類子を指定する UTF-8 文字列のリスト。

- RecrawlPolicy – [RecrawlPolicy](#) オブジェクト。

データセット全体を再度クロールするか、前回のクローラー実行以降に追加されたフォルダのみをクロールするかを指定するポリシー。

- SchemaChangePolicy – [SchemaChangePolicy](#) オブジェクト。

クローラーの更新と削除の動作を指定するポリシー。

- LineageConfiguration – [LineageConfiguration](#) オブジェクト。

クローラーに対してデータシステムを有効にするかどうかを指定する設定。

- State – UTF-8 文字列 (有効な値: READY | RUNNING | STOPPING)。

クローラーが実行中かどうか、あるいは実行が保留中かどうかを示します。

- TablePrefix - UTF-8 文字列。128 バイト長以下。

作成されたテーブルの名前に追加されるプレフィックス。

- Schedule - [スケジュール](#) オブジェクト。

スケジュールされたクローラーの場合、クローラーが実行されるスケジュール。

- CrawlElapsedTime - 数値 (long 型)。

クローラーが実行されている場合は、最後のクロールが開始されてから経過した合計時間が含まれます。

- CreationTime - タイムスタンプ。

クローラーが作成された時刻。

- LastUpdated - タイムスタンプ。

クローラーが最後に更新された時刻。

- LastCrawl - [LastCrawlInfo](#) オブジェクト。

最後のクロールのステータス、およびエラーが発生した場合はエラー情報。

- Version - 数値 (long 型)。

クローラーのバージョン。

- Configuration - UTF-8 文字列。

クローラーの構成情報。このバージョン付きの JSON 文字列では、クローラーの動作特性を指定できません。詳細については、「[クローラー設定オプションの設定](#)」を参照してください。

- CrawlerSecurityConfiguration - UTF-8 文字列。128 バイト長以下。

このクローラーで使用される SecurityConfiguration 構造の名前。

- LakeFormationConfiguration - [LakeFormationConfiguration](#) オブジェクト。

クローラーが IAM ロールの AWS Lake Formation 認証情報ではなく、クローラーの認証情報を使用するかどうかを指定します。

## Schedule 構造

cron ステートメントを使用してイベントをスケジュールするスケジューリングオブジェクト。

フィールド

- `ScheduleExpression` – UTF-8 文字列。

スケジュールを指定するために使用される cron 式 ([ジョブとクローラーの時間ベースのスケジュール](#)を参照してください。) たとえば、毎日 12:15 UTC に何かを実行するには、`cron(15 12 * * ? *)` を指定します。

- `State` – UTF-8 文字列 (有効な値: SCHEDULED | NOT\_SCHEDULED | TRANSITIONING)。

スケジュールの状態。

## CrawlerTargets 構造

クローリングするデータストアを指定します。

フィールド

- `S3Targets` – [S3Target](#) オブジェクトの配列。

Amazon Simple Storage Service (Amazon S3) のターゲットを指定します。

- `JdbcTargets` – [JdbcTarget](#) オブジェクトの配列。

JDBC ターゲットを指定します。

- `MongoDBTargets` – [MongoDBTarget](#) オブジェクトの配列。

Amazon DocumentDB または MongoDB のターゲットを指定します。

- `DynamoDBTargets` – [DynamoDBTarget](#) オブジェクトの配列。

Amazon DynamoDB のターゲットを指定します。

- `CatalogTargets` – [CatalogTarget](#) オブジェクトの配列。

AWS Glue Data Catalog ターゲットを指定します。

- `DeltaTargets` – [DeltaTarget](#) オブジェクトの配列。

Delta データストアのターゲットを指定します。

- IcebergTargets – [IcebergTarget](#) オブジェクトの配列。

Apache Iceberg データストアのターゲットを指定します。

- HudiTargets – [HudiTarget](#) オブジェクトの配列。

Apache Hudi データストアのターゲットを指定します。

## S3Target 構造

Amazon Simple Storage Service (Amazon S3) のデータストアを指定します。

### フィールド

- Path – UTF-8 文字列。

Simple Storage Service (Amazon S3) ターゲットへのパス。

- Exclusions – UTF-8 文字列の配列。

クローラから除外するために使用される glob パターンのリスト。詳細については、「[クローラーを使用したカタログテーブル](#)」を参照してください。

- ConnectionName – UTF-8 文字列。

ジョブまたはクローラーが Amazon Virtual Private Cloud 環境 (Amazon VPC) 内の Amazon S3 のデータにアクセスすることを可能にする接続の名前。

- SampleSize – 数値 (整数)。

データセット内のサンプルファイルをクローラするときクローラされる各リーフフォルダ内のファイル数を設定します。設定されていない場合、すべてのファイルがクローラされます。有効な値は、1 から 249 までの整数です。

- EventQueueArn – UTF-8 文字列。

有効な Amazon SQS の ARN。例えば `arn:aws:sqs:region:account:sqs` です。

- DdqEventQueueArn – UTF-8 文字列。

有効な Amazon デッドレター SQS ARN。例えば `arn:aws:sqs:region:account:deadLetterQueue` です。

## S3DeltaCatalogTarget 構造

AWS Glue データカタログ内の Delta Lake データソースに書き込むターゲットを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- PartitionKeys – UTF-8 文字列の配列。

一連のキーを使用してネイティブパーティショニングを指定します。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベーステーブルの名前。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベースの名前。

- AdditionalOptions – キーバリューペアのマップ配列。

各キーは、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

コネクタの追加接続オプションを指定します。

- SchemaChangePolicy – [CatalogSchemaChangePolicy](#) オブジェクト。

クローラの更新の動作を指定するポリシー。

## S3DeltaDirectTarget 構造

で Delta Lake データソースに書き込むターゲットを指定します Amazon S3。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- PartitionKeys – UTF-8 文字列の配列。

一連のキーを使用してネイティブパーティショニングを指定します。

- Path – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込み先の Delta Lake データソースの Amazon S3 パス。

- Compression – 必須: UTF-8 文字列 (有効な値: uncompressed="UNCOMPRESSED" | snappy="SNAPPY")。

データの圧縮方法を指定します。データに標準のファイル拡張子が付いている場合、このオプションは一般的に不要です。指定できる値は "gzip" および "bzip" です。

- Format – 必須: UTF-8 文字列 (有効な値: json="JSON" | csv="CSV" | avro="AVRO" | orc="ORC" | parquet="PARQUET" | hudi="HUDI" | delta="DELTA")。

ターゲットのデータ出力形式を指定します。

- AdditionalOptions – キーバリューペアのマップ配列。

各キーは、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

コネクタの追加接続オプションを指定します。

- SchemaChangePolicy – [DirectSchemaChangePolicy](#) オブジェクト。

クローラの更新の動作を指定するポリシー。

## JdbcTarget 構造

クローラする JDBC データストアを指定します。

フィールド

- ConnectionName – UTF-8 文字列。

JDBC ターゲットに接続するために使用する接続名。

- Path – UTF-8 文字列。

JDBC ターゲットのパス。

- Exclusions – UTF-8 文字列の配列。

クローラから除外するために使用される glob パターンのリスト。詳細については、「[クローラーを使用したカタログテーブル](#)」を参照してください。

- EnableAdditionalMetadata – UTF-8 文字列の配列。

RAWTYPES または COMMENTS の値を指定して、テーブルのレスポンスでその他のメタデータを有効にできます。RAWTYPES はネイティブレベルのデータ型、COMMENTS はデータベース内の列またはテーブルに関連するコメントを提供します。

その他のメタデータが必要ない場合は、フィールドを空白のままにしてください。

## MongoDBTarget 構造

クローラする Amazon DocumentDB または MongoDB データストアを指定します。

フィールド

- ConnectionName – UTF-8 文字列。

Amazon DocumentDB または MongoDB ターゲットに接続するために使用する接続名。

- Path – UTF-8 文字列。

Amazon DocumentDB または MongoDB ターゲット (データベース/コレクション) のパス。

- ScanAll – ブール。

すべてのレコードをスキャンするか、テーブルから行をサンプリングするかを示します。テーブルが高スループットテーブルではない場合、すべてのレコードのスキャンには時間がかかることがあります。

true 値はすべてのレコードをスキャンすることを意味し、false 値はレコードをサンプリングすることを意味します。値を指定しないと、true 値にデフォルト設定されます。

## DynamoDBTarget 構造

クローलする Amazon DynamoDB テーブルを指定します。

フィールド

- Path – UTF-8 文字列。

クローलする DynamoDB テーブルの名前。

- scanAll – ブール。

すべてのレコードをスキャンするか、テーブルから行をサンプリングするかを示します。テーブルが高スループットテーブルではない場合、すべてのレコードのスキャンには時間がかかることがあります。

true 値はすべてのレコードをスキャンすることを意味し、false 値はレコードをサンプリングすることを意味します。値を指定しないと、true 値にデフォルト設定されます。

- scanRate – 数値 (double)。

AWS Glue クローラーが使用する設定済み読み込みキャパシティユニットの割合。読み取りキャパシティユニットは、DynamoDB で定義されている用語で、テーブルに対して実行できる読み取り回数/秒のレート制限として機能する数値です。

有効な値は NULL または 0.1~1.5 の値です。NULL 値は、ユーザーが値を指定しない場合に使用され、設定済み読み取りキャパシティユニットでは 0.5 (プロビジョニングされたテーブルの場合)、は最大の設定済み読み取りキャパシティユニットの場合は 0.25 (オンデマンドモードを使用するテーブルの場合) にデフォルト設定されます。

## DeltaTarget 構造

1 つ以上の Delta テーブルをクローलする Delta データストアを指定します。

フィールド

- DeltaTables – UTF-8 文字列の配列。

Delta テーブルへの Simple Storage Service (Amazon S3) パスのリスト。

- ConnectionName – UTF-8 文字列。

Delta テーブルターゲットに接続するために使用する接続の名前。

- WriteManifest – ブール。

マニフェストファイルを Delta テーブルパスに書き込むかどうかを指定します。

- CreateNativeDeltaTable – ブール。

クローラーがネイティブテーブルを作成するかどうかを指定します。これにより、Delta トランザクションログの直接クエリをサポートするクエリエンジンとの統合が可能になります。

## IcebergTarget 構造

Amazon S3内の Iceberg テーブルが格納されている Apache Iceberg データソースを指定します。

### フィールド

- Paths – UTF-8 文字列の配列。

Iceberg メタデータフォルダを 1 つ以上の Amazon S3 パス `s3://bucket/prefix`。

- ConnectionName – UTF-8 文字列。

Iceberg ターゲットに接続するために使用する接続の名前。

- Exclusions – UTF-8 文字列の配列。

クロールから除外するために使用される glob パターンのリスト。詳細については、「[クローラーを使用したカタログテーブル](#)」を参照してください。

- MaximumTraversalDepth – 数値 (整数)。

クローラーが Amazon S3 パス内の Iceberg メタデータフォルダを検出するために通過できる Amazon S3 パスの最大深度。クローラーの実行時間を制限するために使用されます。

## HudiTarget 構造

Apache Hudi データソースを指定します。

### フィールド

- Paths – UTF-8 文字列の配列。

Hudi Amazon S3 の場所文字列の配列。各文字列は、Hudi テーブルのメタデータファイルが存在するルートフォルダを示します。Hudi フォルダは、ルートフォルダの子フォルダ内に存在する場合があります。

クローラーは、Hudi フォルダのパス以下にあるすべてのフォルダをスキャンします。

- `ConnectionName` – UTF-8 文字列。

Hudi ターゲットに接続するために使用する接続の名前。Hudi ファイルが VPC 認証を必要とするバケットに格納されている場合、ここで接続プロパティを設定できます。

- `Exclusions` – UTF-8 文字列の配列。

クロールから除外するために使用される glob パターンのリスト。詳細については、「[クローラーを使用したカタログテーブル](#)」を参照してください。

- `MaximumTraversalDepth` – 数値 (整数)。

クローラーが Amazon S3 パス内の Hudi メタデータフォルダを検出するために通過できる Amazon S3 パスの最大深度。クローラーの実行時間を制限するために使用されます。

## CatalogTarget 構造

AWS Glue Data Catalog ターゲットを指定します。

フィールド

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

同期するデータベースの名前。

- `Tables` – 必須: UTF-8 文字列の配列。1 個の以上の文字列。

同期するテーブルのリスト。

- `ConnectionName` – UTF-8 文字列。

NETWORK 接続タイプとペアになっている Catalog 接続タイプを使用するときにクロールのターゲットとなる Amazon S3-backed データカタログテーブルの接続の名前。

- `EventQueueArn` – UTF-8 文字列。

有効な Amazon SQS の ARN。例えば `arn:aws:sqs:region:account:sqs` です。

- `DlqEventQueueArn` – UTF-8 文字列。

有効な Amazon デッドレター SQS ARN。例えば  
`arn:aws:sqs:region:account:deadLetterQueue` です。

## CrawlerMetrics 構造

指定されたクローラーのメトリクス。

フィールド

- `CrawlerName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

クローラー名。

- `TimeLeftSeconds` – 数値 (double)。None 以下。

実行中のクローラーを完了までの予測時間。

- `StillEstimating` – ブール。

クローラーがこの実行を完了するのにどれくらいの時間がかかるかをまだ見積もっている場合は true です。

- `LastRuntimeSeconds` – 数値 (double)。None 以下。

クローラーの最新の実行にかかる時間 (秒単位)。

- `MedianRuntimeSeconds` – 数値 (double)。None 以下。

このクローラーの実行時間の中央値 (秒単位)。

- `TablesCreated` – 数値 (整数)、None 以下。

このクローラーで作成されたテーブルの数。

- `TablesUpdated` – 数値 (整数)、None 以下。

このクローラーで更新されたテーブルの数。

- `TablesDeleted` – 数値 (整数)、None 以下。

このクローラーで削除されたテーブルの数。

## CrawlerHistory 構造

クローラーの実行に関する情報が含まれます。

## フィールド

- **CrawlId** – UTF-8 文字列。  
各クロールの UUID 識別子。
- **State** – UTF-8 文字列 (有効な値: RUNNING | COMPLETED | FAILED | STOPPED)。  
クロールの状態。
- **StartTime** – タイムスタンプ。  
クロールが開始された日時。
- **EndTime** – タイムスタンプ。  
クロールが終了された日時。
- **Summary** – UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。  
特定のクロールに関する JSON 形式の実行サマリー。追加、更新、または削除されたカタログテーブルとパーティションが含まれます。
- **ErrorMessage** – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
エラーが発生した場合は、クロールに関連付けられたエラーメッセージ。
- **LogGroup** - UTF-8 文字列。1 ~ 512 バイト長。[Log group string pattern](#) に一致。  
クロールに関連付けられたロググループ。
- **LogStream** - UTF-8 文字列。1 ~ 512 バイト長。[Log-stream string pattern](#) に一致。  
クロールに関連付けられたログストリーム。
- **MessagePrefix** – UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。  
このクロールに関する CloudWatch メッセージのプレフィックス。
- **DPUHour** – 数値 (double)。None 以下。  
クロールに使用されるデータ処理単位 (DPU) の数 (時間単位)。

## CrawlsFilter 構造

指定されたクローラーのクローラー実行をフィルタリングするために使用できるフィールド、コンパレータ、および値のリスト。

## フィールド

- **FieldName** – UTF-8 文字列 (有効な値: CRAWL\_ID | STATE | START\_TIME | END\_TIME | DPU\_HOUR)。

指定されたクローラーに対するクローラーの実行をフィルタリングするために使用されるキー。各フィールド名に有効な値は次のとおりです。

- **CRAWL\_ID**: クロールの UUID 識別子を表す文字列。
- **STATE**: クロールの状態を表す文字列。
- **START\_TIME** および **END\_TIME**: epoch タイムスタンプ (ミリ秒単位)。
- **DPU\_HOUR**: クロールに使用されるデータ処理単位 (DPU) の数 (時間単位)。
- **FilterOperator** - UTF-8 文字列 (有効値: GT | GE | LT | LE | EQ | NE)。

値を操作する定義済みのコンパレータ。利用できる演算子は次のとおりです。

- **GT**: より大きい。
- **GE**: 以上。
- **LT**: より小さい。
- **LE**: 以下。
- **EQ**: 等しい。
- **NE**: 等しくない。
- **FieldValue** – UTF-8 文字列。

クロールフィールドでの比較のために提供される値。

## SchemaChangePolicy 構造

クローラーの更新と削除の動作を指定するポリシー。

### フィールド

- **UpdateBehavior** – UTF-8 文字列 (有効な値: LOG | UPDATE\_IN\_DATABASE)。

クローラーが変更されたスキーマを検出したときの更新動作。

- **DeleteBehavior** – UTF-8 文字列 (有効な値: LOG | DELETE\_FROM\_DATABASE | DEPRECATE\_IN\_DATABASE)。

クローラーが削除されたオブジェクトを検出したときの削除動作。

## LastCrawlInfo 構造

最新のクローलについてのステータスとエラー情報。

フィールド

- `Status` – UTF-8 文字列 (有効な値: SUCCEEDED | CANCELLED | FAILED)。

最後のクロールのステータス。

- `ErrorMessage` – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

エラーが発生した場合、最後のクロールに関するエラー情報。

- `LogGroup` – UTF-8 文字列。1 ~ 512 バイト長。[Log group string pattern](#) に一致。

最後のクロールのロググループ。

- `LogStream` – UTF-8 文字列。1 ~ 512 バイト長。[Log-stream string pattern](#) に一致。

最後のクロールのログストリーム。

- `MessagePrefix` – UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。

このクロールについてのメッセージのプレフィックス。

- `StartTime` – タイムスタンプ。

クロールが開始された時刻。

## RecrawlPolicy 構造

最初のクロールの完了後に Amazon S3 データソースをクロールするときに、データセット全体を再度クロールするか、前回のクローラーの実行以降に追加されたフォルダのみをクロールするかを指定します。詳細については、デベロッパーガイドの「[AWS Glueの増分クロール](#)」を参照してください。

フィールド

- `RecrawlBehavior` – UTF-8 文字列 (有効な値: CRAWL\_EVERYTHING | CRAWL\_NEW\_FOLDERS\_ONLY | CRAWL\_EVENT\_MODE)。

データセット全体を再度クローलするか、前回のクローラーの実行以降に追加されたフォルダのみをクローलするかを指定します。

CRAWL\_EVERYTHING という値は、データセット全体を再度クロールすることを指定します。

CRAWL\_NEW\_FOLDERS\_ONLY という値は、前回のクローラー実行後に追加されたフォルダのみをクロールすることを指定します。

CRAWL\_EVENT\_MODE の値は Simple Storage Service (Amazon S3) イベントによって識別される変更のみをクロールするように指定します。

## LineageConfiguration 構造

クローラーのデータシステム設定を指定します。

フィールド

- CrawlerLineageSettings – UTF-8 文字列 (有効な値: ENABLE | DISABLE)。

クローラーに対してデータシステムを有効にするかどうかを指定します。有効な値は次のとおりです。

- ENABLE: クローラーのデータシステムを有効にします。
- DISABLE : クローラーのデータシステムを無効にします。

## LakeFormationConfiguration 構造

クローラー AWS Lake Formation の構成設定を指定します。

フィールド

- UseLakeFormationCredentials – ブール。

IAM ロールの AWS Lake Formation 認証情報の代わりにクローラーの認証情報を使用するかどうかを指定します。

- AccountId - UTF-8 文字列。12 バイト長以下。

クロスアカウントクロールに必要です。ターゲットデータと同じアカウントのクロールでは、null のままにすることができます。

## 操作

- [CreateCrawler アクション \(Python: create\\_crawler\)](#)
- [DeleteCrawler アクション \(Python: delete\\_crawler\)](#)
- [GetCrawler アクション \(Python: get\\_crawler\)](#)
- [GetCrawlers アクション \(Python: get\\_crawlers\)](#)
- [GetCrawlerMetrics アクション \(Python: get\\_crawler\\_metrics\)](#)
- [UpdateCrawler アクション \(Python: update\\_crawler\)](#)
- [StartCrawler アクション \(Python: start\\_crawler\)](#)
- [StopCrawler アクション \(Python: stop\\_crawler\)](#)
- [BatchGetCrawlers アクション \(Python: batch\\_get\\_crawlers\)](#)
- [ListCrawlers アクション \(Python: list\\_crawlers\)](#)
- [ListCrawls アクション \(Python: list\\_crawls\)](#)

### CreateCrawler アクション (Python: create\_crawler)

指定されたターゲット、ロール、設定、およびオプションのスケジュールを使用して、新しいクローラーを作成します。s3Targets、jdbcTargets、またはDynamoDBTargets フィールドで、少なくとも 1 つ以上のクローラーを指定する必要があります。

#### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

新しいクローラーの名前。

- Role – 必須: UTF-8 文字列。

新しいクローラーが顧客リソースにアクセスするために使用する IAM ロール、または IAM ロールの Amazon リソースネーム (ARN)。

- DatabaseName – UTF-8 文字列。

次のような結果が書き込まれる AWS Glue データベース。arn:aws:daylight:us-east-1::database/sometable/\*

- Description – 説明文字列、2048 バイト長以下、「[URI address multi-line string pattern](#)」に一致。

新しいクローラーの説明。

- **Targets** – 必須: [CrawlerTargets](#) オブジェクト。  
クローリングするターゲットのコレクションのリスト。
- **Schedule** – UTF-8 文字列。  
スケジュールを指定するために使用される cron 式 ([ジョブとクローラーの時間ベースのスケジュール](#)を参照してください。) たとえば、毎日 12:15 UTC に何かを実行するには、cron(15 12 \* \* ? \*) を指定します。
- **Classifiers** – UTF-8 文字列の配列。  
ユーザーが登録したカスタム分類子のリスト。デフォルトでは、すべての組み込みの分類子がクローリングに含まれますが、これらのカスタム分類子によって常に分類別のデフォルトの分類子が上書きされます。
- **TablePrefix** - UTF-8 文字列。128 バイト長以下。  
作成されたカタログテーブルに使用されるテーブルプレフィックス。
- **SchemaChangePolicy** – [SchemaChangePolicy](#) オブジェクト。  
クローラーの更新と削除動作のためのポリシー。
- **RecrawlPolicy** – [RecrawlPolicy](#) オブジェクト。  
データセット全体を再度クローリングするか、前回のクローラー実行以降に追加されたフォルダのみをクローリングするかを指定するポリシー。
- **LineageConfiguration** – [LineageConfiguration](#) オブジェクト。  
クローラーのデータシステム設定を指定します。
- **LakeFormationConfiguration** – [LakeFormationConfiguration](#) オブジェクト。  
クローラー AWS Lake Formation の構成設定を指定します。
- **Configuration** – UTF-8 文字列。  
クローラーの構成情報。このバージョン付きの JSON 文字列では、クローラーの動作特性を指定できます。詳細については、「[クローラー設定オプションの設定](#)」を参照してください。
- **CrawlerSecurityConfiguration** - UTF-8 文字列。128 バイト長以下。  
このクローラーで使用される SecurityConfiguration 構造の名前。
- **Tags** – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

このクローラーリクエストで使用するタグ。クローラーへのアクセスを制限するためにタグを使用することができます。のタグの詳細については AWS Glue、デベロッパーガイドの[AWS 「のタグ AWS Glue」](#)を参照してください。

## レスポンス

- 応答パラメータはありません。

## エラー

- `InvalidInputException`
- `AlreadyExistsException`
- `OperationTimeoutException`
- `ResourceNumberLimitExceededException`

## DeleteCrawler アクション (Python: `delete_crawler`)

クローラーの状態が `STOPPED` でない限り AWS Glue Data Catalog、指定されたクローラーを `STOPPED` から削除します `RUNNING`。

## リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

削除するクローラーの名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- `EntityNotFoundException`

- `CrawlerRunningException`
- `SchedulerTransitioningException`
- `OperationTimeoutException`

## GetCrawler アクション (Python: `get_crawler`)

指定されたクローラーのメタデータを取得します。

### リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
メタデータを取得するクローラーの名前。

### レスポンス

- `Crawler` – [Crawler](#) オブジェクト。  
指定されたクローラーのメタデータ。

### エラー

- `EntityNotFoundException`
- `OperationTimeoutException`

## GetCrawlers アクション (Python: `get_crawlers`)

顧客アカウントで定義されたすべてのクローラーのメタデータを取得します。

### リクエスト

- `MaxResults` – 1~1000 の数値 (整数)。  
各呼び出しで返されるクローラーの数。
- `NextToken` – UTF-8 文字列。  
継続トークン (これが継続リクエストの場合)。

## 応答

- Crawlers – [Crawler](#) オブジェクトの配列。

クローラーメタデータのリスト。

- NextToken – UTF-8 文字列。

継続トークン (返されるリストがこの顧客アカウントで定義されたリストの最後に達していない場合)。

## エラー

- OperationTimeoutException

## GetCrawlerMetrics アクション (Python: `get_crawler_metrics`)

指定されたクローラーに関するメトリクスを取得します。

### リクエスト

- CrawlerNameList - UTF-8 文字列の配列、文字列 100 個以下。

メトリクスを取得するクローラーの名前のリスト。

- MaxResults – 1 ~ 1000 の数値 (整数)。

返されるリストの最大サイズ。

- NextToken – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

## 応答

- CrawlerMetricsList – [CrawlerMetrics](#) オブジェクトの配列。

指定されたクローラーのメトリクスのリスト。

- NextToken – UTF-8 文字列。

継続トークン (戻されたリストに最後に使用可能なメトリクスが含まれていない場合)。

## エラー

- `OperationTimeoutException`

### UpdateCrawler アクション (Python: `update_crawler`)

クローラーを更新します。クローラーが実行されている場合、クローラーを更新する前に `StopCrawler` を使用してクローラーを停止する必要があります。

#### リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

新しいクローラーの名前。

- `Role` – UTF-8 文字列。

新しいクローラーが顧客リソースにアクセスするために使用する IAM ロール、または IAM ロールの Amazon リソースネーム (ARN)。

- `DatabaseName` – UTF-8 文字列。

結果が保存される AWS Glue データベース。例: `arn:aws:daylight:us-east-1::database/sometable/*`。

- `Description` – UTF-8 文字列。2,048 バイト長以下。[URI address multi-line string pattern](#) に一致。

新しいクローラーの説明。

- `Targets` – [CrawlerTargets](#) オブジェクト。

クロールするターゲットのリスト。

- `Schedule` – UTF-8 文字列。

スケジュールを指定するために使用される cron 式 ([ジョブとクローラーの時間ベースのスケジュール](#)を参照してください。) たとえば、毎日 12:15 UTC に何かを実行するには、`cron(15 12 * * ? *)` を指定します。

- `Classifiers` – UTF-8 文字列の配列。

ユーザーが登録したカスタム分類子のリスト。デフォルトでは、すべての組み込みの分類子がクロールに含まれますが、これらのカスタム分類子によって常に分類別のデフォルトの分類子が上書きされます。

- TablePrefix - UTF-8 文字列。128 バイト長以下。

作成されたカタログテーブルに使用されるテーブルプレフィックス。

- SchemaChangePolicy – [SchemaChangePolicy](#) オブジェクト。

クローラーの更新と削除動作のためのポリシー。

- RecrawlPolicy – [RecrawlPolicy](#) オブジェクト。

データセット全体を再度クロールするか、前回のクローラー実行以降に追加されたフォルダのみをクロールするかを指定するポリシー。

- LineageConfiguration – [LineageConfiguration](#) オブジェクト。

クローラーのデータシステム設定を指定します。

- LakeFormationConfiguration – [LakeFormationConfiguration](#) オブジェクト。

クローラー AWS Lake Formation の構成設定を指定します。

- Configuration – UTF-8 文字列。

クローラーの構成情報。このバージョン付きの JSON 文字列では、クローラーの動作特性を指定できません。詳細については、「[クローラー設定オプションの設定](#)」を参照してください。

- CrawlerSecurityConfiguration - UTF-8 文字列。128 バイト長以下。

このクローラーで使用される SecurityConfiguration 構造の名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- InvalidInputException
- VersionMismatchException
- EntityNotFoundException
- CrawlerRunningException
- OperationTimeoutException

## StartCrawler アクション (Python: start\_crawler)

スケジュールされているものに関係なく、指定されたクローラーを使用してクローリングを開始します。クローラーがすでに実行されている場合、はを返します [CrawlerRunningException](#)。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
開始するクローラーの名前。

### レスポンス

- 応答パラメータはありません。

### エラー

- EntityNotFoundException
- CrawlerRunningException
- OperationTimeoutException

## StopCrawler アクション (Python: stop\_crawler)

指定されたクローラーが実行されている場合は、クローリングを停止します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
停止するクローラーの名前。

### レスポンス

- 応答パラメータはありません。

### エラー

- EntityNotFoundException

- `CrawlerNotRunningException`
- `CrawlerStoppingException`
- `OperationTimeoutException`

## BatchGetCrawlers アクション (Python: `batch_get_crawlers`)

指定されたクローラー名のリストのリソースメタデータのリストを返します。ListCrawlers オペレーションを呼び出した後で、このオペレーションを呼び出すことで、アクセス許可が付与されているデータにアクセスできます。このオペレーションは、タグを使用するアクセス許可条件を含め、すべての IAM のアクセス許可をサポートします。

### リクエスト

- `CrawlerNames` – 必須: UTF-8 文字列の配列。文字列 100 個以下。

クローラー名のリスト。これは ListCrawlers 操作から返された名前でもあります。

### 応答

- `Crawlers` – [Crawler](#) オブジェクトの配列。

クローラー定義のリスト。

- `CrawlersNotFound` - UTF-8 文字列の配列、文字列 100 個以下。

クローラーの名前のリストが見つかりません。

### エラー

- `InvalidInputException`
- `OperationTimeoutException`

## ListCrawlers アクション (Python: `list_crawlers`)

この AWS アカウントのすべてのクローラーリソース、または指定されたタグを持つリソースの名前を取得します。このオペレーションにより、アカウントで利用可能なリソースとその名前を確認できます。

このオペレーションはオプションの Tags フィールドを受け取ります。このフィールドを応答のフィルターとして使用すると、タグ付きリソースをグループとして取得できます。タグフィルタリングの使用を選択した場合は、タグが付いたリソースのみが取得されます。

## リクエスト

- MaxResults – 1 ~ 1000 の数値 (整数)。

返されるリストの最大サイズ。

- NextToken – UTF-8 文字列。

継続トークン (これが継続リクエストの場合)。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1 ~ 128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

これらのタグ付きリソースのみを返すように指定します。

## レスポンス

- CrawlerNames - UTF-8 文字列の配列、文字列 100 個以下。

アカウント内のすべてのクローラーの名前、または指定されたタグを持つクローラーの名前。

- NextToken – UTF-8 文字列。

継続トークン (戻されたリストに最後に使用可能なメトリクスが含まれていない場合)。

## エラー

- OperationTimeoutException

## ListCrawls アクション (Python: list\_crawls)

指定されたクローラーのすべてのクローールを返します。クローラー履歴機能の起動日以降に発生したクローールのみを返し、最大 12 ヶ月分のクローールのみを保持します。古いクローールは返されません。

この API は以下の操作に使用できます。

- 指定されたクローラーのすべてのクロールを取得します。
- 指定されたクローラーのすべてのクロールを、制限されたカウント内に取得します。
- 特定された時間範囲内に、指定されたクローラーのすべてのクロールを取得します。
- 特定の状態、クロール ID、または DPU 時間値を持つ指定されたクローラーのすべてのクロールを取得します。

## リクエスト

- CrawlerName – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

取得する実行のクローラー名。

- MaxResults – 1~1000 の数値 (整数)。

返される結果の最大数。デフォルトは 20 で、最大は 100 です。

- Filters – [CrawlsFilter](#) オブジェクトの配列。

CrawlsFilter オブジェクトのリストで指定した条件で、クロールをフィルタリングします。

- NextToken – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

## 応答

- Crawls – [CrawlerHistory](#) オブジェクトの配列。

条件を満たしたクロール実行を表す CrawlerHistory オブジェクトのリスト。

- NextToken – UTF-8 文字列。

返されたトークンのリストをページ分割するための継続トークン。リストの現在のセグメントが最後のセグメントではない場合に返されます。

## エラー

- EntityNotFoundException
- OperationTimeoutException
- InvalidInputException

## 列統計 API

列統計 API は、テーブル内の列の統計を返す AWS Glue API の説明を記述します。

### データ型

- [ColumnStatisticsTaskRun の構造](#)
- [ColumnStatisticsTaskRunningException の構造](#)
- [ColumnStatisticsTaskNotRunningException の構造](#)
- [ColumnStatisticsTaskStoppingException の構造](#)

### ColumnStatisticsTaskRun の構造

列統計実行の詳細を示すオブジェクト。

#### フィールド

- `CustomerId` - UTF-8 文字列。12 バイト長以下。

AWS アカウント ID。

- `ColumnStatisticsTaskRunId` - UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

特定の列統計タスク実行の識別子。

- `DatabaseName` - UTF-8 文字列。

テーブルが存在するデータベース。

- `TableName` - UTF-8 文字列。

列統計が生成されるテーブルの名前。

- `ColumnNameList` - UTF-8 文字列の配列。

列名のリスト。何も指定されない場合は、テーブルのすべての列名がデフォルトで使用されます。

- `CatalogID` - カタログ ID 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルが存在するデータカタログの ID。提供されない場合は、AWS アカウント ID がデフォルトで使用されます。

- Role – UTF-8 文字列。

統計を生成するためにサービスが引き受ける IAM ロール。

- SampleSize – 数値 (double)。100 以下。

統計の生成に使用される行の割合。何も指定されない場合は、テーブル全体が統計の生成に使用されます。

- SecurityConfiguration - UTF-8 文字列。128 バイト長以下。

列統計タスク実行の CloudWatch ログを暗号化するために使用されるセキュリティ設定の名前。

- NumberOfWorkers - 数値 (整数)。1 以上。

列統計の生成に使用されるワーカーの数。ジョブは、最大 25 個のインスタンスを自動スケールするように事前設定されています。

- WorkerType – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

統計の生成に使用されるワーカーのタイプ。デフォルト: g.1x。

- Status – UTF-8 文字列 (有効な値: STARTING | RUNNING | SUCCEEDED | FAILED | STOPPED)。

タスク実行のステータス。

- CreationTime – タイムスタンプ。

このタスクが作成された時刻。

- LastUpdated – タイムスタンプ。

このタスクが変更された前回の時点。

- StartTime – タイムスタンプ。

タスクの開始時刻。

- EndTime – タイムスタンプ。

タスクの終了時刻。

- ErrorMessage – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

ジョブのエラーメッセージ。

- DPUSeconds – 数値 (double)。None 以下。

すべての自動スケールされたワーカーについて計算された DPU 使用量 (秒)。

## ColumnStatisticsTaskRunningException の構造

列統計生成ジョブの実行中に別のジョブの開始を試行するとスローされる例外。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## ColumnStatisticsTaskNotRunningException の構造

実行中のタスクがないときにタスク実行の停止を試行するとスローされる例外。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## ColumnStatisticsTaskStoppingException の構造

タスク実行の停止を試行するとスローされる例外。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## 操作

- [StartColumnStatisticsTaskRun アクション \(Python: start\\_column\\_statistics\\_task\\_run\)](#)
- [GetColumnStatisticsTaskRun アクション \(Python: get\\_column\\_statistics\\_task\\_run\)](#)
- [GetColumnStatisticsTaskRuns アクション \(Python: get\\_column\\_statistics\\_task\\_runs\)](#)
- [ListColumnStatisticsTaskRuns アクション \(Python: list\\_column\\_statistics\\_task\\_runs\)](#)
- [StopColumnStatisticsTaskRun アクション \(Python: stop\\_column\\_statistics\\_task\\_run\)](#)

## StartColumnStatisticsTaskRun アクション (Python: start\_column\_statistics\_task\_run)

指定されたテーブルと列のために列統計タスク実行を開始します。

### リクエスト

- **DatabaseName** – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
テーブルが存在するデータベースの名前。
- **TableName** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
統計を生成するテーブルの名前。
- **ColumnNameList** – UTF-8 文字列の配列。  
統計を生成する列名のリスト。何も指定されない場合は、テーブルのすべての列名がデフォルトで使用されます。
- **Role** – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
統計を生成するためにサービスが引き受ける IAM ロール。
- **SampleSize** – 数値 (double)。100 以下。  
統計の生成に使用される行の割合。何も指定されない場合は、テーブル全体が統計の生成に使用されます。
- **CatalogID** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルが存在するデータカタログの ID。提供されない場合は、AWS アカウント ID がデフォルトで使用されます。
- **SecurityConfiguration** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
列統計タスク実行の CloudWatch ログを暗号化するために使用されるセキュリティ設定の名前。

### レスポンス

- **ColumnStatisticsTaskRunId** – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
列統計タスク実行の識別子。

## エラー

- AccessDeniedException
- EntityNotFoundException
- ColumnStatisticsTaskRunningException
- OperationTimeoutException
- ResourceNumberLimitExceededException
- InvalidInputException

### GetColumnStatisticsTaskRun アクション (Python: `get_column_statistics_task_run`)

タスク実行 ID を指定して、タスク実行に関連付けられたメタデータ/情報を取得します。

#### リクエスト

- ColumnStatisticsTaskRunId – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

特定の列統計タスク実行の識別子。

#### レスポンス

- ColumnStatisticsTaskRun – [ColumnStatisticsTaskRun](#) オブジェクト。

列統計実行の詳細を表す ColumnStatisticsTaskRun オブジェクト。

## エラー

- EntityNotFoundException
- OperationTimeoutException
- InvalidInputException

### GetColumnStatisticsTaskRuns アクション (Python: `get_column_statistics_task_runs`)

指定されたテーブルに関連付けられたすべての実行に関する情報を取得します。

## リクエスト

- `DatabaseName` – 必須: UTF-8 文字列。

テーブルが存在するデータベースの名前。

- `TableName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの名前。

- `MaxResults` – 1~1000 の数値 (整数)。

応答の最大サイズ。

- `NextToken` – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

## 応答

- `ColumnStatisticsTaskRuns` – [ColumnStatisticsTaskRun](#) オブジェクトの配列。

列統計タスク実行のリスト。

- `NextToken` – UTF-8 文字列。

継続トークン (すべてのタスク実行がまだ返されていない場合)。

## エラー

- `OperationTimeoutException`

## ListColumnStatisticsTaskRuns アクション (Python: `list_column_statistics_task_runs`)

特定のアカウントについてのすべてのタスク実行をリストします。

## リクエスト

- `MaxResults` – 1~1000 の数値 (整数)。

応答の最大サイズ。

- `NextToken` – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

## 応答

- ColumnStatisticsTaskRunIds - UTF-8 文字列の配列、文字列 100 個以下。

列統計タスク実行 ID のリスト。

- NextToken - UTF-8 文字列。

継続トークン (すべてのタスク実行 ID がまだ返されていない場合)。

## エラー

- OperationTimeoutException

## StopColumnStatisticsTaskRun アクション (Python: stop\_column\_statistics\_task\_run)

指定されたテーブルについてのタスク実行を停止します。

## リクエスト

- DatabaseName - 必須: UTF-8 文字列。

テーブルが存在するデータベースの名前。

- TableName - 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

テーブルの名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- ColumnStatisticsTaskNotRunningException
- ColumnStatisticsTaskStoppingException

- `OperationTimeoutException`

## クローラースケジューラ API

クローラースケジューラ API では、AWS Glue クローラーのデータ型と、クローラーを作成、削除、更新、および一覧表示するための API について説明します。

### データ型

- [Schedule 構造](#)

### Schedule 構造

cron ステートメントを使用してイベントをスケジュールするスケジューリングオブジェクト。

#### フィールド

- `ScheduleExpression` – UTF-8 文字列。

スケジュールを指定するために使用される cron 式 ([ジョブとクローラーの時間ベースのスケジュール](#)を参照してください。) たとえば、毎日 12:15 UTC に何かを実行するには、`cron(15 12 * * ? *)` を指定します。

- `State` – UTF-8 文字列 (有効な値: SCHEDULED | NOT\_SCHEDULED | TRANSITIONING)。

スケジュールの状態。

### 操作

- [UpdateCrawlerSchedule アクション \(Python: `update\_crawler\_schedule`\)](#)
- [StartCrawlerSchedule アクション \(Python: `start\_crawler\_schedule`\)](#)
- [StopCrawlerSchedule アクション \(Python: `stop\_crawler\_schedule`\)](#)

### UpdateCrawlerSchedule アクション (Python: `update_crawler_schedule`)

cron 式を使用してクローラーのスケジュールを更新します。

## リクエスト

- CrawlerName – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

スケジュールを更新するクローラーの名前。

- Schedule – UTF-8 文字列。

スケジュールを指定するために使用される更新された cron 式 ([ジョブとクローラーの時間ベースのスケジュール](#)を参照)。たとえば、毎日 12:15 UTC に何かを実行するには、cron(15 12 \* \* ? \*) を指定します。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- VersionMismatchException
- SchedulerTransitioningException
- OperationTimeoutException

## StartCrawlerSchedule アクション (Python: start\_crawler\_schedule)

クローラーがすでに実行中、またはスケジュールの状態がすでに SCHEDULED でなければ、指定されたクローラーのスケジュールの状態を SCHEDULED に変更します。

## リクエスト

- CrawlerName – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

スケジュールするクローラーの名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- SchedulerRunningException
- SchedulerTransitioningException
- NoScheduleException
- OperationTimeoutException

## StopCrawlerSchedule アクション (Python: stop\_crawler\_schedule)

指定されたクローラーのスケジュールの状態を NOT\_SCHEDULED に設定しますが、クローラーがすでに実行中の場合は停止されません。

### リクエスト

- CrawlerName – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
スケジュールの状態を設定するクローラーの名前。

### レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- SchedulerNotRunningException
- SchedulerTransitioningException
- OperationTimeoutException

## ETL スクリプト API の自動生成

ETL スクリプト生成 API では、AWS Glue で ETL スクリプトを生成するためのデータ型と API について説明します。

## データ型

- [CodeGenNode 構造](#)
- [CodeGenNodeArg 構造](#)
- [CodeGenEdge 構造](#)
- [場所の構造](#)
- [CatalogEntry 構造](#)
- [MappingEntry 構造](#)

### CodeGenNode 構造

Directed Acyclic Graph (DAG) でノードを表す

フィールド

- Id – 必須: UTF-8 文字列、1~255 バイト長、[Identifier string pattern](#) に一致。  
ノードのグラフ内で一意のノード識別子。
- NodeType – 必須: UTF-8 文字列。  
このノードのタイプ。
- Args – 必須: [CodeGenNodeArg](#) オブジェクトの配列。構造 50 個以下。  
ノードのプロパティ、名前と値のペアの形式。
- LineNumber – 数値 (整数)。  
ノードの行数。

### CodeGenNodeArg 構造

ノードの引数またはプロパティ。

フィールド

- Name – 必須: UTF-8 文字列。  
引数またはプロパティの名前。

- Value – 必須: UTF-8 文字列。

引数またはプロパティの値。

- Param – ブール。

値がパラメータとして使用される場合は True。

## CodeGenEdge 構造

Directed Acyclic Graph (DAG) で方向のエッジを表します。

フィールド

- Source – 必須: UTF-8 文字列、1~255 バイト長、[Identifier string pattern](#) に一致。

エッジが始まるノードの ID。

- Target – 必須: UTF-8 文字列、1~255 バイト長、[Identifier string pattern](#) に一致。

エッジが終了するノードの ID。

- TargetParameter – UTF-8 文字列。

エッジのターゲット。

## 場所の構造

リソースの場所。

フィールド

- Jdbc – [CodeGenNodeArg](#) オブジェクトの配列。構造 50 個以下。

JDBC の場所。

- S3 – [CodeGenNodeArg](#) オブジェクトの配列。構造 50 個以下。

Amazon Simple Storage Service (Amazon S3) の場所。

- DynamoDB – [CodeGenNodeArg](#) オブジェクトの配列。構造 50 個以下。

Amazon DynamoDB テーブルの場所。

## CatalogEntry 構造

AWS Glue Data Catalog でテーブル定義を指定します。

フィールド

- `DatabaseName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
テーブルメタデータが存在するデータベース。
- `TableName` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
該当するテーブルの名前。

## MappingEntry 構造

マッピングを定義します。

フィールド

- `SourceTable` – UTF-8 文字列。  
ソーステーブルの名前。
- `SourcePath` – UTF-8 文字列。  
ソースパス。
- `SourceType` – UTF-8 文字列。  
ソースタイプ。
- `TargetTable` – UTF-8 文字列。  
ターゲットテーブル。
- `TargetPath` – UTF-8 文字列。  
ターゲットパス。
- `TargetType` – UTF-8 文字列。  
ターゲットのタイプ。

## 操作

- [CreateScript アクション \(Python: create\\_script\)](#)
- [GetDataflowGraph アクション \(Python: get\\_dataflow\\_graph\)](#)
- [GetMapping アクション \(Python: get\\_mapping\)](#)
- [GetPlan アクション \(Python: get\\_plan\)](#)

### CreateScript アクション (Python: create\_script)

Directed Acyclic Graph (DAG) をコードに変換します。

#### リクエスト

- DagNodes – [CodeGenNode](#) オブジェクトの配列。  
DAG 内のノードのリスト。
- DagEdges – [CodeGenEdge](#) オブジェクトの配列。  
DAG 内のエッジのリスト。
- Language – UTF-8 文字列 (有効な値: PYTHON | SCALA)。  
DAG から生成されたコードのプログラミング言語。

#### レスポンス

- PythonScript – UTF-8 文字列。  
DAG から生成された Python スクリプト。
- ScalaCode – UTF-8 文字列。  
DAG から生成された Scala コード。

#### エラー

- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`

## GetDataflowGraph アクション (Python: get\_dataflow\_graph)

Python スクリプトを Directed Acyclic Graph (DAG) に変換します。

### リクエスト

- PythonScript – UTF-8 文字列。

変換する Python スクリプト。

### 応答

- DagNodes – [CodeGenNode](#) オブジェクトの配列。

結果の DAG 内のノードのリスト。

- DagEdges – [CodeGenEdge](#) オブジェクトの配列。

結果の DAG 内のエッジのリスト。

### エラー

- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## GetMapping アクション (Python: get\_mapping)

マッピングを作成します。

### リクエスト

- Source – 必須: [CatalogEntry](#) オブジェクト。

ソーステーブルを指定します。

- Sinks – [CatalogEntry](#) オブジェクトの配列。

ターゲットテーブルのリスト。

- Location – [ロケーション](#) オブジェクト。

マッピングのパラメータ。

## レスポンス

- Mapping – 必須: [MappingEntry](#) オブジェクトの配列。

指定されたターゲットへのマッピングのリスト。

## エラー

- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`
- `EntityNotFoundException`

## GetPlan アクション (Python: `get_plan`)

指定されたマッピングを実行するコードを取得します。

### リクエスト

- Mapping – 必須: [MappingEntry](#) オブジェクトの配列。

ソーステーブルからターゲットテーブルへのマッピングのリスト。

- Source – 必須: [CatalogEntry](#) オブジェクト。

ソーステーブル。

- Sinks – [CatalogEntry](#) オブジェクトの配列。

ターゲットテーブル。

- Location – [ロケーション](#) オブジェクト。

マッピングのパラメータ。

- Language – UTF-8 文字列 (有効な値: PYTHON | SCALA)。

マッピングを実行するコードのプログラミング言語。

- `AdditionalPlanOptionsMap` – キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

追加のオプションのキー値パラメータを保持するマップ。

現在、次のキーと値のペアがサポートされています。

- `inferSchema` – AWS Glue ジョブによって生成されたデフォルトスクリプトの `inferSchema` の設定を `true` にするか `false` にするかを指定します。例えば、`inferSchema` を `true` に設定するには、次のキーと値のペアを渡します。

```
--additional-plan-options-map '{"inferSchema":"true"}
```

## レスポンス

- `PythonScript` – UTF-8 文字列。

マッピングを実行する Python スクリプト。

- `ScalaCode` – UTF-8 文字列。

マッピングを実行する Scala コード。

## エラー

- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`

## ビジュアルジョブ API

ビジュアルジョブ API では、ジョブのビジュアル設定を表す JSON オブジェクトから AWS Glue API を使用して、データ統合 AWS Glue ジョブを作成できます。

のリスト `CodeGenConfigurationNodes` は、作成されたジョブの DAG を AWS Glue Studio に登録し、関連するコードを生成するために、ジョブの作成または更新 API に提供されます。

## データ型

- [CodeGenConfigurationNode 構造](#)
- [JDBC ConnectorOptions 構造](#)
- [StreamingDataPreviewOptions 構造](#)
- [AthenaConnectorSource 構造](#)
- [JDBC ConnectorSource 構造](#)
- [SparkConnectorSource 構造](#)
- [CatalogSource 構造](#)
- [MySQLCatalogSource 構造](#)
- [PostgreSQLCatalogSource 構造](#)
- [OracleSQLCatalogSource 構造](#)
- [MicrosoftSQLServerCatalogSource 構造](#)
- [CatalogKinesisSource 構造](#)
- [DirectKinesisSource 構造](#)
- [KinesisStreamingSourceOptions 構造](#)
- [CatalogKafkaSource 構造](#)
- [DirectKafkaSource 構造](#)
- [KafkaStreamingSourceOptions 構造](#)
- [RedshiftSource 構造](#)
- [AmazonRedshiftSource 構造](#)
- [AmazonRedshiftNodeData 構造](#)
- [AmazonRedshiftAdvancedOption 構造](#)
- [Option 構造](#)
- [S3CatalogSource 構造](#)
- [S3SourceAdditionalOptions 構造](#)
- [S3CsvSource 構造](#)
- [DirectJDBCSource 構造](#)
- [S3DirectSourceAdditionalOptions 構造](#)

- [S3JsonSource 構造](#)
- [S3ParquetSource 構造](#)
- [S3DeltaSource 構造](#)
- [S3CatalogDeltaSource 構造](#)
- [CatalogDeltaSource 構造](#)
- [S3HudiSource 構造](#)
- [S3CatalogHudiSource 構造](#)
- [CatalogHudiSource 構造](#)
- [DynamoDBCatalogSource 構造](#)
- [RelationalCatalogSource 構造](#)
- [JDBC ConnectorTarget 構造](#)
- [SparkConnectorTarget 構造](#)
- [BasicCatalogTarget 構造](#)
- [MySQLCatalogTarget 構造](#)
- [PostgreSQLCatalogTarget 構造](#)
- [OracleSQLCatalogTarget 構造](#)
- [MicrosoftSQLServerCatalogTarget 構造](#)
- [RedshiftTarget 構造](#)
- [AmazonRedshiftTarget 構造](#)
- [UpsertRedshiftTargetOptions 構造](#)
- [S3CatalogTarget 構造](#)
- [S3GlueParquetTarget 構造](#)
- [CatalogSchemaChangePolicy 構造](#)
- [S3DirectTarget 構造](#)
- [S3HudiCatalogTarget 構造](#)
- [S3HudiDirectTarget 構造](#)
- [S3DeltaCatalogTarget 構造](#)
- [S3DeltaDirectTarget 構造](#)

- [DirectSchemaChangePolicy 構造](#)
- [ApplyMapping 構造](#)
- [Mapping 構造](#)
- [SelectFields 構造](#)
- [DropFields 構造](#)
- [RenameField 構造](#)
- [スピゴット構造](#)
- [Join 構造](#)
- [JoinColumn 構造](#)
- [SplitFields 構造](#)
- [SelectFromCollection 構造](#)
- [FillMissingValues 構造](#)
- [Filter 構造](#)
- [FilterExpression 構造](#)
- [FilterValue 構造](#)
- [CustomCode 構造](#)
- [SparkSQL 構造](#)
- [SqlAlias 構造](#)
- [DropNullFields 構造](#)
- [NullCheckBoxList 構造](#)
- [NullValueField 構造](#)
- [Datatype 構造](#)
- [Merge 構造](#)
- [Union 構造](#)
- [PIIDetection 構造](#)
- [Aggregate 構造](#)
- [DropDuplicates 構造](#)
- [GovernedCatalogTarget 構造](#)
- [GovernedCatalogSource 構造](#)

- [AggregateOperation](#) 構造
- [GlueSchema](#) 構造
- [GlueStudioSchemaColumn](#) 構造
- [GlueStudioColumn](#) 構造
- [DynamicTransform](#) 構造
- [TransformConfigParameter](#) 構造
- [EvaluateDataQuality](#) 構造
- [DQ ResultsPublishingOptions](#) 構造
- [DQ StopJobOnFailureOptions](#) 構造
- [EvaluateDataQualityMultiFrame](#) 構造
- [Recipe](#) 構造
- [RecipeReference](#) 構造
- [SnowflakeNodeData](#) 構造
- [SnowflakeSource](#) 構造
- [SnowflakeTarget](#) 構造
- [ConnectorDataSource](#) 構造
- [ConnectorDataTarget](#) 構造

## CodeGenConfigurationNode 構造

CodeGenConfigurationNode は、すべての有効なノードタイプを列挙します。そのメンバー変数は 1 つしか入力できません。

### フィールド

- AthenaConnectorSource – [AthenaConnectorSource](#) オブジェクト。

Amazon Athena データソースへのコネクタを指定します。

- JDBCConectorSource – [JDBCConectorSource](#) オブジェクト。

JDBC データソースへのコネクタを指定します。

- SparkConnectorSource – [SparkConnectorSource](#) オブジェクト。

Apache Spark データソースへのコネクタを指定します。

- CatalogSource – [CatalogSource](#) オブジェクト。

AWS Glue データカタログ内のデータストアを指定します。

- RedshiftSource – [RedshiftSource](#) オブジェクト。

Amazon Redshift データストアを指定します。

- S3CatalogSource – [S3CatalogSource](#) オブジェクト。

データカタログ内の Amazon S3 AWS Glue データストアを指定します。

- S3CsvSource – [S3CsvSource](#) オブジェクト。

Amazon S3 に格納されているコマンド区切り値 (CSV) データストアを指定します。

- S3JsonSource – [S3JsonSource](#) オブジェクト。

Amazon S3 の JSON データストアを指定します。

- S3ParquetSource – [S3ParquetSource](#) オブジェクト。

Amazon S3 に保存されている Apache Parquet データストアを指定します。

- RelationalCatalogSource – [RelationalCatalogSource](#) オブジェクト。

AWS Glue データカタログ内のリレーショナルカタログデータストアを指定します。

- DynamoDBCatalogSource – [DynamoDBCatalogSource](#) オブジェクト。

データカタログ内の DynamoDB Catalog AWS Glue データストアを指定します。

- JDBCConnectorTarget – [JDBCConnectorTarget](#) オブジェクト。

Apache Parquet 列指向ストレージで Amazon S3 に書き込むデータターゲットを指定します。

- SparkConnectorTarget – [SparkConnectorTarget](#) オブジェクト。

Apache Spark コネクタを使用するターゲットを指定します。

- CatalogTarget – [BasicCatalogTarget](#) オブジェクト。

AWS Glue データカタログテーブルを使用するターゲットを指定します。

- RedshiftTarget – [RedshiftTarget](#) オブジェクト。

Amazon Redshift を使用するターゲットを指定します。

- `S3CatalogTarget` – [S3CatalogTarget](#) オブジェクト。

Data Catalog を使用して Amazon S3 に書き込む AWS Glue データターゲットを指定します。

- `S3GlueParquetTarget` – [S3GlueParquetTarget](#) オブジェクト。

Apache Parquet 列指向ストレージで Amazon S3 に書き込むデータターゲットを指定します。

- `S3DirectTarget` – [S3DirectTarget](#) オブジェクト。

Amazon S3 に書き込むデータターゲットを指定します。

- `ApplyMapping` – [ApplyMapping](#) オブジェクト。

データソースのマップデータプロパティキーを、データターゲットのデータプロパティキーに変換指定します。キーの名前を変更したり、データタイプを変更したり、データセットから削除するキーを選択できます。

- `SelectFields` – [SelectFields](#) オブジェクト。

保持するデータプロパティキーの選択変換を指定します。

- `DropFields` – [DropFields](#) オブジェクト。

削除するデータプロパティキーを選択する変換を指定します。

- `RenameField` – [RenameField](#) オブジェクト。

1 つのデータプロパティキーの名前を変更する変換を指定します。

- `Spigot` – [スピゴット](#) オブジェクト。

Amazon S3 バケットにデータのサンプルを書き込むための変換を指定します。

- `Join` – [Join](#) オブジェクト。

指定したデータプロパティキーの比較フレーズを使用して、2 つのデータセットを 1 つに結合する変換を指定します。結合タイプは、内部結合、外部結合、左結合、右結合、左半結合、左反結合を使用できます。

- `SplitFields` – [SplitFields](#) オブジェクト。

データプロパティキーを 2 つの `DynamicFrames` に分割する変換を指定します。出力は `DynamicFrames` のコレクションです。一方は選択したデータプロパティキー、他方は残っている方のデータプロパティキーを持ちます。

- `SelectFromCollection` – [SelectFromCollection](#) オブジェクト。

DynamicFrame のコレクションから 1 つの DynamicFrames を選択するトランスフォームを指定します。出力は選択された DynamicFrame です。

- FillMissingValues – [FillMissingValues](#) オブジェクト。

変換を使用して、データセット内に欠落値があるレコードを検索し、補完により決定する値を持つ新しいフィールドを追加します。入力データセットは、欠落値を決定する機械学習モデルのトレーニングに使用されます。

- Filter – [フィルター](#) オブジェクト。

フィルター条件に基づいて、データセットを 2 つに分割する変換を指定します。

- CustomCode – [CustomCode](#) オブジェクト。

データ変換を実行するためにカスタムコードを使用する変換を指定します。出力は のコレクションです DynamicFrames。

- SparkSQL – [SparkSQL](#) オブジェクト。

データを変換するために Spark SQL 構文を使用して、SQL クエリを入力する変換を指定します。出力は、単一の DynamicFrame です。

- DirectKinesisSource – [DirectKinesisSource](#) オブジェクト。

直接 Amazon Kinesis データソースを指定します。

- DirectKafkaSource – [DirectKafkaSource](#) オブジェクト。

Apache Kafka データストアを指定します。

- CatalogKinesisSource – [CatalogKinesisSource](#) オブジェクト。

データカタログ内の Kinesis AWS Glue データソースを指定します。

- CatalogKafkaSource – [CatalogKafkaSource](#) オブジェクト。

データカタログで Apache Kafka データストアを指定します。

- DropNullFields – [DropNullFields](#) オブジェクト。

列のすべての値が Null である場合に、データセットから列を削除する変換を指定します。デフォルトでは、AWS Glue Studio は null オブジェクトを認識しますが、空の文字列、「null」である文字列、-1 整数、またはゼロなどの他のプレースホルダーなどの一部の値は、自動的に null として認識されません。

- Merge – [マージ](#) オブジェクト。

レコードを識別するために、DynamicFrame プライマリキーに基づく DynamicFrame ステージングに結合変換を指定します。重複レコード ( 同じプライマリキーを持つレコード ) は重複除外されません。

- Union – [Union](#) オブジェクト。

2 つ以上のデータセットの行を 1 つの結果に結合する変換を指定します。

- PII Detection – [PIIDetection](#) オブジェクト。

PII データを識別、削除、またはマスクする変換を指定します。

- Aggregate – [集計](#) オブジェクト。

選択したフィールドによって行をグループ化し、指定された関数を使用して集計値を計算する変換を指定します。

- DropDuplicates – [DropDuplicates](#) オブジェクト。

繰り返しデータの行をデータセットから削除する変換を指定します。

- GovernedCatalogTarget – [GovernedCatalogTarget](#) オブジェクト。

管理されたカタログに書き込むデータターゲットを指定します。

- GovernedCatalogSource – [GovernedCatalogSource](#) オブジェクト。

管理されたデータカタログ内のデータソースを指定します。

- MicrosoftSQLServerCatalogSource – [MicrosoftSQLServerCatalogSource](#) オブジェクト。

AWS Glue データカタログ内の Microsoft SQL Server データソースを指定します。

- MySQLCatalogSource – [MySQLCatalogSource](#) オブジェクト。

AWS Glue データカタログ内の MySQL データソースを指定します。

- OracleSQLCatalogSource – [OracleSQLCatalogSource](#) オブジェクト。

Data Catalog で Oracle AWS Glue データソースを指定します。

- PostgreSQLCatalogSource – [PostgreSQLCatalogSource](#) オブジェクト。

AWS Glue データカタログ内の PostgreSQL データソースを指定します。

- MicrosoftSQLServerCatalogTarget – [MicrosoftSQLServerCatalogTarget](#) オブジェクト。

Microsoft SQL を使用するターゲットを指定します。

- MySQLCatalogTarget – [MySQLCatalogTarget](#) オブジェクト。

MySQL を使用するターゲットを指定します。

- OracleSQLCatalogTarget – [OracleSQLCatalogTarget](#) オブジェクト。

Oracle SQL を使用するターゲットを指定します。

- PostgreSQLCatalogTarget – [PostgreSQLCatalogTarget](#) オブジェクト。

Postgres SQL を使用するターゲットを指定します。

- DynamicTransform – [DynamicTransform](#) オブジェクト。

ユーザーが作成したカスタムビジュアル変換を指定します。

- EvaluateDataQuality – [EvaluateDataQuality](#) オブジェクト。

データ品質評価基準を指定します。

- S3CatalogHudiSource – [S3CatalogHudiSource](#) オブジェクト。

AWS Glue データカタログに登録されている Hudi データソースを指定します。データソースは に保存する必要があります Amazon S3。

- CatalogHudiSource – [CatalogHudiSource](#) オブジェクト。

AWS Glue データカタログに登録されている Hudi データソースを指定します。

- S3HudiSource – [S3HudiSource](#) オブジェクト。

に保存されている Hudi データソースを指定します Amazon S3。

- S3HudiCatalogTarget – [S3HudiCatalogTarget](#) オブジェクト。

AWS Glue データカタログ内の Hudi データソースに書き込むターゲットを指定します。

- S3HudiDirectTarget – [S3HudiDirectTarget](#) オブジェクト。

で Hudi データソースに書き込むターゲットを指定します Amazon S3。

- S3CatalogDeltaSource – [S3CatalogDeltaSource](#) オブジェクト。

AWS Glue データカタログに登録されている Delta Lake データソースを指定します。データソースは に保存する必要があります Amazon S3。

- CatalogDeltaSource – [CatalogDeltaSource](#) オブジェクト。

AWS Glue データカタログに登録されている Delta Lake データソースを指定します。

- S3DeltaSource – [S3DeltaSource](#) オブジェクト。

に保存されている Delta Lake データソースを指定します Amazon S3。

- S3DeltaCatalogTarget – [S3DeltaCatalogTarget](#) オブジェクト。

AWS Glue データカタログ内の Delta Lake データソースに書き込むターゲットを指定します。

- S3DeltaDirectTarget – [S3DeltaDirectTarget](#) オブジェクト。

で Delta Lake データソースに書き込むターゲットを指定します Amazon S3。

- AmazonRedshiftSource – [AmazonRedshiftSource](#) オブジェクト。

Amazon Redshift 内のデータソースに書き込むターゲットを指定します。

- AmazonRedshiftTarget – [AmazonRedshiftTarget](#) オブジェクト。

Amazon Redshift 内のデータターゲットに書き込むターゲットを指定します。

- EvaluateDataQualityMultiFrame – [EvaluateDataQualityMultiFrame](#) オブジェクト。

データ品質評価基準を指定します。複数の入力データを許可し、DynamicFrames のコレクションを返します。

- Recipe – [レシピ](#) オブジェクト。

AWS Glue DataBrew レシピノードを指定します。

- SnowflakeSource – [SnowflakeSource](#) オブジェクト。

Snowflake データソースを指定します。

- SnowflakeTarget – [SnowflakeTarget](#) オブジェクト。

Snowflake データソースに書き込むターゲットを指定します。

- ConnectorDataSource – [ConnectorDataSource](#) オブジェクト。

標準の接続オプションを使用して生成されたソースを指定します。

- ConnectorDataTarget – [ConnectorDataTarget](#) オブジェクト。

標準の接続オプションを使用して生成されたターゲットを指定します。

## JDBC ConnectorOptions 構造

コネクタの追加接続オプション。

## フィールド

- FilterPredicate– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

ソースからのデータをフィルタリングする追加の条件句。例:

```
BillingCity='Mountain View'
```

テーブル名の代わりにクエリを使用する場合は、指定された filterPredicate でクエリが動作することを確認します。

- PartitionColumn– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

パーティション化に使用される整数カラムの名前を示す文字列。このオプションは、lowerBound、upperBound、および numPartitions に含まれている場合にのみ機能します。このオプションの機能は、Spark SQL JDBC リーダーのものと同様です。

- LowerBound – 数値 (long)。None 以下。

パーティションストライドを決定するために使用される partitionColumn の最小値を示す整数。

- UpperBound – 数値 (long)。None 以下。

パーティションストライドを決定するために使用される partitionColumn の最大値を示す整数。

- NumPartitions – 数値 (long)。None 以下。

ターゲットパーティション数。この値は、(範囲に含まれる) lowerBound と (範囲に含まれない) upperBound とともに使用され、partitionColumn の分割で使用するために生成された WHERE 句の式のための、パーティションストライドを形成します。

- JobBookmarkKeys – UTF-8 文字列の配列。

ソートするジョブブックマークキーの名前。

- JobBookmarkKeysSortOrder– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

昇順または降順のソート順を指定します。

- DataTypeMapping – キーバリューペアのマップ配列。

各キーは UTF-8 文字列です。(有効な値: ARRAY | BIGINT | BINARY | BIT | BLOB | BOOLEAN | CHAR | CLOB | DATALINK | DATE | DECIMAL | DISTINCT | DOUBLE | FLOAT | INTEGER | JAVA\_OBJECT | LONGNVARCHAR | LONGVARBINARY | LONGVARCHAR | NCHAR | NCLOB | NULL

| NUMERIC | NVARCHAR | OTHER | REAL | REF | REF\_CURSOR | ROWID | SMALLINT | SQLXML  
| STRUCT | TIME | TIME\_WITH\_TIMEZONE | TIMESTAMP | TIMESTAMP\_WITH\_TIMEZONE |  
TINYINT | VARBINARY | VARCHAR )

各値は UTF-8 文字列です (有効な値: DATE | STRING | TIMESTAMP | INT | FLOAT | LONG |  
BIGDECIMAL | BYTE | SHORT | DOUBLE)。

JDBC データタイプ から AWS Glue データタイプ に対するマッピングを構築する、カスタムのデータタイプマッピング。例えば、オプションは、ドライバーの `ResultSet.getString()` メソッドを呼び出し `FLOAT` で JDBC タイプのデータフィールドを `Java String` タイプに `"dataTypeMapping":{"FLOAT":"STRING"}` マッピングし、それを使用して AWS Glue レコードを構築します。 `ResultSet` オブジェクトは各ドライバによって実装されるため、その動作は使用するドライバにより決定されます。ドライバによる変換の実行方法については、JDBC ドライバのドキュメントを参照してください。

## StreamingDataPreviewOptions 構造

データのサンプルを表示するためのデータプレビューに関連するオプションを指定します。

フィールド

- `PollingTime` – 10 以上の数値 (long)。  
ミリ秒単位のポーリング時間。
- `RecordPollingLimit` – 1 以上の数値 (long)。  
ポーリングされるレコード数の制限。

## AthenaConnectorSource 構造

Amazon Athena データソースへのコネクタを指定します。

フィールド

- `Name` – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
データソースの名前。
- `ConnectionName` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

コネクタに関連付けられている接続の名前。

- ConnectorName – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

AWS Glue Studio のデータストアへのアクセスを支援するコネクタの名前。

- ConnectionType – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

Amazon Athena データストアへの接続を指定する marketplace.athena や custom.athena など接続のタイプ。

- ConnectionTable – UTF-8 文字列、「[Custom string pattern #41](#)」に一致。

データソース内のテーブルの名前。

- SchemaName – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取り元となる Cloudwatch ロググループの名前。例えば、/aws-glue/jobs/output。

- OutputSchemas – [GlueSchema](#) オブジェクトの配列。

カスタム Athena ソース用のデータスキーマを指定します。

## JDBC ConnectorSource 構造

JDBC データソースへのコネクタを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データソースの名前。

- ConnectionName – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

コネクタに関連付けられている接続の名前。

- ConnectorName – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

AWS Glue Studio のデータストアへのアクセスを支援するコネクタの名前。

- ConnectionType – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

JDBC データストアへの接続を指定する marketplace.jdbc や custom.jdbc など接続のタイプ。

- AdditionalOptions – [JDBCConnectorOptions](#) オブジェクト。

コネクタの追加接続オプション。

- ConnectionTable – UTF-8 文字列、「[Custom string pattern #41](#)」に一致。

データソース内のテーブルの名前。

- Query – UTF-8 文字列、「[Custom string pattern #42](#)」に一致。

データを取得するテーブルまたは SQL クエリ。ConnectionTable または query を指定できません。両方を指定することはできません。

- OutputSchemas – [GlueSchema](#) オブジェクトの配列。

カスタム JDBC ソース用のデータスキーマを指定します。

## SparkConnectorSource 構造

Apache Spark データソースへのコネクタを指定します。

フィールド

- Name – 必須: UTF-8 文字列。[Custom string pattern #43](#) に一致。

データソースの名前。

- ConnectionName – 必須: UTF-8 文字列。[Custom string pattern #40](#) に一致。

コネクタに関連付けられている接続の名前。

- ConnectorName – 必須: UTF-8 文字列。[Custom string pattern #40](#) に一致。

AWS Glue Studio のデータストアへのアクセスを支援するコネクタの名前。

- ConnectionType – 必須: UTF-8 文字列。[Custom string pattern #40](#) に一致。

Apache Spark データストアへの接続を指定する marketplace.spark や custom.spark などの接続のタイプ。

- AdditionalOptions – キーバリュペアのマップ配列。

各キーは、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

コネクタの追加接続オプション。

- OutputSchemas – [GlueSchema](#) オブジェクトの配列。  
カスタム Spark ソース用のデータスキーマを指定します。

## CatalogSource 構造

AWS Glue データカタログ内のデータストアを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
データストアの名前。
- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースの名前。
- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースのテーブルの名前。

## MySQLCatalogSource 構造

AWS Glue データカタログ内の MySQL データソースを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
データソースの名前。
- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースの名前。
- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースのテーブルの名前。

## PostgreSQLCatalogSource 構造

AWS Glue データカタログ内の PostgreSQL データソースを指定します。

## フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
データソースの名前。
- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースの名前。
- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースのテーブルの名前。

## OracleSQLCatalogSource 構造

Data Catalog で Oracle AWS Glue データソースを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
データソースの名前。
- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースの名前。
- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースのテーブルの名前。

## MicrosoftSQLServerCatalogSource 構造

AWS Glue データカタログ内の Microsoft SQL Server データソースを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
データソースの名前。
- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースの名前。

- `Table` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースのテーブルの名前。

## CatalogKinesisSource 構造

データカタログ内の Kinesis AWS Glue データソースを指定します。

フィールド

- `Name` – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データソースの名前。

- `WindowSize` – 数値 (整数)、None 以下。

各マイクロバッチの処理にかかる時間。

- `DetectSchema` – ブール。

受信データからスキーマを自動的に決定するかどうか。

- `Table` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースのテーブルの名前。

- `Database` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースの名前。

- `StreamingOptions` – [KinesisStreamingSourceOptions](#) オブジェクト。

Kinesis ストリーミングデータソースの追加オプション。

- `DataPreviewOptions` – [StreamingDataPreviewOptions](#) オブジェクト。

データプレビューの追加オプション。

## DirectKinesisSource 構造

直接 Amazon Kinesis データソースを指定します。

フィールド

- `Name` – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データソースの名前。

- WindowSize – 数値 (整数)、None 以下。

各マイクロバッチの処理にかかる時間。

- DetectSchema – ブール。

受信データからスキーマを自動的に決定するかどうか。

- StreamingOptions – [KinesisStreamingSourceOptions](#) オブジェクト。

Kinesis ストリーミングデータソースの追加オプション。

- DataPreviewOptions – [StreamingDataPreviewOptions](#) オブジェクト。

データプレビューの追加オプション。

## KinesisStreamingSourceOptions 構造

Amazon Kinesis ストリーミングデータソースの追加オプション。

フィールド

- EndpointUrl – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

Kinesis エンドポイントの URL。

- StreamName – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

Kinesis データストリームの名前。

- Classification – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

オプションの分類。

- Delimiter – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

区切り記号文字を指定します。

- StartingPosition – UTF-8 文字列 (有効な値: latest="LATEST" | trim\_horizon="TRIM\_HORIZON" | earliest="EARLIEST" | timestamp="TIMESTAMP")。

Kinesis データストリーム内の、データの読み取り開始位置。指定できる値は "latest"、"trim\_horizon"、"earliest"、または UTC 形式のタイムスタンプ文字列であ

り、この文字列のパターンは yyyy-mm-ddTHH:MM:SSZ です (Z は UTC タイムゾーンのオフセットを +/- で表します。例: 「2023-04-04T08:00:00-04:00」)。デフォルト値は、"latest"です。

注:startingPosition」に UTC 形式のタイムスタンプ文字列である値は、AWS Glue バージョン 4.0 以降でのみサポートされています。

- MaxFetchTimeInMs – 数値 (long)。None 以下。

ジョブエグゼキューターが Kinesis データストリームから現在のバッチのレコードを読み取るために費やした最大時間は、ミリ秒 (ms) 単位で指定されます。この時間内に複数の GetRecords API コールを行うことができます。デフォルト値は、1000です。

- MaxFetchRecordsPerShard – 数値 (long)。None 以下。

1 マイクロバッチあたりに Kinesis データストリームでシャードごとにフェッチするレコードの最大数。メモ: ストリーミングジョブが既に Kinesis (同じ get-records 呼び出しで) から余分なレコードを読み取っている場合、クライアントはこの制限を超えることができません。MaxFetchRecordsPerShard が厳密である必要がある場合、MaxRecordPerRead の倍数にする必要があります。デフォルト値は、100000です。

- MaxRecordPerRead – 数値 (long)。None 以下。

各 getRecords オペレーションごとに、Kinesis データストリームからフェッチするレコードの最大数。デフォルト値は、10000です。

- AddIdleTimeBetweenReads – ブール。

2 つの連続する getRecords オペレーション間の遅延時間を追加します。デフォルト値は、"False"です。このオプションは、Glue バージョン 2.0 以降でのみ設定可能です。

- IdleTimeBetweenReadsInMs – 数値 (long)。None 以下。

2 つの連続する getRecords オペレーション間での、最短の遅延時間 (ミリ秒単位で指定)。デフォルト値は、1000です。このオプションは、Glue バージョン 2.0 以降でのみ設定可能です。

- DescribeShardInterval – 数値 (long)。None 以下。

スクリプトがリシャードイングを検討するための 2 つの ListShards API コール間の最小時間間隔。デフォルト値は、1sです。

- NumRetries – 数値 (整数)、None 以下。

Kinesis Data Streams API リクエストを再試行する最大の回数。デフォルト値は、3です。

- RetryIntervalMs – 数値 (long)。None 以下。

Kinesis Data Streams API 呼び出しを再試行するまでのクールオフ期間 (ミリ秒単位で指定)。デフォルト値は、1000です。

- MaxRetryIntervalMs – 数値 (long)。None 以下。

再試行で 2 つの Kinesis Data Streams API を呼び出す間の最大クールオフ期間 (ミリ秒単位で指定)。デフォルト値は、10000です。

- AvoidEmptyBatches – ブール。

バッチ処理を開始する前に、Kinesis データストリームで未読のデータをチェックすることで、空のマイクロバッチジョブを作成しないようにします。デフォルト値は、"False"です。

- StreamArn– UTF-8 文字列、[「Custom string pattern #40」](#)に一致。

Kinesis データストリームの Amazon リソースネーム (ARN)。

- RoleArn– UTF-8 文字列、[「Custom string pattern #40」](#)に一致。

AWS Security Token Service (AWS STS) の使用を引き受けるロールの、Amazon リソースネーム (ARN)。このロールには、Kinesis データストリームのレコードの説明操作または読み取り操作の権限が必要です。このパラメーターは、別のアカウントのデータストリームにアクセスするときに使用する必要があります。"awsSTSSessionName" と組み合わせて使用します。

- RoleSessionName– UTF-8 文字列、[「Custom string pattern #40」](#)に一致。

AWS STS を使用するロールを引き受ける、セッションの識別子。このパラメーターは、別のアカウントのデータストリームにアクセスするときに使用する必要があります。"awsSTSRoleARN" と組み合わせて使用します。

- AddRecordTimestamp– UTF-8 文字列、[「Custom string pattern #40」](#)に一致。

このオプションが「true」に設定されている場合、データ出力には、対応するレコードがストリームによって受信された時刻を表示する「\_\_src\_timestamp」という名前が付けられた追加の列が含まれます。デフォルト値は、「false」です。このオプションは、AWS Glue バージョン 4.0 以降でサポートされています。

- EmitConsumerLagMetrics– UTF-8 文字列、[「Custom string pattern #40」](#)に一致。

このオプションを「true」に設定すると、バッチごとに、ストリームが受信した最も古いレコードからに到着するまでの期間のメトリクスが出力されます AWS Glue CloudWatch。メトリクスの名前は「glue.driver.streaming」ですmaxConsumerLagInMs。デフォルト値は、「false」です。このオプションは AWS Glue バージョン 4.0 以降でサポートされています。

- StartingTimestamp – UTF-8 文字列。

データの読み取りを開始する Kinesis データストリーム内のレコードのタイムスタンプ。指定できる値は、UTC 形式のタイムスタンプ文字列です。この文字列のパターンは yyyy-mm-ddTHH:MM:SSZ です (Z は UTC タイムゾーンのオフセットを +/- で表します。例: 「2023-04-04T08:00:00+08:00」)。

## CatalogKafkaSource 構造

データカタログで Apache Kafka データストアを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データストアの名前。

- WindowSize – 数値 (整数)、None 以下。

各マイクロバッチの処理にかかる時間。

- DetectSchema – ブール。

受信データからスキーマを自動的に決定するかどうか。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースのテーブルの名前。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースの名前。

- StreamingOptions – [KafkaStreamingSourceOptions](#) オブジェクト。

ストリーミングオプションを指定します。

- DataPreviewOptions – [StreamingDataPreviewOptions](#) オブジェクト。

データのサンプルを表示するためのデータプレビューに関連するオプションを指定します。

## DirectKafkaSource 構造

Apache Kafka データストアを指定します。

## フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
データストアの名前。
- StreamingOptions – [KafkaStreamingSourceOptions](#) オブジェクト。  
ストリーミングオプションを指定します。
- WindowSize – 数値 (整数)、None 以下。  
各マイクロバッチの処理にかかる時間。
- DetectSchema – ブール。  
受信データからスキーマを自動的に決定するかどうか。
- DataPreviewOptions – [StreamingDataPreviewOptions](#) オブジェクト。  
データのサンプルを表示するためのデータプレビューに関連するオプションを指定します。

## KafkaStreamingSourceOptions 構造

ストリーミングの追加オプション。

### フィールド

- BootstrapServers– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。  
ブートストラップサーバーの URL のリスト (例: b-1.vpc-test-2.o4q88o.c6.kafka.us-east-1.amazonaws.com:9094)。このオプションは API 呼び出しで指定するか、データカタログ内のテーブルメタデータで定義する必要があります。
- SecurityProtocol– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。  
ブローカーと通信するために使用されるプロトコル。使用できる値は、"SSL" または "PLAINTEXT" です。
- ConnectionName– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。  
コレクションの名前。
- TopicName– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

Apache Kafka で指定されたトピック名。少なくとも

"topicName"、"assign"、"subscribePattern" の内いずれかを指定する必要があります。

- Assign– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

消費する特有の TopicPartitions。少なくとも

"topicName"、"assign"、"subscribePattern" の内いずれかを指定する必要があります。

- SubscribePattern– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

サブスクライブする先のトピックリストを識別する Java の正規表現文字列。少なくとも

"topicName"、"assign"、"subscribePattern" の内いずれかを指定する必要があります。

- Classification– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

オプションの分類。

- Delimiter– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

区切り記号文字を指定します。

- StartingOffsets– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

Kafka トピック内で、データの読み取りを開始する位置 使用できる値は、"earliest" または "latest" です。デフォルト値は "latest" です。

- EndingOffsets– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

バッチクエリの終了位置。設定が可能な値は、"latest" または、各 TopicPartition の終了オフセットを指定する JSON 文字列のいずれかです。

- PollTimeoutMs – 数値 (long)。None 以下。

Spark ジョブエグゼキュータで、Kafka からデータをポーリングする際のタイムアウト値 (ミリ秒単位)。デフォルト値は、512です。

- NumRetries – 数値 (整数)、None 以下。

Kafka オフセットのフェッチが失敗したと判断される前の再試行回数。デフォルト値は、3です。

- RetryIntervalMs – 数値 (long)。None 以下。

Kafka オフセットのフェッチを開始するまでの待機時間 (ミリ秒)。デフォルト値は、10です。

- MaxOffsetsPerTrigger – 数値 (long)。None 以下。

処理されるオフセットの最大数を、トリガー間隔ごとのレート上限で指定する値。指定されたオフセットの合計数は、異なるポリュームの topicPartitions 間で均等に分割されます。デフォルト値は「null」です。この場合、コンシューマーは既知の最新のオフセットまで、すべてのオフセットを読み取ります。

- MinPartitions – 数値 (整数)、None 以下。

Kafka から読み取ることを想定する、最小のパーティション数。デフォルト値は「null」です。これは、Spark パーティションの数が Kafka パーティションの数に等しいことを意味します。

- IncludeHeaders – ブール。

Kafka ヘッダーを含めるかどうかを決定します。このオプションが「true」に設定されている場合、データ出力には、「glue\_streaming\_kafka\_headers」という名前で Array[Struct(key: String, value: String)] 型の列が追加されます。デフォルト値は「false」です。このオプションは、AWS Glue バージョン 3.0 以降でのみ使用できます。

- AddRecordTimestamp– UTF-8 文字列、[「Custom string pattern #40」](#) に一致。

このオプションが「true」に設定されている場合、データ出力には、対応するレコードがトピックによって受信された時刻を表示する「\_\_src\_timestamp」という名前が付けられた追加の列が含まれます。デフォルト値は、「false」です。このオプションは、AWS Glue バージョン 4.0 以降でサポートされています。

- EmitConsumerLagMetrics– UTF-8 文字列、[「Custom string pattern #40」](#) に一致。

このオプションを「true」に設定すると、バッチごとに、トピックによって受信された最も古いレコードからに到着した時点までの期間のメトリクスが出力されます AWS Glue CloudWatch。メトリクスの名前は「glue.driver.streaming」ですmaxConsumerLagInMs。デフォルト値は、「false」です。このオプションは AWS Glue バージョン 4.0 以降でサポートされています。

- StartingTimestamp – UTF-8 文字列。

データの読み取りを開始する Kafka トピック内のレコードのタイムスタンプ。指定できる値は、UTC 形式のタイムスタンプ文字列です。この文字列のパターンは yyyy-mm-ddTHH:MM:SSZ です (Z は UTC タイムゾーンのオフセットを +/- で表します。例: 「2023-04-04T08:00:00+08:00」)。

StartingTimestamp または StartingOffsets のいずれかのみを設定する必要があります。

## RedshiftSource 構造

Amazon Redshift データストアを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

Amazon Redshift データストアの名前。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み込むデータベース。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取り元のデータベーステーブル。

- RedshiftTmpDir– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

データベースからコピーするときに一時データをステージングできる Amazon S3 パス。

- TmpDirIAMRole– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

アクセス許可を持つ IAM ロール。

## AmazonRedshiftSource 構造

Amazon Redshift ソースを指定します。

フィールド

- Name– UTF-8 文字列、「[Custom string pattern #43](#)」に一致。

Amazon Redshift ソースの名前。

- Data – [AmazonRedshiftNodeData](#) オブジェクト。

Amazon Redshift ソースノードのデータを指定します。

## AmazonRedshiftNodeData 構造

Amazon Redshift ノードを指定します。

## フィールド

- `AccessType`– UTF-8 文字列、「[Custom string pattern #39](#)」に一致。

Redshift 接続のアクセスタイプ。直接接続またはカタログ接続が可能です。

- `SourceType`– UTF-8 文字列、「[Custom string pattern #39](#)」に一致。

特定のテーブルがソースかカスタムクエリかを指定するソースタイプ。

- `Connection` – [オプション](#) オブジェクト。

Redshift クラスター AWS Glue への接続。

- `Schema` – [オプション](#) オブジェクト。

直接接続で作業するときの Redshift スキーマの名前。

- `Table` – [オプション](#) オブジェクト。

直接接続で作業するときの Redshift テーブルの名前。

- `CatalogDatabase` – [オプション](#) オブジェクト。

データカタログを使用する場合の AWS Glue Data Catalog データベースの名前。

- `CatalogTable` – [オプション](#) オブジェクト。

AWS Glue データカタログを使用する場合のデータカタログテーブル名。

- `CatalogRedshiftSchema` – UTF-8 文字列。

データカタログで作業するときの Redshift スキーマの名前。

- `CatalogRedshiftTable` – UTF-8 文字列。

読み取り元のデータベーステーブル。

- `TempDir`– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

データベースからコピーするとき一時データをステージングできる Amazon S3 パス。

- `IamRole` – [オプション](#) オブジェクト。

オプション。S3 に接続するとき使用するロールの名前。空欄のままにすると、IAM ロールはデフォルトでジョブのロールになります。

- `AdvancedOptions` – [AmazonRedshiftAdvancedOption](#) オブジェクトの配列。

Redshift クラスターに接続するときのオプションの値。

- SampleQuery – UTF-8 文字列。

が「クエリ」の場合 SourceType に Redshift ソースからデータを取得するために使用される SQL。

- PreAction – UTF-8 文字列。

upsert を用いる MERGE または APPEND を実行する前に使用される SQL。

- PostAction – UTF-8 文字列。

upsert を用いる MERGE または APPEND を実行する前に使用される SQL。

- Action – UTF-8 文字列。

Redshift クラスターへの書き込み方法を指定します。

- TablePrefix – UTF-8 文字列、「[Custom string pattern #39](#)」に一致。

テーブルへのプレフィックスを指定します。

- Upsert – ブール。

APPEND を実行するときに Redshift シンクで使用するアクション。

- MergeAction – UTF-8 文字列、「[Custom string pattern #39](#)」に一致。

Redshift シンク内の MERGE の処理方法を決定するときに使用するアクション。

- MergeWhenMatched – UTF-8 文字列、「[Custom string pattern #39](#)」に一致。

既存のレコードが新しいレコードと一致する場合、Redshift シンク内の MERGE の処理方法を決定するときに使用するアクション。

- MergeWhenNotMatched – UTF-8 文字列、「[Custom string pattern #39](#)」に一致。

既存のレコードが新しいレコードと一致しない場合、Redshift シンク内の MERGE の処理方法を決定するときに使用するアクション。

- MergeClause – UTF-8 文字列。

一致するレコードを処理するためにカスタムマージで使用される SQL。

- CrawlerConnection – UTF-8 文字列。

使用するカタログテーブルに関連する接続の名前を指定します。

- TableSchema – [オプション](#) オブジェクトの配列。

特定のノードにおけるスキーマ出力の配列。

- StagingTable – UTF-8 文字列。

upsert を用いる MERGE または APPEND を実行するときに使用する一時的なステージングテーブルの名前。

- SelectedColumns – [オプション](#) オブジェクトの配列。

upsert を用いる MERGE または APPEND を実行するときに、一致するレコードを決定するために使用する列の名前のリスト。

## AmazonRedshiftAdvancedOption 構造

Redshift クラスターに接続するときのオプションの値を指定します。

フィールド

- Key – UTF-8 文字列。

追加接続オプションのキー。

- Value – UTF-8 文字列。

追加接続オプションの値。

## Option 構造

オプションの値を指定します。

フィールド

- Value– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

オプションの値を指定します。

- Label– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

オプションのラベルを指定します。

- Description– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

オプションの説明を指定します。

## S3CatalogSource 構造

データカタログ内の Amazon S3 AWS Glue データストアを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データストアの名前。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み込むデータベース。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取り元のデータベーステーブル。

- PartitionPredicate – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

この述語を満たすパーティションは削除されます。これらのパーティションの保存期間内のファイルは削除されません。"" を設定 – デフォルトでは空です。

- AdditionalOptions – [S3SourceAdditionalOptions](#) オブジェクト。

追加の接続オプションを指定します。

## S3SourceAdditionalOptions 構造

Amazon S3 データストアの追加の接続オプションを指定します。

フィールド

- BoundedSize – 数値 (long 型)。

処理されるバイトのデータセットのターゲットサイズの上限を設定します。

- BoundedFiles – 数値 (long 型)。

処理されるファイルのターゲット数の上限を設定します。

## S3CsvSource 構造

Amazon S3 に格納されているコマンド区切り値 (CSV) データストアを指定します。

## フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データストアの名前。

- Paths – 必須: UTF-8 文字列の配列。

読み取りのソースとなる Amazon S3 パスのリスト。

- CompressionType – UTF-8 文字列 (有効な値: gzip="GZIP" | bzip2="BZIP2")。

データの圧縮方法を指定します。データに標準のファイル拡張子が付いている場合、このオプションは一般的に不要です。指定できる値は "gzip" および "bzip" です。

- Exclusions – UTF-8 文字列の配列。

除外する Unix スタイルの glob パターンの JSON リストを含む文字列。たとえば、`"[\"**/*.pdf\"]"` はすべての PDF ファイルを除外します。

- GroupSize – UTF-8 文字列、 [「Custom string pattern #40」](#) に一致。

ターゲットグループのサイズ (バイト単位)。デフォルトは、入力データのサイズとクラスターのサイズに基づいて計算されます。入力ファイルが 50,000 個未満の場合、このオプションを有効にするには、"groupFiles" を "inPartition" に設定する必要があります。

- GroupFiles – UTF-8 文字列、 [「Custom string pattern #40」](#) に一致。

入力ファイルが 50,000 個を超える場合、デフォルトでファイルのグループ化が有効化されます。入力ファイルが 50,000 個未満の場合にグループ化を有効化するには、このパラメータに "inPartition" を設定します。入力ファイルが 50,000 個を超える場合に、グループ化を無効にするには、このパラメータを "none" に設定します。

- Recurse – ブール。

true に設定した場合は、指定したパスの下にあるすべてのサブディレクトリ内のファイルを再帰的に読み取ります。

- MaxBand – 数値 (整数)、None 以下。

このオプションでは、s3 リストの一貫性が認められるまでの期間をミリ秒単位で指定します。変更タイムスタンプが最後の maxBand ミリ秒以内のファイルは、を使用して JobBookmarks Amazon S3 の結果整合性を考慮すると、特に追跡されます。ほとんどのユーザーはこのオプションを設定する必要はありません。デフォルトは 900000 ミリ秒 (15 分) です。

- MaxFilesInBand – 数値 (整数)、None 以下。

このオプションは、直前の maxBand 秒間に保存するファイルの最大数を指定します。この数を越えた場合、余分なファイルはスキップされ、次のジョブ実行時にのみ処理されます。

- AdditionalOptions – [S3DirectSourceAdditionalOptions](#) オブジェクト。

追加の接続オプションを指定します。

- Separator – 必須: UTF-8 文字列 (有効な値: comma="COMMA" | ctrlA="CTRLA" | pipe="PIPE" | semicolon="SEMICOLON" | tab="TAB")。

区切り記号文字を指定します。デフォルトではカンマ:"," ですが、他の任意の文字を指定できます。

- Escaper – UTF-8 文字列、「[Custom string pattern #41](#)」に一致。

エスケープに使用する文字を指定します。このオプションは、CSV ファイルを読み取る場合にのみ使用します。デフォルト値は、none です。有効にすると、直後の文字はそのまま使用されます。ただし、よく知られている小さいエスケープセット (\n、\r、\t、\0) を除きます。

- QuoteChar – 必須: UTF-8 文字列 (有効な値: quote="QUOTE" | quillemet="QUILLEMET" | single\_quote="SINGLE\_QUOTE" | disabled="DISABLED")。

引用に使用する文字を指定します。デフォルト値は二重引用符 '"' です。これに -1 を設定すると、全体的に引用が無効になります。

- Multiline – ブール。

単一のレコードが複数行にまたがることができるかどうかを指定するブール値。これが発生するのは、フィールドに引用符で囲まれた改行文字がある場合などです。複数行にまたがるレコードがある場合は、このオプションを True に設定する必要があります。デフォルト値は False であり、解析時によりアグレッシブなファイル分割を可能にします。

- WithHeader – ブール。

最初の行をヘッダーとして扱うかどうかを指定するブール値。デフォルト値は、False です。

- WriteHeader – ブール。

ヘッダーを出力に書き込むかどうかを指定するブール値。デフォルト値は、True です。

- SkipFirst – ブール。

最初のデータ行をスキップするかどうかを指定するブール値。デフォルト値は、False です。

- OptimizePerformance – ブール。

高度な SIMD CSV リーダーで、Apache Arrow ベースの列指向メモリ形式を使用するかどうかを指定するブール値。AWS Glue バージョン 3.0 でのみ使用できます。

- OutputSchemas – [GlueSchema](#) オブジェクトの配列。

S3 CSV ソース用のデータスキーマを指定します。

## DirectJDBCSource 構造

直接 JDBC ソース接続を指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

JDBC ソース接続の名前。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

JDBC ソース接続のデータベース。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

JDBC ソース接続のテーブル。

- ConnectionName – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

JDBC ソースの接続名。

- ConnectionType – 必須: UTF-8 文字列 (有効な値: sqlserver | mysql | oracle | postgresql | redshift)。

JDBC ソースの接続タイプ。

- RedshiftTmpDir – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

JDBC Redshift ソースの一時ディレクトリ。

## S3DirectSourceAdditionalOptions 構造

Amazon S3 データストアの追加の接続オプションを指定します。

## フィールド

- BoundedSize – 数値 (long 型)。

処理されるバイトのデータセットのターゲットサイズの上限を設定します。

- BoundedFiles – 数値 (long 型)。

処理されるファイルのターゲット数の上限を設定します。

- EnableSamplePath – ブール。

オプションを設定しサンプルパスを有効にします。

- SamplePath – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

有効にした場合は、サンプルパスを指定します。

## S3JsonSource 構造

Amazon S3 の JSON データストアを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。[Custom string pattern #43](#) に一致。

データストアの名前。

- Paths – 必須: UTF-8 文字列の配列。

読み取りのソースとなる Amazon S3 パスのリスト。

- CompressionType – UTF-8 文字列 (有効な値: gzip="GZIP" | bzip2="BZIP2")。

データの圧縮方法を指定します。データに標準のファイル拡張子が付いている場合、このオプションは一般的に不要です。指定できる値は "gzip" および "bzip" です。

- Exclusions – UTF-8 文字列の配列。

除外する Unix スタイルの glob パターンの JSON リストを含む文字列。たとえば、"[\"\*\*.\*pdf\"]" はすべての PDF ファイルを除外します。

- GroupSize – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

ターゲットグループのサイズ (バイト単位)。デフォルトは、入力データのサイズとクラスターのサイズに基づいて計算されます。入力ファイルが 50,000 個未満の場合、このオプションを有効にするには、"groupFiles" を "inPartition" に設定する必要があります。

- GroupFiles – UTF-8 文字列、[「Custom string pattern #40」](#) に一致。

入力ファイルが 50,000 個を超える場合、デフォルトでファイルのグループ化が有効化されます。入力ファイルが 50,000 個未満の場合にグループ化を有効化するには、このパラメータに "inPartition" を設定します。入力ファイルが 50,000 個を超える場合に、グループ化を無効にするには、このパラメータを "none" に設定します。

- Recurse – ブール。

true に設定した場合は、指定したパスの下にあるすべてのサブディレクトリ内のファイルを再帰的に読み取ります。

- MaxBand – 数値 (整数)、None 以下。

このオプションでは、s3 リストの一貫性が認められるまでの期間をミリ秒単位で指定します。変更タイムスタンプが最後の maxBand ミリ秒以内のファイルは、を使用して JobBookmarks Amazon S3 の結果整合性を考慮すると、特に追跡されます。ほとんどのユーザーはこのオプションを設定する必要はありません。デフォルトは 900000 ミリ秒 (15 分) です。

- MaxFilesInBand – 数値 (整数)、None 以下。

このオプションは、直前の maxBand 秒間に保存するファイルの最大数を指定します。この数を超えた場合、余分なファイルはスキップされ、次のジョブ実行時にのみ処理されます。

- AdditionalOptions – [S3DirectSourceAdditionalOptions](#) オブジェクト。

追加の接続オプションを指定します。

- JsonPath – UTF-8 文字列、[「Custom string pattern #40」](#) に一致。

JSON データを定義する JsonPath 文字列。

- Multiline – ブール。

単一のレコードが複数行にまたがるかどうかを指定するブール値。これが発生するのは、フィールドに引用符で囲まれた改行文字がある場合などです。複数行にまたがるレコードがある場合は、このオプションを True に設定する必要があります。デフォルト値は False であり、解析時によりアグレッシブなファイル分割を可能にします。

- OutputSchemas – [GlueSchema](#) オブジェクトの配列。

S3 JSON ソース用のデータスキーマを指定します。

## S3ParquetSource 構造

Amazon S3 に保存されている Apache Parquet データストアを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データストアの名前。

- Paths – 必須: UTF-8 文字列の配列。

読み取りのソースとなる Amazon S3 パスのリスト。

- CompressionType – UTF-8 文字列 (有効な値: snappy="SNAPPY" | lzo="LZO" | gzip="GZIP" | uncompressed="UNCOMPRESSED" | none="NONE")。

データの圧縮方法を指定します。データに標準のファイル拡張子が付いている場合、このオプションは一般的に不要です。指定できる値は "gzip" および "bzip" です。

- Exclusions – UTF-8 文字列の配列。

除外する Unix スタイルの glob パターンの JSON リストを含む文字列。たとえば、["/\*.pdf"] はすべての PDF ファイルを除外します。

- GroupSize – UTF-8 文字列、 [「Custom string pattern #40」](#) に一致。

ターゲットグループのサイズ (バイト単位)。デフォルトは、入力データのサイズとクラスターのサイズに基づいて計算されます。入力ファイルが 50,000 個未満の場合、このオプションを有効にするには、"groupFiles" を "inPartition" に設定する必要があります。

- GroupFiles – UTF-8 文字列、 [「Custom string pattern #40」](#) に一致。

入力ファイルが 50,000 個を超える場合、デフォルトでファイルのグループ化が有効化されます。入力ファイルが 50,000 個未満の場合にグループ化を有効化するには、このパラメータに "inPartition" を設定します。入力ファイルが 50,000 個を超える場合に、グループ化を無効にするには、このパラメータを "none" に設定します。

- Recurse – ブール。

true に設定した場合は、指定したパスの下にあるすべてのサブディレクトリ内のファイルを再帰的に読み取ります。

- MaxBand – 数値 (整数)、None 以下。

このオプションでは、s3 リストの一貫性が認められるまでの期間をミリ秒単位で指定します。変更タイムスタンプが最後の maxBand ミリ秒以内のファイルは、を使用して JobBookmarks Amazon S3 の結果整合性を考慮すると、特に追跡されます。ほとんどのユーザーはこのオプションを設定する必要はありません。デフォルトは 900000 ミリ秒 (15 分) です。

- MaxFilesInBand – 数値 (整数)、None 以下。

このオプションは、直前の maxBand 秒間に保存するファイルの最大数を指定します。この数を超えた場合、余分なファイルはスキップされ、次のジョブ実行時にのみ処理されます。

- AdditionalOptions – [S3DirectSourceAdditionalOptions](#) オブジェクト。

追加の接続オプションを指定します。

- OutputSchemas – [GlueSchema](#) オブジェクトの配列。

S3 Parquet ソース用のデータスキーマを指定します。

## S3DeltaSource 構造

に保存されている Delta Lake データソースを指定します Amazon S3。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

Delta Lake ソースの名前。

- Paths – 必須: UTF-8 文字列の配列。

読み取りのソースとなる Amazon S3 パスのリスト。

- AdditionalDeltaOptions – キーバリュペアのマップ配列。

各キーは、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

追加の接続オプションを指定します。

- `AdditionalOptions` – [S3DirectSourceAdditionalOptions](#) オブジェクト。

コネクタの追加オプションを指定します。

- `OutputSchemas` – [GlueSchema](#) オブジェクトの配列。

Delta Lake ソース用のデータスキーマを指定します。

## S3CatalogDeltaSource 構造

AWS Glue データカタログに登録されている Delta Lake データソースを指定します。データソースはに保存する必要があります Amazon S3。

### フィールド

- `Name` – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

Delta Lake データソースの名前。

- `Database` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースの名前。

- `Table` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースのテーブルの名前。

- `AdditionalDeltaOptions` – キーバリュペアのマップ配列。

各キーは、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

追加の接続オプションを指定します。

- `OutputSchemas` – [GlueSchema](#) オブジェクトの配列。

Delta Lake ソース用のデータスキーマを指定します。

## CatalogDeltaSource 構造

AWS Glue データカタログに登録されている Delta Lake データソースを指定します。

## フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
Delta Lake データソースの名前。
- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースの名前。
- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースのテーブルの名前。
- AdditionalDeltaOptions – キーバリューペアのマップ配列。  
各キーは、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。  
各値は、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。  
追加の接続オプションを指定します。
- OutputSchemas – [GlueSchema](#) オブジェクトの配列。  
Delta Lake ソース用のデータスキーマを指定します。

## S3HudiSource 構造

に保存されている Hudi データソースを指定します Amazon S3。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
Hudi ソースの名前。
- Paths – 必須: UTF-8 文字列の配列。  
読み取りのソースとなる Amazon S3 パスのリスト。
- AdditionalHudiOptions – キーバリューペアのマップ配列。  
各キーは、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。  
各値は、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

追加の接続オプションを指定します。

- `AdditionalOptions` – [S3DirectSourceAdditionalOptions](#) オブジェクト。

コネクタの追加オプションを指定します。

- `OutputSchemas` – [GlueSchema](#) オブジェクトの配列。

Hudi ソース用のデータスキーマを指定します。

## S3CatalogHudiSource 構造

AWS Glue データカタログに登録されている Hudi データソースを指定します。Hudi データソースはに保存する必要があります Amazon S3。

フィールド

- `Name` – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

Hudi データソースの名前。

- `Database` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースの名前。

- `Table` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースのテーブルの名前。

- `AdditionalHudiOptions` – キーバリューペアのマップ配列。

各キーは、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

追加の接続オプションを指定します。

- `OutputSchemas` – [GlueSchema](#) オブジェクトの配列。

Hudi ソース用のデータスキーマを指定します。

## CatalogHudiSource 構造

AWS Glue データカタログに登録されている Hudi データソースを指定します。

## フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
Hudi データソースの名前。
- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースの名前。
- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースのテーブルの名前。
- AdditionalHudiOptions – キーバリューペアのマップ配列。  
各キーは、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。  
各値は、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。  
追加の接続オプションを指定します。
- OutputSchemas – [GlueSchema](#) オブジェクトの配列。  
Hudi ソース用のデータスキーマを指定します。

## DynamoDBCatalogSource 構造

データカタログ内の DynamoDB AWS Glue データソースを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
データソースの名前。
- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースの名前。
- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
読み取りデータベースのテーブルの名前。

## RelationalCatalogSource 構造

AWS Glue データカタログ内の、リレーショナルデータベースデータソースを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データソースの名前。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースの名前。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取りデータベースのテーブルの名前。

## JDBC ConnectorTarget 構造

Apache Parquet 列指向ストレージで Amazon S3 に書き込むデータターゲットを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- ConnectionName – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

コネクタに関連付けられている接続の名前。

- ConnectionTable – 必須: UTF-8 文字列。 [Custom string pattern #41](#) に一致。

データターゲットのテーブルの名前。

- ConnectorName – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

使用されるコネクタの名前。

- ConnectionType – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

JDBC データターゲットへの接続を指定する `marketplace.jdbc` や `custom.jdbc` などの接続のタイプ。

- `AdditionalOptions` – キーバリューペアのマップ配列。

各キーは、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

コネクタの追加接続オプション。

- `OutputSchemas` – [GlueSchema](#) オブジェクトの配列。

JDBC ターゲット用のデータスキーマを指定します。

## SparkConnectorTarget 構造

Apache Spark コネクタを使用するターゲットを指定します。

フィールド

- `Name` – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- `Inputs` – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- `ConnectionName` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

Apache Spark コネクタの接続の名前。

- `ConnectorName` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

Apache Spark コネクタの名前。

- `ConnectionType` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

Apache Spark データストアへの接続を指定する `marketplace.spark` や `custom.spark` などの接続のタイプ。

- `AdditionalOptions` – キーバリューペアのマップ配列。

各キーは、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

コネクタの追加接続オプション。

- OutputSchemas – [GlueSchema](#) オブジェクトの配列。

カスタム Spark ターゲット用のデータスキーマを指定します。

## BasicCatalogTarget 構造

AWS Glue データカタログテーブルを使用するターゲットを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

ターゲットとして使用するテーブルを含むデータベース。このデータベースは、データカタログに既に存在している必要があります。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

出力データのスキーマを定義するテーブル。このテーブルは、のデータカタログに既に存在している必要があります。

## MySQLCatalogTarget 構造

MySQL を使用するターゲットを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベースの名前。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベーステーブルの名前。

## PostgreSQLCatalogTarget 構造

Postgres SQL を使用するターゲットを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベースの名前。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベーステーブルの名前。

## OracleSQLCatalogTarget 構造

Oracle SQL を使用するターゲットを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベースの名前。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベーステーブルの名前。

## MicrosoftSQLServerCatalogTarget 構造

Microsoft SQL を使用するターゲットを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベースの名前。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベーステーブルの名前。

## RedshiftTarget 構造

Amazon Redshift を使用するターゲットを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベースの名前。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベーステーブルの名前。

- RedshiftTmpDir– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

データベースからコピーするときに一時データをステージングできる Amazon S3 パス。

- TmpDirIAMRole– UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

アクセス許可を持つ IAM ロール。

- UpsertRedshiftOptions – [UpsertRedshiftTargetOptions](#) オブジェクト。

Redshift ターゲットに書き込む際の upsert 処理を設定するためのオプションセット。

## AmazonRedshiftTarget 構造

Amazon Redshift ターゲットを指定します。

フィールド

- Name– UTF-8 文字列、「[Custom string pattern #43](#)」に一致。

Amazon Redshift ターゲットの名前。

- Data – [AmazonRedshiftNodeData](#) オブジェクト。

Amazon Redshift ターゲットノードのデータを指定します。

- Inputs – UTF-8 文字列の配列、1 個の文字列。

データターゲットへの入力であるノード。

## UpsertRedshiftTargetOptions 構造

Redshift ターゲットに書き込む際の upsert 処理を設定するオプション。

## フィールド

- TableLocation – UTF-8 文字列、[「Custom string pattern #40」](#) に一致。  
Redshift テーブルの物理的な場所。
- ConnectionName – UTF-8 文字列、[「Custom string pattern #40」](#) に一致。  
Redshift に書き込むために使用する接続名。
- UpsertKeys – UTF-8 文字列の配列。  
更新または挿入のどちらを実行するかを決定するためのキー。

## S3CatalogTarget 構造

Data Catalog を使用して Amazon S3 に書き込む AWS Glue データターゲットを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。[Custom string pattern #43](#) に一致。  
データターゲットの名前。
- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。  
データターゲットへの入力であるノード。
- PartitionKeys – UTF-8 文字列の配列。  
一連のキーを使用してネイティブパーティショニングを指定します。
- Table – 必須: UTF-8 文字列。[Custom string pattern #40](#) に一致。  
書き込むデータベーステーブルの名前。
- Database – 必須: UTF-8 文字列。[Custom string pattern #40](#) に一致。  
書き込むデータベースの名前。
- SchemaChangePolicy – [CatalogSchemaChangePolicy](#) オブジェクト。  
クローラの更新の動作を指定するポリシー。

## S3GlueParquetTarget 構造

Apache Parquet 列指向ストレージで Amazon S3 に書き込むデータターゲットを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- PartitionKeys – UTF-8 文字列の配列。

一連のキーを使用してネイティブパーティショニングを指定します。

- Path – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込む単一の Amazon S3 パス。

- Compression – UTF-8 文字列 (有効な値: snappy="SNAPPY" | lzo="LZO" | gzip="GZIP" | uncompressed="UNCOMPRESSED" | none="NONE")。

データの圧縮方法を指定します。データに標準のファイル拡張子が付いている場合、このオプションは一般的に不要です。指定できる値は "gzip" および "bzip" です。

- SchemaChangePolicy – [DirectSchemaChangePolicy](#) オブジェクト。

クローラの更新の動作を指定するポリシー。

## CatalogSchemaChangePolicy 構造

クローラの更新の動作を指定するポリシー。

### フィールド

- EnableUpdateCatalog – ブール。

クローラが変更されたスキーマを検出したとき、指定の更新動作を使用するかどうか。

- UpdateBehavior – UTF-8 文字列 (有効な値: UPDATE\_IN\_DATABASE | LOG)。

クローラーが変更されたスキーマを検出したときの更新動作。

## S3DirectTarget 構造

Amazon S3 に書き込むデータターゲットを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- PartitionKeys – UTF-8 文字列の配列。

一連のキーを使用してネイティブパーティショニングを指定します。

- Path – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込む単一の Amazon S3 パス。

- Compression – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

データの圧縮方法を指定します。データに標準のファイル拡張子が付いている場合、このオプションは一般的に不要です。指定できる値は "gzip" および "bzip" です。

- Format – 必須: UTF-8 文字列 (有効な値: json="JSON" | csv="CSV" | avro="AVRO" | orc="ORC" | parquet="PARQUET" | hudi="HUDI" | delta="DELTA")。

ターゲットのデータ出力形式を指定します。

- SchemaChangePolicy – [DirectSchemaChangePolicy](#) オブジェクト。

クローラの更新の動作を指定するポリシー。

## S3HudiCatalogTarget 構造

AWS Glue データカタログ内の Hudi データソースに書き込むターゲットを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。  
データターゲットへの入力であるノード。
- PartitionKeys – UTF-8 文字列の配列。  
一連のキーを使用してネイティブパーティショニングを指定します。
- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
書き込むデータベーステーブルの名前。
- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
書き込むデータベースの名前。
- AdditionalOptions – 必須: キーバリューペアのマップ配列。  
各キーは、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。  
各値は、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。  
コネクタの追加接続オプションを指定します。
- SchemaChangePolicy – [CatalogSchemaChangePolicy](#) オブジェクト。  
クローラの更新の動作を指定するポリシー。

## S3HudiDirectTarget 構造

で Hudi データソースに書き込むターゲットを指定します Amazon S3。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。  
データターゲットの名前。
- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。  
データターゲットへの入力であるノード。
- Path – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
書き込み先の Hudi データソースの Amazon S3 パス。

- **Compression** – 必須: UTF-8 文字列 (有効な値: `gzip="GZIP" | lzo="LZO" | uncompressed="UNCOMPRESSED" | snappy="SNAPPY"`)。

データの圧縮方法を指定します。データに標準のファイル拡張子が付いている場合、このオプションは一般的に不要です。指定できる値は `"gzip"` および `"bzip"` です。

- **PartitionKeys** – UTF-8 文字列の配列。

一連のキーを使用してネイティブパーティショニングを指定します。

- **Format** – 必須: UTF-8 文字列 (有効な値: `json="JSON" | csv="CSV" | avro="AVRO" | orc="ORC" | parquet="PARQUET" | hudi="HUDI" | delta="DELTA"`)。

ターゲットのデータ出力形式を指定します。

- **AdditionalOptions** – 必須: キーバリュールペアのマップ配列。

各キーは、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

コネクタの追加接続オプションを指定します。

- **SchemaChangePolicy** – [DirectSchemaChangePolicy](#) オブジェクト。

クローラの更新の動作を指定するポリシー。

## S3DeltaCatalogTarget 構造

AWS Glue データカタログ内の Delta Lake データソースに書き込むターゲットを指定します。

### フィールド

- **Name** – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- **Inputs** – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- **PartitionKeys** – UTF-8 文字列の配列。

一連のキーを使用してネイティブパーティショニングを指定します。

- **Table** – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベーステーブルの名前。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベースの名前。

- AdditionalOptions – キーバリューペアのマップ配列。

各キーは、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、 [Custom string pattern #40](#) に適合する UTF-8 文字列です。

コネクタの追加接続オプションを指定します。

- SchemaChangePolicy – [CatalogSchemaChangePolicy](#) オブジェクト。

クローラの更新の動作を指定するポリシー。

## S3DeltaDirectTarget 構造

で Delta Lake データソースに書き込むターゲットを指定します Amazon S3。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- PartitionKeys – UTF-8 文字列の配列。

一連のキーを使用してネイティブパーティショニングを指定します。

- Path – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込み先の Delta Lake データソースの Amazon S3 パス。

- Compression – 必須: UTF-8 文字列 (有効な値: uncompressed="UNCOMPRESSED" | snappy="SNAPPY")。

データの圧縮方法を指定します。データに標準のファイル拡張子が付いている場合、このオプションは一般的に不要です。指定できる値は "gzip" および "bzip" です。

- **Format** – 必須: UTF-8 文字列 (有効な値: json="JSON" | csv="CSV" | avro="AVRO" | orc="ORC" | parquet="PARQUET" | hudi="HUDI" | delta="DELTA")。

ターゲットのデータ出力形式を指定します。

- **AdditionalOptions** – キーバリューペアのマップ配列。

各キーは、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

コネクタの追加接続オプションを指定します。

- **SchemaChangePolicy** – [DirectSchemaChangePolicy](#) オブジェクト。

クローラの更新の動作を指定するポリシー。

## DirectSchemaChangePolicy 構造

クローラの更新の動作を指定するポリシー。

フィールド

- **EnableUpdateCatalog** – ブール。

クローラが変更されたスキーマを検出したとき、指定の更新動作を使用するかどうか。

- **UpdateBehavior** – UTF-8 文字列 (有効な値: UPDATE\_IN\_DATABASE | LOG)。

クローラが変更されたスキーマを検出したときの更新動作。

- **Table** – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

スキーマ変更ポリシーが適用されるデータベース内のテーブルを指定します。

- **Database** – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

スキーマ変更ポリシーを適用するデータベースを指定します。

## ApplyMapping 構造

データソースのマップデータプロパティキーを、データターゲットのデータプロパティキーに変換指定します。キーの名前を変更したり、データタイプを変更したり、データセットから削除するキーを選択できます。

## フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- Mapping – 必須: [Mapping](#) オブジェクトの配列。

データソースのデータプロパティキーを、データターゲットのデータプロパティキーにマッピングします。

## Mapping 構造

データプロパティキーのマッピングを指定します。

### フィールド

- ToKey – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

マッピングを適用した後に、列名を何にするのかを示します。FromPath と同じでも構いません。

- FromPath – UTF-8 文字列の配列。

変更するテーブルまたは列。

- FromType – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

修正されるデータのタイプ。

- ToType – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

データの修正先のデータタイプ。

- Dropped – ブール。

true の場合、列は削除されます。

- Children – [Mapping](#) オブジェクトの配列。

ネストされたデータ構造にのみ適用されます。親構造を変更し、その子構造を変更する場合は、このデータ構造に記入できます。それはまた Mapping であり、FromPath はこの構造から親の FromPath プラス FromPath です。

子部分に、次のような構造があるとします。

```
{ "FromPath": "OuterStructure", "ToKey": "OuterStructure", "ToType":
"Struct", "Dropped": false, "Children": [{ "FromPath": "inner", "ToKey":
"inner", "ToType": "Double", "Dropped": false, }] }
```

次のような Mapping を指定します。

```
{ "FromPath": "OuterStructure", "ToKey": "OuterStructure", "ToType":
"Struct", "Dropped": false, "Children": [{ "FromPath": "inner", "ToKey":
"inner", "ToType": "Double", "Dropped": false, }] }
```

## SelectFields 構造

保持するデータプロパティキーの選択変換を指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- Paths – 必須: UTF-8 文字列の配列。

データ構造内の変数への JSON パス。

## DropFields 構造

削除するデータプロパティキーを選択する変換を指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- Paths – 必須: UTF-8 文字列の配列。

データ構造内の変数への JSON パス。

## RenameField 構造

1 つのデータプロパティキーの名前を変更する変換を指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- SourcePath – 必須: UTF-8 文字列の配列。

ソースデータのデータ構造内の変数への JSON パス。

- TargetPath – 必須: UTF-8 文字列の配列。

ターゲットデータのデータ構造内の変数への JSON パス。

## スピゴット構造

Amazon S3 バケットにデータのサンプルを書き込むための変換を指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- Path – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

Amazon S3 パスでは、データセットから Amazon S3 バケットの JSON ファイルにレコードのサブセットが書き出されます。

- Topk – 数値 (integer)。100 以下。

データセットの先頭から書き込むレコードの数を指定します。

- Prob – 数値 (double)。1 以下。

特定のレコードを選ぶ確率 (最大値が 1 の 10 進値)。値 1 は、データセットから読み込まれた各行をサンプル出力に含めることを示します。

## Join 構造

指定したデータプロパティキーの比較フレーズを使用して、2 つのデータセットを 1 つに結合する変換を指定します。結合タイプは、内部結合、外部結合、左結合、右結合、左半結合、左反結合を使用できます。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。2 個の文字列。

ノード名で識別されるデータ入力。

- JoinType – 必須: UTF-8 文字列 (有効な値: equijoin="EQUIJOIN" | left="LEFT" | right="RIGHT" | outer="OUTER" | leftsemi="LEFT\_SEMI" | leftanti="LEFT\_ANTI")。

データセットで実行する結合の種類を指定します。

- Columns – 必須: [JoinColumn](#) オブジェクトの配列。2 個の構造。

結合する 2 つの列のリスト。

## JoinColumn 構造

結合する列を指定します。

フィールド

- From – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

結合する列。

- Keys – 必須: UTF-8 文字列の配列。

結合される列のキー。

## SplitFields 構造

データプロパティキーを 2 つの DynamicFrames に分割する変換を指定します。出力は DynamicFrames のコレクションです。一方は選択したデータプロパティキー、他方は残っている方のデータプロパティキーを持ちます。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- Paths – 必須: UTF-8 文字列の配列。

データ構造内の変数への JSON パス。

## SelectFromCollection 構造

DynamicFrame のコレクションから 1 つの DynamicFrames を選択するトランスフォームを指定します。出力は選択された DynamicFrame です。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- Index – 必須: 数値 (integer)。None 以下。

選択した のインデックス DynamicFrame 。

## FillMissingValues 構造

変換を使用して、データセット内に欠落値があるレコードを検索し、補完により決定する値を持つ新しいフィールドを追加します。入力データセットは、欠落値を決定する機械学習モデルのトレーニングに使用されます。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- ImputedPath – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

帰属するデータセットのデータ構造内の変数への JSON パス。

- FilledPath– UTF-8 文字列、「 [Custom string pattern #40](#) 」に一致。

データセットのデータ構造内の変数への JSON パスを入力します。

## Filter 構造

フィルター条件に基づいて、データセットを 2 つに分割する変換を指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- LogicalOperator – 必須: UTF-8 文字列 (有効な値: AND | OR)。

キー値を指定値と比較して行をフィルタリングするために使用される演算子。

- Filters – 必須: [FilterExpression](#) オブジェクトの配列。

フィルタ式を指定します。

## FilterExpression 構造

フィルタ式を指定します。

フィールド

- Operation – 必須: UTF-8 文字列 (有効な値: EQ | LT | GT | LTE | GTE | REGEX | ISNULL)。

表現で実行するオペレーションの種類。

- Negated – ブール。

その表現を無効にするかどうか。

- Values – 必須: [FilterValue](#) オブジェクトの配列。

フィルタ値のリスト。

## FilterValue 構造

FilterExpression の値リストにある単一のエントリを表します。

フィールド

- Type – 必須: UTF-8 文字列 (有効な値: COLUMNEXTRACTED | CONSTANT)。

フィルタ値のタイプ。

- Value – 必須: UTF-8 文字列の配列。

関連させる値。

## CustomCode 構造

データ変換を実行するためにカスタムコードを使用する変換を指定します。出力は のコレクションです DynamicFrames。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の以上の文字列。

ノード名で識別されるデータ入力。

- Code – 必須: UTF-8 文字列。 [Custom string pattern #35](#) に一致。

データ変換を実行するために使用されるカスタムコード。

- ClassName – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

カスタムコードノードクラスに定義された名前。

- OutputSchemas – [GlueSchema](#) オブジェクトの配列。

カスタムコード変換用のデータスキーマを指定します。

## SparkSQL 構造

データを変換する Spark SQL 構文を使用して、SQL クエリを入力する変換を指定します。出力は、単一の DynamicFrame です。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の以上の文字列。

ノード名で識別されるデータ入力。SQL クエリで使用する各入力ノードにテーブル名を関連付けることができます。選択する名前は、Spark SQL の規則を満たす必要があります。

- SqlQuery – 必須: UTF-8 文字列。 [Custom string pattern #42](#) に一致。

Spark SQL 構文を使用し、単一のデータセットを返す SQL クエリ。

- `SqlAliases` – 必須: [SqlAlias](#) オブジェクトの配列。

エイリアスのリスト。エイリアスを使用すると、特定の入力に対して SQL で使用する名前を指定できます。例えば、「」という名前のデータソースがあるとします `MyDataSource`。 `From` を `MyDataSource` に、 `Alias` に指定した場合 `SqlName`、SQL で次のことを実行できます。

```
select * from SqlName
```

からデータを取得する および `MyDataSource`。

- `OutputSchemas` – [GlueSchema](#) オブジェクトの配列。

SparkSQL 変換用のデータスキーマを指定します。

## SqlAlias 構造

`SqlAliases` の値リストにある単一のエントリを表します。

フィールド

- `From` – 必須: UTF-8 文字列。 [Custom string pattern #39](#) に一致。

テーブルまたはテーブル内の列。

- `Alias` – 必須: UTF-8 文字列。 [Custom string pattern #41](#) に一致。

テーブルまたはテーブル内の列に与えられた一時的な名前。

## DropNullFields 構造

列のすべての値が `Null` である場合に、データセットから列を削除する変換を指定します。デフォルトでは、AWS Glue Studio は `null` オブジェクトを認識しますが、空の文字列、「`null`」である文字列、`-1` 整数、またはゼロなどの他のプレースホルダーなどの一部の値は、自動的に `null` として認識されません。

フィールド

- `Name` – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- NullCheckBoxList – [NullCheckBoxList](#) オブジェクト。

特定の値が削除のために Null 値として認識されるかどうかを表す構造。

- NullTextList – [NullValueField](#) オブジェクトの配列。構造 50 個以下。

データセットに固有の Null プレースホルダーとして使用される 0 などのカスタム Null 値を表す NullValueField 構造のリストを指定する構造。

Null プレースホルダの値とデータタイプの両方がデータと一致する場合にのみ、DropNullFields 変換でカスタム NULL 値を削除します。

## NullCheckBoxList 構造

特定の値が削除の Null 値として認識されるかどうかを表します。

フィールド

- IsEmpty – ブール。

空の文字列を Null 値と見なすことを指定します。

- IsNullString – ブール。

Null の単語を綴る値を Null 値と見なすことを指定します。

- IsNegOne – ブール。

-1 の整数値が Null 値と見なすことを指定します。

## NullValueField 構造

データセットに固有の Null プレースホルダとして使用される 0 や他の値などのカスタムの Null 値を表します。

フィールド

- Value – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

Null プレースホルダの値。

- Datatype – 必須: [Datatype](#) オブジェクト。

値のデータタイプ。

## Datatype 構造

値のデータタイプを表す構造。

フィールド

- Id – 必須: UTF-8 文字列。 [Custom string pattern #39](#) に一致。

値のデータタイプ。

- Label – 必須: UTF-8 文字列。 [Custom string pattern #39](#) に一致。

データタイプに割り当てられたラベル。

## Merge 構造

レコードを識別するために、DynamicFrame プライマリキーに基づく DynamicFrame ステージングに結合変換を指定します。重複レコード ( 同じプライマリキーを持つレコード ) は重複除外されません。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。 2 個の文字列。

ノード名で識別されるデータ入力。

- Source – 必須: UTF-8 文字列。 [Custom string pattern #39](#) に一致。

DynamicFrame ステージングと結合する DynamicFrame ソース。

- PrimaryKeys – 必須: UTF-8 文字列の配列。

ソースおよびステージング動的フレームからのレコードを照合するプライマリキーフィールドのリスト。

## Union 構造

2 つ以上のデータセットの行を 1 つの結果に結合する変換を指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。2 個の文字列。

変換用のノード ID 入力。

- UnionType – 必須: UTF-8 文字列 (有効な値: ALL | DISTINCT)。

Union 変換のタイプを示します。

データソースのすべての行 ALL を結果の に結合するには、 を指定します DynamicFrame。結果として生じるユニオンでは、重複する行は削除されません。

DISTINCT 結果の で重複する行を削除するには、 を指定します DynamicFrame。

## PIIDetection 構造

PII データを識別、削除、またはマスクする変換を指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

変換用のノード ID 入力。

- PiiType – 必須: UTF-8 文字列 (有効な値: RowAudit | RowMasking | ColumnAudit | ColumnMasking)。

PIIDetection 変換のタイプを示します。

- EntityTypesToDetect – 必須: UTF-8 文字列の配列。

PIIDetection 変換が PII データとして識別するエンティティのタイプを示します。

PII タイプのエンティティには以下が含まれま

す。PERSON\_NAME、DATE、USA\_SNN、EMAIL、USA\_ITIN、USA\_PASSPORT\_NUMBER、PHONE

- OutputColumnName – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

その行で検出されたエンティティタイプを含む、すべての出力列名を示します。

- SampleFraction – 数値 (double)。1 以下。

PII エンティティのスキャン時にサンプリングするデータの割合を示します。

- ThresholdFraction – 数値 (double)。1 以下。

PII データとして識別されるために、列内で適合する必要があるデータの割合を示します。

- MaskValue – UTF-8 文字列、256 バイト長以下、[Custom string pattern #37](#) に一致。

検出されたエンティティを置き換える値を示します。

## Aggregate 構造

選択したフィールドによって行をグループ化し、指定された関数を使用して集計値を計算する変換を指定します。

フィールド

- Name – 必須: UTF-8 文字列。[Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

集約変換の入力として使用するフィールドと行を指定します。

- Groups – 必須: UTF-8 文字列の配列。

グループ化に使用するフィールドを指定します。

- Aggs – 必須: [AggregateOperation](#) オブジェクトの配列、1~30 個の構造。

指定したフィールドで実行する集計関数を指定します。

## DropDuplicates 構造

繰り返しデータの行をデータセットから削除する変換を指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

変換ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

ノード名で識別されるデータ入力。

- Columns – UTF-8 文字列の配列。

繰り返しがある場合に、それをマージまたは削除する列の名前。

## GovernedCatalogTarget 構造

Data Catalog を使用して Amazon S3 に書き込む AWS Glue データターゲットを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データターゲットの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

データターゲットへの入力であるノード。

- PartitionKeys – UTF-8 文字列の配列。

一連のキーを使用してネイティブパーティショニングを指定します。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベーステーブルの名前。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

書き込むデータベースの名前。

- SchemaChangePolicy – [CatalogSchemaChangePolicy](#) オブジェクト。

管理されたカタログを更新する際の動作を指定するポリシー。

## GovernedCatalogSource 構造

管理対象データカタログ内の AWS Glue データストアを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データストアの名前。

- Database – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み込むデータベース。

- Table – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

読み取り元のデータベーステーブル。

- PartitionPredicate – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。

この述語を満たすパーティションは削除されます。これらのパーティションの保存期間内のファイルは削除されません。"" を設定 – デフォルトでは空です。

- AdditionalOptions – [S3SourceAdditionalOptions](#) オブジェクト。

追加の接続オプションを指定します。

## AggregateOperation 構造

集約変換で集約を実行するために必要なパラメータのセットを指定します。

フィールド

- Column – 必須: UTF-8 文字列の配列。

集計関数を適用するデータセットの列を指定します。

- AggFunc – 必須: UTF-8 文字列 (有効な値: avg | countDistinct | count | first | last | kurtosis | max | min | skewness | stddev\_samp | stddev\_pop | sum | sumDistinct | var\_samp | var\_pop)。

適用する集計関数を指定します。

使用可能な集計関数には、avg

countDistinct、count、first、last、kurtosis、max、min、skewness、stddev\_samp、stddev\_pop、sum、s  
などがあります。

## GlueSchema 構造

AWS Glue でスキーマを決定できない場合に、ユーザー定義のスキーマを指定します。

フィールド

- Columns – [GlueStudioSchemaColumn](#) オブジェクトの配列。

AWS Glue スキーマを構成する列定義を指定します。

## GlueStudioSchemaColumn 構造

AWS Glue スキーマ定義の 1 つの列を指定します。

フィールド

- Name – 必須: バイト長が 1024 以下で [Single-line string pattern](#) に適合する、UTF-8 文字列。

AWS Glue Studio スキーマの列の名前。

- Type - UTF-8 文字列。131,072 バイト長以下。 [Single-line string pattern](#) に一致。

AWS Glue Studio スキーマ内のこの列の hive タイプ。

## GlueStudioColumn 構造

AWS Glue Studio で 1 つの列を指定します。

フィールド

- Key – 必須: UTF-8 文字列。 [Custom string pattern #41](#) に一致。

AWS Glue Studio の列のキー。

- FullPath – 必須: UTF-8 文字列の配列。

AWS Glue Studio TThe の列の完全な URL。

- Type – 必須: UTF-8 文字列 (有効な値: array="ARRAY" | bigint="BIGINT" | bigint array="BIGINT\_ARRAY" | binary="BINARY" | binary array="BINARY\_ARRAY" | boolean="BOOLEAN" | boolean array="BOOLEAN\_ARRAY" | byte="BYTE" | byte array="BYTE\_ARRAY" | char="CHAR" | char array="CHAR\_ARRAY" | choice="CHOICE" | choice array="CHOICE\_ARRAY" | date="DATE" | date array="DATE\_ARRAY" | decimal="DECIMAL" | decimal array="DECIMAL\_ARRAY" | double="DOUBLE" | double array="DOUBLE\_ARRAY" | enum="ENUM" | enum array="ENUM\_ARRAY" | float="FLOAT" | float array="FLOAT\_ARRAY" | int="INT" | int array="INT\_ARRAY" | interval="INTERVAL" | interval array="INTERVAL\_ARRAY" | long="LONG" | long array="LONG\_ARRAY" | object="OBJECT" | short="SHORT" | short array="SHORT\_ARRAY" | smallint="SMALLINT" | smallint array="SMALLINT\_ARRAY" | string="STRING" | string array="STRING\_ARRAY" | timestamp="TIMESTAMP" | timestamp array="TIMESTAMP\_ARRAY" | tinyint="TINYINT" | tinyint array="TINYINT\_ARRAY" | varchar="VARCHAR" | varchar array="VARCHAR\_ARRAY" | null="NULL" | unknown="UNKNOWN" | unknown array="UNKNOWN\_ARRAY")。

AWS Glue Studio TThe の列のタイプ。

- Children – 構造の配列。

AWS Glue Studio TThe の親列の子。

## DynamicTransform 構造

動的変換を実行するために必要なパラメータのセットを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

動的変換の名前を指定します。

- TransformName – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

AWS Glue Studio ビジュアルエディタに表示される動的変換の名前を指定します。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。  
必要な動的変換の入力を指定します。
- Parameters – [TransformConfigParameter](#) オブジェクトの配列。  
動的変換のパラメータを指定します。
- FunctionName – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
動的変換の関数の名前を指定します。
- Path – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
動的変換ソースファイルおよび設定ファイルのパスを指定します。
- Version – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。  
このフィールドは使用されず、将来のリリースで非推奨となります。
- OutputSchemas – [GlueSchema](#) オブジェクトの配列。  
動的変換用のデータスキーマを指定します。

## TransformConfigParameter 構造

動的変換の設定ファイル内のパラメータを指定します。

### フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。  
動的変換の設定ファイル内のパラメータの名前を指定します。
- Type – 必須: UTF-8 文字列 (有効な値: str="STR" | int="INT" | float="FLOAT" | complex="COMPLEX" | bool="BOOL" | list="LIST" | null="NULL")。  
動的変換の設定ファイル内のパラメータタイプを指定します。
- ValidationRule – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。  
動的変換の設定ファイル内の検証ルールを指定します。
- ValidationMessage – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。  
動的変換の設定ファイル内の検証メッセージを指定します。
- Value – UTF-8 文字列の配列。

動的変換の設定ファイル内のパラメータの値を指定します。

- `ListType` – UTF-8 文字列 (有効な値: `str="STR" | int="INT" | float="FLOAT" | complex="COMPLEX" | bool="BOOL" | list="LIST" | null="NULL"`)。

動的変換の設定ファイル内のパラメータのリスト型を指定します。

- `IsOptional` – ブール。

動的変換の設定ファイル内のパラメータがオプションかどうかを指定します。

## EvaluateDataQuality 構造

データ品質評価基準を指定します。

フィールド

- `Name` – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

データ品質評価の名前。

- `Inputs` – 必須: UTF-8 文字列の配列。1 個の文字列。

データ品質評価の入力。

- `Ruleset` – 必須: UTF-8 文字列、1 ~ 65536 バイト長、 [Custom string pattern #38](#) に一致。

データ品質評価のルールセット。

- `Output` – UTF-8 文字列 (有効な値: `PrimaryInput | EvaluationResults`)。

データ品質評価の出力。

- `PublishingOptions` – [DQResultsPublishingOptions](#) オブジェクト。

結果の発行方法を設定するオプション。

- `StopJobOnFailureOptions` – [DQStopJobOnFailureOptions](#) オブジェクト。

データ品質評価に失敗した場合にジョブを停止する方法を設定するオプション。

## DQ ResultsPublishingOptions 構造

データ品質評価の結果の発行方法を設定するオプション。

## フィールド

- `EvaluationContext` – UTF-8 文字列、「[Custom string pattern #39](#)」に一致。  
評価のコンテキスト。
- `ResultsS3Prefix` – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。  
結果に追加された Amazon S3 プレフィックス。
- `CloudWatchMetricsEnabled` – ブール。  
データ品質結果のメトリクスを有効にします。
- `ResultsPublishingEnabled` – ブール。  
データ品質結果の発行を有効にします。

## DQ StopJobOnFailureOptions 構造

データ品質評価に失敗した場合にジョブを停止する方法を設定するオプション。

### フィールド

- `StopJobOnFailureTiming` – UTF-8 文字列 (有効な値: `Immediate` | `AfterDataLoad`)。  
データ品質評価が失敗した場合にジョブを停止するタイミング。オプションは即時または `AfterDataLoad`。

## EvaluateDataQualityMultiFrame 構造

データ品質評価基準を指定します。

### フィールド

- `Name` – 必須: UTF-8 文字列。[Custom string pattern #43](#) に一致。  
データ品質評価の名前。
- `Inputs` – 必須: UTF-8 文字列の配列。1 個の以上の文字列。  
データ品質評価の入力。このリストにおける最初の入力はプライマリデータソースです。
- `AdditionalDataSources` – キーバリュペアのマップ配列。

各キーは、[Custom string pattern #43](#) に適合する UTF-8 文字列です。

各値は、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

プライマリを除くすべてのデータソースのエイリアス。

- Ruleset – 必須: UTF-8 文字列、1~65536 バイト長、[Custom string pattern #38](#) に一致。

データ品質評価のルールセット。

- PublishingOptions – [DQResultsPublishingOptions](#) オブジェクト。

結果の発行方法を設定するオプション。

- AdditionalOptions – キーバリュールペアのマッピング配列。

各キーは UTF-8 文字列 (有効な値: performanceTuning.caching="CacheOption" | observations.scope="ObservationsOption")。

各値は UTF-8 文字列。

変換のランタイム動作を設定するオプション。

- StopJobOnFailureOptions – [DQStopJobOnFailureOptions](#) オブジェクト。

データ品質評価に失敗した場合にジョブを停止する方法を設定するオプション。

## Recipe 構造

AWS Glue ジョブで AWS Glue DataBrew recipe を使用する AWS Glue Studio ノード。

フィールド

- Name – 必須: UTF-8 文字列。[Custom string pattern #43](#) に一致。

AWS Glue Studio ノードの名前。

- Inputs – 必須: UTF-8 文字列の配列。1 個の文字列。

レシピノードへの入力となるノード。これは ID によって識別されます。

- RecipeReference – 必須: [RecipeReference](#) オブジェクト。

ノードで使用される DataBrew レシピへの参照。

## RecipeReference 構造

AWS Glue DataBrew レシピへの参照。

フィールド

- `RecipeArn` – 必須: UTF-8 文字列。[Custom string pattern #40](#) に一致。  
recipe の ARN DataBrew 。
- `RecipeVersion` – 必須: UTF-8 文字列、1 ~ 16 バイト長。  
DataBrew レシピ RecipeVersion の 。

## SnowflakeNodeData 構造

AWS Glue Studio の Snowflake ノードの設定を指定します。

フィールド

- `SourceType` – UTF-8 文字列、「[Custom string pattern #39](#)」に一致。  
取得したデータの指定方法を指定します。有効な値: "table"、 "query"。
- `Connection` – [オプション](#) オブジェクト。  
Snowflake エンドポイントへの AWS Glue データカタログ接続を指定します。
- `Schema` – UTF-8 文字列。  
ノードが使用する Snowflake データベーススキーマを指定します。
- `Table` – UTF-8 文字列。  
ノードが使用する Snowflake テーブルを指定します。
- `Database` – UTF-8 文字列。  
ノードが使用する Snowflake データベースを指定します。
- `TempDir` – UTF-8 文字列、「[Custom string pattern #40](#)」に一致。  
現在使用されていません。
- `IamRole` – [オプション](#) オブジェクト。  
現在使用されていません。

- `AdditionalOptions` – キーバリューペアのマップ配列。

各キーは、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

各値は、[Custom string pattern #40](#) に適合する UTF-8 文字列です。

Snowflake コネクタに渡される追加オプションを指定します。ノードの他の場所でオプションが指定されている場合、こちらが優先されます。

- `SampleQuery` – UTF-8 文字列。

query `SourceType` でデータを取得するために使用する SQL 文字列。

- `PreAction` – UTF-8 文字列。

Snowflake コネクタが標準アクションを実行する前に実行される SQL 文字列。

- `PostAction` – UTF-8 文字列。

Snowflake コネクタが標準アクションを実行した後に実行される SQL 文字列。

- `Action` – UTF-8 文字列。

既存のデータを持つテーブルに書き込むときに実行するアクションを指定します。有効な値は、`append`、`merge`、`truncate`、`drop` です。

- `Upsert` – ブール。

アクションが `append` の場合に使用します。行が既に存在する場合の解決動作を指定します。true の場合、既存の行が更新されます。false の場合、それらの行が挿入されます。

- `MergeAction` – UTF-8 文字列、「[Custom string pattern #39](#)」に一致。

マージアクションを指定します。有効な値: `simple`、`custom`。simple の場合、マージ動作は `MergeWhenMatched` と `MergeWhenNotMatched` によって定義されます。custom の場合、`MergeClause` によって定義されます。

- `MergeWhenMatched` – UTF-8 文字列、「[Custom string pattern #39](#)」に一致。

マージ時に既存のデータと一致するレコードを解決する方法を指定します。有効な値: `update`、`delete`。

- `MergeWhenNotMatched` – UTF-8 文字列、「[Custom string pattern #39](#)」に一致。

マージ時に既存のデータと一致しないレコードを処理する方法を指定します。有効な値: `insert`、`none`。

- MergeClause – UTF-8 文字列。

カスタムマージ動作を指定する SQL ステートメント。

- StagingTable – UTF-8 文字列。

merge または upsert を行う append アクションを実行するときに使用されるステージングテーブルの名前。データはこのテーブルに書き込まれ、生成されたポストアクションによって table に移動されます。

- SelectedColumns – [オプション](#) オブジェクトの配列。

マージや upsert の一致を検出するときに、レコードを識別するために組み合わせる列を指定します。value、label、description キーを使用する構造のリストです。各構造は列を記述します。

- AutoPushdown – ブール。

自動クエリプッシュダウンを有効にするかどうかを指定します。プッシュダウンが有効になっている場合、Spark でクエリを実行すると、クエリの一部が Snowflake サーバーに「プッシュダウン」できる場合にクエリがプッシュダウンされます。これにより、一部のクエリのパフォーマンスが向上します。

- TableSchema – [オプション](#) オブジェクトの配列。

ノードのターゲットスキーマを手動で定義します。value、label、description キーを使用する構造のリストです。各構造は列を定義します。

## SnowflakeSource 構造

Snowflake データソースを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

Snowflake データソースの名前。

- Data – 必須: [SnowflakeNodeData](#) オブジェクト。

Snowflake データソースの設定。

- OutputSchemas – [GlueSchema](#) オブジェクトの配列。

出力データのユーザー定義スキーマを指定します。

## SnowflakeTarget 構造

Snowflake ターゲットを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

Snowflake ターゲットの名前。

- Data – 必須: [SnowflakeNodeData](#) オブジェクト。

Snowflake ターゲットノードのデータを指定します。

- Inputs – UTF-8 文字列の配列、1 個の文字列。

データターゲットへの入力であるノード。

## ConnectorDataSource 構造

標準の接続オプションを使用して生成されたソースを指定します。

フィールド

- Name – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

このソースノードの名前。

- ConnectionType – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

基盤となる AWS Glue ライブラリ connectionType に提供される。ノードタイプは、次の接続タイプをサポートします。

- opensearch
- azuresql
- azurecosmos
- bigquery
- saphana
- teradata

- `vertica`
- `Data` – 必須: キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

ノードの接続オプションを指定するマップ。対応する接続タイプの標準接続オプションは、AWS Glue ドキュメントの「[接続パラメータ](#)」セクションにあります。

- `OutputSchemas` – [GlueSchema](#) オブジェクトの配列。

このソース用のデータスキーマを指定します。

## ConnectorDataTarget 構造

標準の接続オプションを使用して生成されたターゲットを指定します。

フィールド

- `Name` – 必須: UTF-8 文字列。 [Custom string pattern #43](#) に一致。

このターゲットノードの名前。

- `ConnectionType` – 必須: UTF-8 文字列。 [Custom string pattern #40](#) に一致。

基盤となる AWS Glue ライブラリ `connectionType` に提供される。ノードタイプは、次の接続タイプをサポートします。

- `opensearch`
- `azuresql`
- `azurecosmos`
- `bigquery`
- `saphana`
- `teradata`
- `vertica`
- `Data` – 必須: キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

ノードの接続オプションを指定するマップ。対応する接続タイプの標準接続オプションは、AWS Glue ドキュメントの「[接続パラメータ](#)」セクションにあります。

- Inputs – UTF-8 文字列の配列、1 個の文字列。

データターゲットへの入力であるノード。

## ジョブ API

Jobs API では、ジョブのデータ型について説明しており、これには、AWS Glue でジョブ、ジョブ実行、およびトリガーを操作するための API が含まれます。

トピック

- [ジョブ](#)
- [ジョブ実行](#)
- [トリガー](#)

## ジョブ

Jobs API は、でのジョブの作成、更新、削除、表示に関連するデータ型と API について説明します AWS Glue。

データ型

- [Job 構造](#)
- [ExecutionProperty 構造](#)
- [NotificationProperty 構造](#)
- [JobCommand 構造](#)
- [ConnectionsList 構造](#)
- [JobUpdate 構造](#)
- [SourceControlDetails 構造](#)

## Job 構造

ジョブ定義を指定します。

### フィールド

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このジョブ定義に割り当てる名前。

- JobMode – UTF-8 文字列 (有効な値: SCRIPT="" | VISUAL="" | NOTEBOOK="")。

ジョブの作成方法を説明するモード。有効な値は次のとおりです。

- SCRIPT - ジョブは AWS Glue Studio スクリプトエディタを使用して作成されました。
- VISUAL - ジョブは AWS Glue Studio ビジュアルエディタを使用して作成されました。
- NOTEBOOK - ジョブはインタラクティブセッションノートブックを使用して作成されました。

JobMode フィールドが欠落しているか null の場合、SCRIPT がデフォルト値として割り当てられます。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

ジョブの説明。

- LogUri – UTF-8 文字列。

このフィールドは、将来の利用のために予約されています。

- Role – UTF-8 文字列。

このジョブに関連付けられている IAM ロールの名前または Amazon リソースネーム (ARN)。

- CreatedOn – タイムスタンプ。

このジョブ定義を作成した日時。

- LastModifiedOn – タイムスタンプ。

このジョブ定義を変更した最後の時点。

- ExecutionProperty – [ExecutionProperty](#) オブジェクト。

このジョブに許可される同時実行の最大数を指定する ExecutionProperty。

- Command – [JobCommand](#) オブジェクト。

このジョブを実行する JobCommand。

- `DefaultArguments` – キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

名前と値のペアとして指定された、このジョブを実行することのデフォルトの引数。

ここで独自のジョブ実行スクリプトが消費する引数と、AWS Glue それ自体が消費する引数を指定できます。

ジョブ引数はログに記録される場合があります。プレーンテキストのシークレットを引数として渡さないでください。ジョブ内に保持する場合は、AWS Glue 接続 AWS Secrets Manager、またはその他のシークレット管理メカニズムからシークレットを取得します。

独自のジョブ引数を指定および使用方法については、デベロッパーガイドのトピック「[Calling AWS Glue APIs in Python](#)」を参照してください。

Spark ジョブの設定時にこのフィールドに指定できる引数については、デベロッパーガイドのトピック「[Special Parameters Used by AWS Glue](#)」を参照してください。

Ray ジョブの設定時にこのフィールドに指定できる引数については、デベロッパーガイドの「[Ray ジョブでジョブパラメータを使用する](#)」を参照してください。

- `NonOverridableArguments` – キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

名前と値のペアとして指定された、ジョブ実行時にジョブ引数を指定しても上書きされないこのジョブの引数。

- `Connections` – [ConnectionsList](#) オブジェクト。

このジョブに使用される接続。

- `MaxRetries` – 数値 (整数)。

が JobRun 失敗した後にこのジョブを再試行する最大回数。

- `AllocatedCapacity` – 数値 (整数)。

~~このフィールドは廃止されました。代わりに `MaxCapacity` を使用します。~~

このジョブの実行に割り当てられた AWS Glue データ処理ユニット (DPU の数。最低 2 つの DPU を割り当てることができます。デフォルトは 10 です。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

- Timeout - 数値 (整数)。1 以上。

ジョブのタイムアウト (分)。ジョブ実行が終了済みになって TIMEOUT ステータスに入るまでに、ジョブ実行でリソースを消費できる最大時間です。バッチジョブのデフォルト値は 2,880 分 (48 時間) です。

ストリーミングジョブのタイムアウト値は 7 日または 10080 分未満である必要があります。値を空白のままにすると、メンテナンスウィンドウを設定していない場合、ジョブは 7 日後に再起動されます。メンテナンスウィンドウを設定している場合、7 日後のメンテナンスウィンドウ中に再起動されます。

- MaxCapacity - 数値 (double)。

Glue バージョン 1.0 以前のジョブの場合、標準ワーカータイプを使用して、このジョブの実行時に割り当てることができる AWS Glue データ処理ユニット (DPU の数。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

一方、Glue バージョン 2.0 以降のジョブでは、Maximum capacity を指定できません。代わりに、Worker type と Number of workers を指定する必要があります。

WorkerType および NumberOfWorkers を使用している場合は MaxCapacity を設定しないでください。

MaxCapacity に割り当てることができる値は、Python シェルジョブ、Apache Spark ETL ジョブ、Apache Spark ストリーミング ETL ジョブのいずれを実行しているかによって異なります。

- Python シェルジョブを指定すると (JobCommand.Name="pythonshell")、0.0625 または 1 DPU のいずれかを割り当てることができます。デフォルトは 0.0625 DPU です。
- Apache Spark ETL ジョブ (JobCommand.Name="glueetl") または Apache Spark ストリーミング ETL ジョブ (JobCommand.Name="gluestreaming") を指定した場合は、2 ~ 100 の DPU を割り当てることができます。デフォルトでは 10 DPU になっています。このジョブタイプには、小数の DPU 割り当てを指定できません。

- WorkerType – UTF-8 文字列 (有効な値: Standard="" | G.1X="" | G.2X="" | G.025X="" | G.4X="" | G.8X="" | Z.2X="")。

ジョブの実行時に割り当てられる事前定義済みのワーカーの種類。Spark ジョブに使用できる値は G.1X、G.2X、G.4X、G.8X、または G.025X です。Ray ジョブに使用できる値は Z.2X です。

- G.1X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 1 DPU (4 vCPU、16 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- G.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 77 GB の空き容量) を備えた 2 DPU (8 vCPU、32 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- G.4X ワーカータイプでは、各ワーカーは 256 GB のディスク (約 235 GB の空き容量) を備えた 4 DPU (16 vCPU、64 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、米国東部 (オハイオ)、米国東部 (バージニア北部)、米国西部 (オレゴン)、アジアパシフィック (シンガポール)、アジアパシフィック (シドニー)、アジアパシフィック (東京)、カナダ (中部)、欧州 (フランクフルト)、欧州 (アイルランド)、欧州 (ストックホルム) の各 AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL ジョブでのみ使用できます。
- G.8X ワーカータイプでは、各ワーカーは 512 GB のディスク (約 487 GB の空き容量) を備えた 8 DPU (32 vCPU、128 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、ワーカータイプでサポートされているのと同じ AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL G.4X ジョブでのみ使用できます。
- G.025X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 0.25 DPU (2 vCPU、4 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。少量のストリーミングジョブには、このワーカータイプをお勧めします。このワーカータイプは、AWS Glue バージョン 3.0 ストリーミングジョブでのみ使用できます。

- Z.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 120 GB の空き容量) を備えた 2 M-DPU (8 vCPU、64 GB のメモリ) にマッピングされており、オートスケーラーに基づき最大 8 個の Ray ワーカーを提供します。
- NumberOfWorkers – 数値 (整数)。

ジョブの実行時に割り当てられた、定義済みの workerType ワーカー数。

- SecurityConfiguration – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このジョブで使用される SecurityConfiguration 構造の名前。

- NotificationProperty – [NotificationProperty](#) オブジェクト。

ジョブ通知の設定プロパティを指定します。

- Running – ブール。

このフィールドは、将来の利用のために予約されています。

- GlueVersion – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

Spark ジョブでは、はジョブ AWS Glue で使用できる Apache Spark と Python のバージョン GlueVersion を決定します。Python バージョンは、Spark タイプのジョブでサポートされるバージョンを示します。

Ray ジョブの場合、GlueVersion を 4.0 以降に設定する必要があります。ただし、Ray ジョブで使用できる Ray、Python、および追加ライブラリのバージョンは、Job コマンドの Runtime パラメータによって決まります。

利用可能な AWS Glue バージョンと対応する Spark および Python バージョンの詳細については、デベロッパーガイドの「[Glue バージョン](#)」を参照してください。

Glue バージョンを指定せずに作成されたジョブは、デフォルトで Glue 0.9 に設定されます。

- CodeGenConfigurationNodes – キーバリューペアのマップ配列。

各キーは、[Custom string pattern #39](#) に適合する UTF-8 文字列です。

各値は [CodeGenConfigurationNode](#) オブジェクトです。

Glue Studio ビジュアルコンポーネントと、Glue Studio によるコード生成の両方がベースとする、有向非巡回グラフの表現。

- `ExecutionClass` - UTF-8 文字列。16 バイト長以下 (有効値: FLEX="" | STANDARD="")。

ジョブが標準実行クラスまたは柔軟な実行クラスのどちらで実行されるのかを示します。標準実行クラスは、素早くジョブを起動する必要があり、専用のリソースが必要な時間的な制約のあるワークロードに最適です。

柔軟な実行クラスは、開始時刻と完了時刻が変化する時間的な制約が厳しいジョブに適しています。

AWS Glue バージョン 3.0 以降およびコマンドタイプのジョブのみを `ExecutionClass glueetl` に設定できます FLEX。柔軟な実行クラスは Spark ジョブで使用できます。

- `SourceControlDetails` - [SourceControlDetails](#) オブジェクト。

ジョブのソース管理設定の詳細。これにより、リモートリポジトリとの間でジョブアーティファクトを同期できます。

- `MaintenanceWindow` - UTF-8 文字列、「[Custom string pattern #30](#)」に一致。

このフィールドは、ストリーミングジョブのメンテナンスウィンドウの曜日と時間を指定します。AWS Glue は定期的にメンテナンスアクティビティを実行します。これらのメンテナンスウィンドウ中に、AWS Glue はストリーミングジョブを再起動する必要があります。

AWS Glue は、指定されたメンテナンスウィンドウから 3 時間以内にジョブを再起動します。たとえば、月曜日の午前 10 時 (GMT) にメンテナンスウィンドウを設定すると、ジョブは午前 10 時 (GMT) から午後 1 時 (GMT) までの間に再起動されます。

- `ProfileName` - UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ジョブに関連付けられた AWS Glue 使用プロファイルの名前。

## ExecutionProperty 構造

ジョブの実行プロパティ。

### フィールド

- `MaxConcurrentRuns` - 数値 (整数)。

このジョブで許可される同時実行の最大数。デフォルトは 1 です。このしきい値に達すると、エラーが返されます。指定できる最大値は、サービスの制限によってコントロールされます。

## NotificationProperty 構造

通知の構成プロパティを指定します。

フィールド

- NotifyDelayAfter - 数値 (整数)。1 以上。

ジョブの実行が開始された後、ジョブ実行遅延通知を送信するまでの待機時間 (分単位)。

## JobCommand 構造

ジョブの実行時に実行するコードを指定します。

フィールド

- Name – UTF-8 文字列。

ジョブコマンドの名前。Apache Spark ETL ジョブの場合は、`glueetl` を指定する必要があります。Python シェルジョブの場合は、`pythonshell` を指定する必要があります。Apache Spark ストリーミング ETL ジョブの場合は、`gluestreaming` を指定する必要があります。Ray ジョブの場合は、`glueray` を指定する必要があります。

- ScriptLocation - UTF-8 文字列。400000 バイト長以下。

ジョブを実行するスクリプトへの Amazon Simple Storage Service (Amazon S3) パスを指定します。

- PythonVersion– UTF-8 文字列、[「Custom string pattern #21」](#) に一致。

Python シェルジョブを実行するために使用中の Python のバージョン。指定できる値は、2 または 3 です。

- Runtime - UTF-8 文字列、64 バイト長以下、[Custom string pattern #29](#) に一致。

Ray ジョブではランタイムを使用して、環境で使用可能な Ray、Python、および他のライブラリのバージョンを指定します。このフィールドは、他のジョブタイプでは使用されません。サポートされているランタイム環境値については、「[AWS Glue デベロッパーガイド](#)」の「[サポートされている Ray ランタイム環境](#)」を参照してください。

## ConnectionsList 構造

ジョブが使用する接続を指定します。

フィールド

- Connections – UTF-8 文字列の配列。

ジョブが使用する接続のリスト。

## JobUpdate 構造

既存のジョブ定義を更新するための情報を指定します。以前のジョブ定義はこの情報によって完全に上書きされます。

フィールド

- JobMode – UTF-8 文字列 (有効な値: SCRIPT="" | VISUAL="" | NOTEBOOK="")。

ジョブの作成方法を説明するモード。有効な値は次のとおりです。

- SCRIPT - ジョブは AWS Glue Studio スクリプトエディタを使用して作成されました。
- VISUAL - ジョブは AWS Glue Studio ビジュアルエディタを使用して作成されました。
- NOTEBOOK - ジョブはインタラクティブセッションノートブックを使用して作成されました。

JobMode フィールドが欠落しているか null の場合、SCRIPT がデフォルト値として割り当てられます。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

定義するジョブの説明。

- LogUri – UTF-8 文字列。

このフィールドは、将来の利用のために予約されています。

- Role – UTF-8 文字列。

このジョブに関連付けられている IAM ロールの名前または Amazon リソースネーム (ARN) (必須)。

- ExecutionProperty – [ExecutionProperty](#) オブジェクト。

このジョブに許可される同時実行の最大数を指定する ExecutionProperty。

- Command – [JobCommand](#) オブジェクト。

このジョブを実行する JobCommand (必須)。

- DefaultArguments – キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

名前と値のペアとして指定された、このジョブを実行するごとのデフォルトの引数。

ここで独自のジョブ実行スクリプトが消費する引数と、AWS Glue それ自体が消費する引数を指定できます。

ジョブ引数はログに記録される場合があります。プレーンテキストのシークレットを引数として渡さないでください。ジョブ内に保持する場合は、AWS Glue 接続 AWS Secrets Manager、またはその他のシークレット管理メカニズムからシークレットを取得します。

独自のジョブ引数を指定および使用する方法については、デベロッパーガイドのトピック「[Calling AWS Glue APIs in Python](#)」を参照してください。

Spark ジョブの設定時にこのフィールドに指定できる引数については、デベロッパーガイドのトピック「[Special Parameters Used by AWS Glue](#)」を参照してください。

Ray ジョブの設定時にこのフィールドに指定できる引数については、デベロッパーガイドの「[Ray ジョブでジョブパラメータを使用する](#)」を参照してください。

- NonOverridableArguments – キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

名前と値のペアとして指定された、ジョブ実行時にジョブ引数を指定しても上書きされないこのジョブの引数。

- Connections – [ConnectionsList](#) オブジェクト。

このジョブに使用される接続。

- MaxRetries – 数値 (整数)。

失敗した場合にこのジョブを再試行する最大回数。

- `AllocatedCapacity` – 数値 (整数)。

このフィールドは廃止されました。代わりに `MaxCapacity` を使用します。

このジョブに割り当てる AWS Glue データ処理ユニット (DPUsの数。最低 2 つの DPUを割り当てることができます。デフォルトは 10 です。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

- `Timeout` - 数値 (整数)。1 以上。

ジョブのタイムアウト (分)。ジョブ実行が終了済みになって `TIMEOUT` ステータスに入るまでに、ジョブ実行でリソースを消費できる最大時間です。バッチジョブのデフォルト値は 2,880 分 (48 時間) です。

ストリーミングジョブのタイムアウト値は 7 日または 10080 分未満である必要があります。値を空白のままにすると、メンテナンスウィンドウを設定していない場合、ジョブは 7 日後に再起動されます。メンテナンスウィンドウを設定している場合、7 日後のメンテナンスウィンドウ中に再起動されます。

- `MaxCapacity` – 数値 (double)。

Glue バージョン 1.0 以前のジョブの場合、標準ワーカータイプを使用して、このジョブの実行時に割り当てることのできる AWS Glue データ処理ユニット (DPUsの数。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

一方、Glue バージョン 2.0 以降のジョブでは、`Maximum capacity` を指定できません。代わりに、`Worker type` と `Number of workers` を指定する必要があります。

`WorkerType` および `NumberOfWorkers` を使用している場合は `MaxCapacity` を設定しないでください。

`MaxCapacity` に割り当てることのできる値は、Python シェルジョブ、Apache Spark ETL ジョブ、Apache Spark ストリーミング ETL ジョブのいずれを実行しているかによって異なります。

- Python シェルジョブを指定すると (`JobCommand.Name="pythonshell"`)、0.0625 または 1 DPU のいずれかを割り当てることができます。デフォルトは 0.0625 DPU です。
- Apache Spark ETL ジョブ (`JobCommand.Name="glueetl"`) または Apache Spark ストリーミング ETL ジョブ (`JobCommand.Name="gluestreaming"`) を指定した場合は、2 ~ 100 の DPU を割り

当てることができます。デフォルトでは 10 DPU になっています。このジョブタイプには、小数の DPU 割り当てを指定できません。

- `WorkerType` – UTF-8 文字列 (有効な値: `Standard=""` | `G.1X=""` | `G.2X=""` | `G.025X=""` | `G.4X=""` | `G.8X=""` | `Z.2X=""`)。

ジョブの実行時に割り当てられる事前定義済みのワーカーの種類。Spark ジョブに使用できる値は G.1X、G.2X、G.4X、G.8X、または G.025X です。Ray ジョブに使用できる値は Z.2X です。

- G.1X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 1 DPU (4 vCPU、16 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- G.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 77 GB の空き容量) を備えた 2 DPU (8 vCPU、32 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- G.4X ワーカータイプでは、各ワーカーは 256 GB のディスク (約 235 GB の空き容量) を備えた 4 DPU (16 vCPU、64 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、米国東部 (オハイオ)、米国東部 (バージニア北部)、米国西部 (オレゴン)、アジアパシフィック (シンガポール)、アジアパシフィック (シドニー)、アジアパシフィック (東京)、カナダ (中部)、欧州 (フランクフルト)、欧州 (アイルランド)、欧州 (ストックホルム) の各 AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL ジョブでのみ使用できます。
- G.8X ワーカータイプでは、各ワーカーは 512 GB のディスク (約 487 GB の空き容量) を備えた 8 DPU (32 vCPU、128 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、ワーカータイプでサポートされているのと同じ AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL G.4X ジョブでのみ使用できます。
- G.025X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 0.25 DPU (2 vCPU、4 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。少量のストリーミングジョブには、このワーカータイプをお勧めします。このワーカータイプは、AWS Glue バージョン 3.0 ストリーミングジョブでのみ使用できます。

- Z.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 120 GB の空き容量) を備えた 2 M-DPU (8 vCPU、64 GB のメモリ) にマッピングされており、オートスケーラーに基づき最大 8 個の Ray ワーカーを提供します。
- NumberOfWorkers – 数値 (整数)。

ジョブの実行時に割り当てられた、定義済みの workerType ワーカー数。

- SecurityConfiguration – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このジョブで使用される SecurityConfiguration 構造の名前。

- NotificationProperty – [NotificationProperty](#) オブジェクト。

ジョブ通知の設定プロパティを指定します。

- GlueVersion – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

Spark ジョブでは、はジョブ AWS Glue で使用できる Apache Spark と Python のバージョン GlueVersion を決定します。Python バージョンは、Spark タイプのジョブでサポートされるバージョンを示します。

Ray ジョブの場合、GlueVersion を 4.0 以降に設定する必要があります。ただし、Ray ジョブで使用できる Ray、Python、および追加ライブラリのバージョンは、Job コマンドの Runtime パラメータによって決まります。

利用可能な AWS Glue バージョンと対応する Spark および Python バージョンの詳細については、デベロッパーガイドの「[Glue バージョン](#)」を参照してください。

Glue バージョンを指定せずに作成されたジョブは、デフォルトで Glue 0.9 に設定されます。

- CodeGenConfigurationNodes – キーバリューペアのマップ配列。

各キーは、[Custom string pattern #39](#) に適合する UTF-8 文字列です。

各値は [CodeGenConfigurationNode](#) オブジェクトです。

Glue Studio ビジュアルコンポーネントと、Glue Studio によるコード生成の両方がベースとする、有向非巡回グラフの表現。

- ExecutionClass - UTF-8 文字列。16 バイト長以下 (有効値: FLEX="" | STANDARD="")。

ジョブが標準実行クラスまたは柔軟な実行クラスのどちらで実行されるのかを示します。標準の実行クラスは、素早くジョブを起動する必要があり、専用のリソースが必要な時間的な制約のあるワークロードに最適です。

柔軟な実行クラスは、開始時刻と完了時刻が変化する時間的な制約が厳しいジョブに適しています。

AWS Glue バージョン 3.0 以降およびコマンドタイプのジョブのみを ExecutionClass glueetl に設定できます FLEX。柔軟な実行クラスは Spark ジョブで使用できます。

- SourceControlDetails – [SourceControlDetails](#) オブジェクト。

ジョブのソース管理設定の詳細。これにより、リモートリポジトリとの間でジョブアーティファクトを同期できます。

- MaintenanceWindow – UTF-8 文字列、「[Custom string pattern #30](#)」に一致。

このフィールドは、ストリーミングジョブのメンテナンスウィンドウの曜日と時間を指定します。AWS Glue は定期的にメンテナンスアクティビティを実行します。これらのメンテナンスウィンドウ中に、AWS Glue はストリーミングジョブを再起動する必要があります。

AWS Glue は、指定されたメンテナンスウィンドウから 3 時間以内にジョブを再起動します。たとえば、月曜日の午前 10 時 (GMT) にメンテナンスウィンドウを設定すると、ジョブは午前 10 時 (GMT) から午後 1 時 (GMT) までの間に再起動されます。

- ProfileName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ジョブに関連付けられた AWS Glue 使用プロファイルの名前。

## SourceControlDetails 構造

ジョブのソース管理設定の詳細。これにより、リモートリポジトリとの間でジョブアーティファクトを同期できます。

### フィールド

- Provider – UTF-8 文字列。

リモートリポジトリのプロバイダ。

- Repository – UTF-8 文字列。1~512 バイト長。

ジョブのアーティファクトを含むリモートリポジトリの名前。

- **Owner** – UTF-8 文字列。1~512 バイト長。

ジョブのアーティファクトを含むリモートリポジトリの所有者。

- **Branch** – UTF-8 文字列。1~512 バイト長。

リモートリポジトリ内のオプションのブランチ。

- **Folder** – UTF-8 文字列。1~512 バイト長。

リモートリポジトリ内のオプションのフォルダ。

- **LastCommitId** – UTF-8 文字列。1~512 バイト長。

リモートリポジトリ内のコミットの最後のコミット ID。

- **LastSyncTimestamp** – UTF-8 文字列。1~512 バイト長。

ジョブ同期が最後に実行された日時。

- **AuthStrategy** – UTF-8 文字列。

認証のタイプ。AWS Secrets Manager に保存されている認証トークン、または個人用アクセストークンです。

- **AuthToken** – UTF-8 文字列。1~512 バイト長。

認証トークンの値。

## 操作

- [CreateJob アクション \(Python: create\\_job\)](#)
- [UpdateJob アクション \(Python: update\\_job\)](#)
- [GetJob アクション \(Python: get\\_job\)](#)
- [GetJobs アクション \(Python: get\\_jobs\)](#)
- [DeleteJob アクション \(Python: delete\\_job\)](#)
- [ListJobs アクション \(Python: list\\_jobs\)](#)
- [BatchGetJobs アクション \(Python: batch\\_get\\_jobs\)](#)

## CreateJob アクション (Python: create\_job)

新しいジョブ定義を作成します。

## リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

このジョブ定義に割り当てる名前。アカウント内で一意にする必要があります。

- JobMode – UTF-8 文字列 (有効な値: SCRIPT="" | VISUAL="" | NOTEBOOK="")。

ジョブの作成方法を説明するモード。有効な値は次のとおりです。

- SCRIPT - ジョブは AWS Glue Studio スクリプトエディタを使用して作成されました。
- VISUAL - ジョブは AWS Glue Studio ビジュアルエディタを使用して作成されました。
- NOTEBOOK - ジョブはインタラクティブセッションノートブックを使用して作成されました。

JobMode フィールドが欠落しているか null の場合、SCRIPT がデフォルト値として割り当てられます。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

定義するジョブの説明。

- LogUri – UTF-8 文字列。

このフィールドは、将来の利用のために予約されています。

- Role – 必須: UTF-8 文字列。

このジョブに関連付けられている IAM ロールの名前または Amazon リソースネーム (ARN)。

- ExecutionProperty – [ExecutionProperty](#) オブジェクト。

このジョブに許可される同時実行の最大数を指定する ExecutionProperty。

- Command – 必須: [JobCommand](#) オブジェクト。

このジョブを実行する JobCommand。

- DefaultArguments – キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

名前と値のペアとして指定された、このジョブを実行するごとのデフォルトの引数。

ここで独自のジョブ実行スクリプトが消費する引数と、AWS Glue それ自体が消費する引数を指定できます。

ジョブ引数はログに記録される場合があります。プレーンテキストのシークレットを引数として渡さないでください。ジョブ内に保持する場合は、AWS Glue 接続 AWS Secrets Manager、またはその他のシークレット管理メカニズムからシークレットを取得します。

独自のジョブ引数を指定および使用方法については、デベロッパーガイドのトピック「[Calling AWS Glue APIs in Python](#)」を参照してください。

Spark ジョブの設定時にこのフィールドに指定できる引数については、デベロッパーガイドのトピック「[Special Parameters Used by AWS Glue](#)」を参照してください。

Ray ジョブの設定時にこのフィールドに指定できる引数については、デベロッパーガイドの「[Ray ジョブでジョブパラメータを使用する](#)」を参照してください。

- NonOverridableArguments – キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

名前と値のペアとして指定された、ジョブ実行時にジョブ引数を指定しても上書きされないこのジョブの引数。

- Connections – [ConnectionsList](#) オブジェクト。

このジョブに使用される接続。

- MaxRetries – 数値 (整数)。

失敗した場合にこのジョブを再試行する最大回数。

- AllocatedCapacity – 数値 (整数)。

このパラメータは廃止されました。代わりに MaxCapacity を使用します。

このジョブに割り当てる AWS Glue データ処理ユニット (DPUsの数。最低 2 つの DPU を割り当てることができます。デフォルトは 10 です。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

- Timeout - 数値 (整数)。1 以上。

ジョブのタイムアウト (分)。ジョブ実行が終了済みになって TIMEOUT ステータスに入るまでに、ジョブ実行でリソースを消費できる最大時間です。バッチジョブのデフォルト値は 2,880 分 (48 時間) です。

ストリーミングジョブのタイムアウト値は 7 日または 10080 分未満である必要があります。値を空白のままにすると、メンテナンスウィンドウを設定していない場合、ジョブは 7 日後に再起動されます。メンテナンスウィンドウを設定している場合、7 日後のメンテナンスウィンドウ中に再起動されます。

- MaxCapacity – 数値 (double)。

Glue バージョン 1.0 以前のジョブの場合、標準ワーカータイプを使用して、このジョブの実行時に割り当てることができる AWS Glue データ処理ユニット (DPUs) の数。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

一方、Glue バージョン 2.0 以降のジョブでは、Maximum capacity を指定できません。代わりに、Worker type と Number of workers を指定する必要があります。

WorkerType および NumberOfWorkers を使用している場合は MaxCapacity を設定しないでください。

MaxCapacity に割り当てることができる値は、Python シェルジョブ、Apache Spark ETL ジョブ、Apache Spark ストリーミング ETL ジョブのいずれを実行しているかによって異なります。

- Python シェルジョブを指定すると (JobCommand.Name="pythonshell")、0.0625 または 1 DPU のいずれかを割り当てることができます。デフォルトは 0.0625 DPU です。
- Apache Spark ETL ジョブ (JobCommand.Name="glueetl") または Apache Spark ストリーミング ETL ジョブ (JobCommand.Name="gluestreaming") を指定した場合は、2 ~ 100 の DPU を割り当てることができます。デフォルトでは 10 DPU になっています。このジョブタイプには、小数の DPU 割り当てを指定できません。
- SecurityConfiguration – UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。

このジョブで使用される SecurityConfiguration 構造の名前。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1 ~ 128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

このジョブと一緒に使用するタグです。ジョブへのアクセスを制限するためにタグを使用することができます。のタグの詳細については AWS Glue、デベロッパーガイドの[AWS 「のタグ AWS Glue」](#)を参照してください。

- NotificationProperty – [NotificationProperty](#) オブジェクト。

ジョブ通知の設定プロパティを指定します。

- GlueVersion – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

Spark ジョブでは、はジョブ AWS Glue で使用できる Apache Spark と Python のバージョン GlueVersion を決定します。Python バージョンは、Spark タイプのジョブでサポートされるバージョンを示します。

Ray ジョブの場合、GlueVersion を 4.0 以降に設定する必要があります。ただし、Ray ジョブで使用できる Ray、Python、および追加ライブラリのバージョンは、Job コマンドの Runtime パラメータによって決まります。

利用可能な AWS Glue バージョンと対応する Spark および Python バージョンの詳細については、デベロッパーガイドの「[Glue バージョン](#)」を参照してください。

Glue バージョンを指定せずに作成されたジョブは、デフォルトで Glue 0.9 に設定されます。

- NumberOfWorkers – 数値 (整数)。

ジョブの実行時に割り当てられた、定義済みの workerType ワーカー数。

- WorkerType – UTF-8 文字列 (有効な値: Standard="" | G.1X="" | G.2X="" | G.025X="" | G.4X="" | G.8X="" | Z.2X="")。

ジョブの実行時に割り当てられる事前定義済みのワーカーの種類。Spark ジョブに使用できる値は G.1X、G.2X、G.4X、G.8X、または G.025X です。Ray ジョブに使用できる値は Z.2X です。

- G.1X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 1 DPU (4 vCPU、16 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケーラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- G.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 77 GB の空き容量) を備えた 2 DPU (8 vCPU、32 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼ

キューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。

- G.4X ワーカータイプでは、各ワーカーは 256 GB のディスク (約 235 GB の空き容量) を備えた 4 DPU (16 vCPU、64 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、米国東部 (オハイオ)、米国東部 (バージニア北部)、米国西部 (オレゴン)、アジアパシフィック (シンガポール)、アジアパシフィック (シドニー)、アジアパシフィック (東京)、カナダ (中部)、欧州 (フランクフルト)、欧州 (アイルランド)、欧州 (ストックホルム) の各 AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL ジョブでのみ使用できます。
- G.8X ワーカータイプでは、各ワーカーは 512 GB のディスク (約 487 GB の空き容量) を備えた 8 DPU (32 vCPU、128 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、ワーカータイプでサポートされているのと同じ AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL G.4X ジョブでのみ使用できます。
- G.025X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 0.25 DPU (2 vCPU、4 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。少量のストリーミングジョブには、このワーカータイプをお勧めします。このワーカータイプは、AWS Glue バージョン 3.0 ストリーミングジョブでのみ使用できます。
- Z.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 120 GB の空き容量) を備えた 2 M-DPU (8 vCPU、64 GB のメモリ) にマッピングされており、オートスケララーに基づき最大 8 個の Ray ワーカーを提供します。
- CodeGenConfigurationNodes – キーバリューペアのマップ配列。

各キーは、[Custom string pattern #39](#) に適合する UTF-8 文字列です。

各値は [CodeGenConfigurationNode](#) オブジェクトです。

Glue Studio ビジュアルコンポーネントと、Glue Studio によるコード生成の両方がベースとする、有向非巡回グラフの表現。

- ExecutionClass - UTF-8 文字列。16 バイト長以下 (有効値: FLEX="" | STANDARD="")。

ジョブが標準実行クラスまたは柔軟な実行クラスのどちらで実行されるのかを示します。標準の実行クラスは、素早くジョブを起動する必要があり、専用のリソースが必要な時間的な制約のあるワークロードに最適です。

柔軟な実行クラスは、開始時刻と完了時刻が変化する時間的な制約が厳しいジョブに適しています。

AWS Glue バージョン 3.0 以降およびコマンドタイプのジョブのみを ExecutionClass glueetl に設定できます FLEX。柔軟な実行クラスは Spark ジョブで使用できます。

- SourceControlDetails – [SourceControlDetails](#) オブジェクト。

ジョブのソース管理設定の詳細。これにより、リモートリポジトリとの間でジョブアーティファクトを同期できます。

- MaintenanceWindow – UTF-8 文字列、「[Custom string pattern #30](#)」に一致。

このフィールドは、ストリーミングジョブのメンテナンスウィンドウの曜日と時間を指定します。AWS Glue は定期的にメンテナンスアクティビティを実行します。これらのメンテナンスウィンドウ中に、AWS Glue はストリーミングジョブを再起動する必要があります。

AWS Glue は、指定されたメンテナンスウィンドウから 3 時間以内にジョブを再起動します。たとえば、月曜日の午前 10 時 (GMT) にメンテナンスウィンドウを設定すると、ジョブは午前 10 時 (GMT) から午後 1 時 (GMT) までの間に再起動されます。

- ProfileName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ジョブに関連付けられた AWS Glue 使用プロファイルの名前。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

このジョブ定義に指定された一意の名前。

## エラー

- InvalidInputException
- IdempotentParameterMismatchException
- AlreadyExistsException
- InternalServiceException

- `OperationTimeoutException`
- `ResourceNumberLimitExceededException`
- `ConcurrentModificationException`

## UpdateJob アクション (Python: `update_job`)

既存のジョブ定義を更新します。以前のジョブ定義はこの情報によって完全に上書きされます。

### リクエスト

- `JobName` – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

更新するジョブ定義の名前。

- `JobUpdate` – 必須: [JobUpdate](#) オブジェクト。

ジョブ定義の更新に使用する値を指定します。指定されていない設定は削除されるか、デフォルト値にリセットされます。

- `ProfileName` – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

ジョブに関連付けられた AWS Glue 使用プロファイルの名前。

### レスポンス

- `JobName` – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

更新されたジョブ定義の名前を返します。

### エラー

- `InvalidInputException`
- `EntityNotFoundException`
- `InternalServiceException`
- `OperationTimeoutException`
- `ConcurrentModificationException`

## GetJob アクション (Python: get\_job)

既存のジョブ定義を取得します。

### リクエスト

- JobName – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
取得するジョブ定義の名前。

### レスポンス

- Job – [Job](#) オブジェクト。  
リクエストされたジョブ定義。

### エラー

- InvalidInputException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException

## GetJobs アクション (Python: get\_jobs)

すべての現在のジョブ定義を取得します。

### リクエスト

- NextToken – UTF-8 文字列。  
継続トークン (これが継続呼び出しの場合)。
- MaxResults – 1~1000 の数値 (整数)。  
応答の最大サイズ。

### 応答

- Jobs – [Job](#) オブジェクトの配列。

ジョブ定義のリスト。

- NextToken – UTF-8 文字列。

継続トークン (一部のジョブ定義がまだ返されていない場合)。

## エラー

- InvalidInputException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException

## DeleteJob アクション (Python: delete\_job)

指定したジョブ定義を削除します。ジョブ定義が見つからない場合、例外はスローされません。

## リクエスト

- JobName – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

削除するジョブ定義の名前。

## レスポンス

- JobName – UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

削除されたジョブ定義の名前。

## エラー

- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## ListJobs アクション (Python: list\_jobs)

この AWS アカウントのすべてのジョブリソースの名前、または指定されたタグを持つリソースを取得します。このオペレーションにより、アカウントで利用可能なリソースとその名前を確認できます。

このオペレーションはオプションの Tags フィールドを受け取ります。このフィールドを応答のフィルターとして使用すると、タグ付きリソースをグループとして取得できます。タグフィルタリングの使用を選択した場合は、タグが付いたリソースのみが取得されます。

### リクエスト

- NextToken – UTF-8 文字列。

継続トークン (これが継続リクエストの場合)。

- MaxResults – 1 ~ 1000 の数値 (整数)。

返されるリストの最大サイズ。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1 ~ 128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

これらのタグ付きリソースのみを返すように指定します。

### レスポンス

- JobNames – UTF-8 文字列の配列。

アカウント内のすべてのジョブの名前、または指定されたタグを持つジョブの名前。

- NextToken – UTF-8 文字列。

継続トークン (戻されたリストに最後に使用可能なメトリクスが含まれていない場合)。

### エラー

- InvalidInputException
- EntityNotFoundException

- `InternalServiceException`
- `OperationTimeoutException`

## BatchGetJobs アクション (Python: `batch_get_jobs`)

指定されたジョブ名のリストのリソースメタデータのリストを返します。ListJobs オペレーションを呼び出した後で、このオペレーションを呼び出すことで、アクセス許可が付与されているデータにアクセスできます。このオペレーションは、タグを使用するアクセス許可条件を含め、すべてのIAM のアクセス許可をサポートします。

### リクエスト

- `JobNames` – 必須: UTF-8 文字列の配列。

ジョブ名のリスト。これは ListJobs 操作から返された名前であることもあります。

### 応答

- `Jobs` – [Job](#) オブジェクトの配列。

ジョブ定義のリスト。

- `JobsNotFound` – UTF-8 文字列の配列。

ジョブの名前のリストが見つかりません。

### エラー

- `InternalServiceException`
- `OperationTimeoutException`
- `InvalidInputException`

## ジョブ実行

Jobs Runs API は、でのジョブ実行の開始、停止、表示、およびジョブブックマークのリセットに関連するデータ型と API について説明します AWS Glue。ジョブ実行履歴には、ワークフローとジョブ実行のために 90 日間アクセスできます。

## データ型

- [JobRun 構造](#)
- [Predecessor 構造](#)
- [JobBookmarkEntry 構造](#)
- [BatchStopJobRunSuccessfulSubmission 構造](#)
- [BatchStopJobRunError 構造](#)
- [NotificationProperty 構造](#)

## JobRun 構造

ジョブ実行についての情報が含まれています。

### フィールド

- Id – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

このジョブ実行の ID。

- Attempt – 数値 (整数)。

このジョブを実行しようとした回数。

- PreviousRunId – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

このジョブの以前の実行の ID。例えば、StartJobRun アクションで指定された JobRunId を表します。

- TriggerName – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

このジョブ実行を開始したトリガーの名前。

- JobName – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

この実行で使用されているジョブ定義の名前。

- JobMode – UTF-8 文字列 (有効な値: SCRIPT="" | VISUAL="" | NOTEBOOK="")。

ジョブの作成方法を説明するモード。有効な値は次のとおりです。

- SCRIPT - ジョブは AWS Glue Studio スクリプトエディタを使用して作成されました。
- VISUAL - ジョブは AWS Glue Studio ビジュアルエディタを使用して作成されました。

- NOTEBOOK - ジョブはインタラクティブセッションノートブックを使用して作成されました。

JobMode フィールドが欠落しているか null の場合、SCRIPT がデフォルト値として割り当てられます。

- StartedOn - タイムスタンプ。

このジョブ実行が開始された日付と時刻。

- LastModifiedOn - タイムスタンプ。

このジョブ実行が最後に変更された時刻。

- CompletedOn - タイムスタンプ。

このジョブ実行が完了した日付と時刻。

- JobRunState - UTF-8 文字列 (有効な値: STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT | ERROR | WAITING | EXPIRED)。

現在のジョブ実行の状態。異常終了したジョブのステータスの詳細については、「[AWS Glue Job Run のステータス](#)」を参照してください。

- Arguments - キーバリュペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

この実行に関連付けられているジョブの引数。このジョブ実行では、ジョブ定義自体に設定されているデフォルト引数を置き換えます。

ここで独自のジョブ実行スクリプトが消費する引数と、AWS Glue それ自体が消費する引数を指定できます。

ジョブ引数はログに記録される場合があります。プレーンテキストのシークレットを引数として渡さないでください。ジョブ内に保持する場合は、AWS Glue 接続 AWS Secrets Manager、またはその他のシークレット管理メカニズムからシークレットを取得します。

独自のジョブ引数を指定および使用する方法については、デベロッパーガイドのトピック「[Calling AWS Glue APIs in Python](#)」を参照してください。

Spark ジョブの設定時にこのフィールドに指定できる引数については、デベロッパーガイドのトピック「[Special Parameters Used by AWS Glue](#)」を参照してください。

Ray ジョブの設定時にこのフィールドに指定できる引数については、デベロッパーガイドの「[Ray ジョブでジョブパラメータを使用する](#)」を参照してください。

- ErrorMessage – UTF-8 文字列。

このジョブ実行に関連付けられているエラーメッセージ。

- PredecessorRuns – [Predecessor](#) オブジェクトの配列。

このジョブ実行に先行するもののリスト。

- AllocatedCapacity – 数値 (整数)。

このフィールドは廃止されました。代わりに MaxCapacity を使用します。

この に割り当てられた AWS Glue データ処理ユニット (DPUs) の数 JobRun。2~100 DPU の範囲で割り当てることができます。デフォルト値は 10 です。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

- ExecutionTime – 数値 (整数)。

ジョブ実行でリソースを消費した時間 (秒)。

- Timeout - 数値 (整数)。1 以上。

JobRun のタイムアウト (分)。ジョブ実行が終了済みになって TIMEOUT ステータスに入るまでに、ジョブ実行でリソースを消費できる最大時間です。この値は、親ジョブで設定したタイムアウト値を上書きします。

ストリーミングジョブのタイムアウト値は 7 日または 10080 分未満である必要があります。値を空白のままにすると、メンテナンスウィンドウを設定していない場合、ジョブは 7 日後に再起動されます。メンテナンスウィンドウを設定している場合、7 日後のメンテナンスウィンドウ中に再起動されます。

- MaxCapacity – 数値 (double)。

Glue バージョン 1.0 以前のジョブの場合、標準ワーカータイプを使用して、このジョブの実行時に割り当てることができる AWS Glue データ処理ユニット (DPUs) の数。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

一方、Glue バージョン 2.0 以降のジョブでは、Maximum capacity を指定できません。代わりに、Worker type と Number of workers を指定する必要があります。

WorkerType および NumberOfWorkers を使用している場合は MaxCapacity を設定しないでください。

MaxCapacity に割り当てることができる値は、Python シェルジョブ、Apache Spark ETL ジョブ、Apache Spark ストリーミング ETL ジョブのいずれを実行しているかによって異なります。

- Python シェルジョブを指定すると (JobCommand.Name="pythonshell")、0.0625 または 1 DPU のいずれかを割り当てることができます。デフォルトは 0.0625 DPU です。
- Apache Spark ETL ジョブ (JobCommand.Name="glueetl") または Apache Spark ストリーミング ETL ジョブ (JobCommand.Name="gluestreaming") を指定した場合は、2~100 の DPU を割り当てることができます。デフォルトでは 10 DPU になっています。このジョブタイプには、小数の DPU 割り当てを指定できません。
- WorkerType – UTF-8 文字列 (有効な値: Standard="" | G.1X="" | G.2X="" | G.025X="" | G.4X="" | G.8X="" | Z.2X="")。

ジョブの実行時に割り当てられる事前定義済みのワーカーの種類。Spark ジョブに使用できる値は G.1X、G.2X、G.4X、G.8X、または G.025X です。Ray ジョブに使用できる値は Z.2X です。

- G.1X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 1 DPU (4 vCPU、16 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- G.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 77 GB の空き容量) を備えた 2 DPU (8 vCPU、32 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- G.4X ワーカータイプでは、各ワーカーは 256 GB のディスク (約 235 GB の空き容量) を備えた 4 DPU (16 vCPU、64 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、米国東部 (オハイオ)、米国東部 (バージニア北部)、米国西部 (オレゴン)、アジアパシフィック (シンガポール)、アジアパシフィック (シドニー)、アジアパシフィック (東京)、カナダ (中部)、欧州 (フランクフルト)、欧州 (アイルランド)、欧州 (ストックホルム) の各 AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL ジョブでのみ使用できます。

- G.8X ワーカータイプでは、各ワーカーは 512 GB のディスク (約 487 GB の空き容量) を備えた 8 DPU (32 vCPU、128 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、ワーカータイプでサポートされているのと同じ AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL G.4X ジョブでのみ使用できます。
- G.025X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 0.25 DPU (2 vCPU、4 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。少量のストリーミングジョブには、このワーカータイプをお勧めします。このワーカータイプは、AWS Glue バージョン 3.0 ストリーミングジョブでのみ使用できます。
- Z.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 120 GB の空き容量) を備えた 2 M-DPU (8 vCPU、64 GB のメモリ) にマッピングされており、オートスケーラーに基づき最大 8 個の Ray ワーカーを提供します。
- NumberOfWorkers – 数値 (整数)。

ジョブの実行時に割り当てられた、定義済みの workerType ワーカー数。

- SecurityConfiguration – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このジョブ実行で使用される SecurityConfiguration 構造の名前。

- LogGroupName – UTF-8 文字列。

CloudWatch を使用して Amazon でサーバー側で暗号化できる安全なログ記録用のロググループの名前 AWS KMS。この名前には /aws-glue/jobs/ を指定できます。その場合、デフォルトの暗号化は NONE です。ロール名および SecurityConfiguration 名 (つまり /aws-glue/jobs-yourRoleName-yourSecurityConfigurationName/) を追加すると、そのセキュリティ設定がロググループの暗号化に使用されます。

- NotificationProperty – [NotificationProperty](#) オブジェクト。

ジョブ実行通知の構成プロパティを指定します。

- GlueVersion – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

Spark ジョブでは、はジョブ AWS Glue で使用できる Apache Spark と Python のバージョン GlueVersion を決定します。Python バージョンは、Spark タイプのジョブでサポートされるバージョンを示します。

Ray ジョブの場合、GlueVersion を 4.0 以降に設定する必要があります。ただし、Ray ジョブで使用できる Ray、Python、および追加ライブラリのバージョンは、Job コマンドの Runtime パラメータによって決まります。

利用可能な AWS Glue バージョンと対応する Spark および Python バージョンの詳細については、デベロッパーガイドの「[Glue バージョン](#)」を参照してください。

Glue バージョンを指定せずに作成されたジョブは、デフォルトで Glue 0.9 に設定されます。

- DPUSeconds – 数値 (double)。

このフィールドは、実行クラス FLEX で実行するジョブまたは自動スケーリングが有効になっているときに実行するジョブのいずれかについて設定できます。これは、ジョブ実行のライフサイクル中に各エグゼキューターが実行された合計時間 (秒単位) を DPU 係数 (G.1X ワーカーの場合は 1、G.2X ワーカーの場合は 2、G.025X ワーカーの場合は 0.25) で乗じたものを表します。Auto Scaling ジョブの場合に、この値は  $\text{executionEngineRuntime} * \text{MaxCapacity}$  とは異なる場合がありますが、これは特定の時間に実行されるエグゼキューターの数  $\text{MaxCapacity}$  より少ないことがあるためです。この場合、DPUSeconds の値が  $\text{executionEngineRuntime} * \text{MaxCapacity}$  より小さくなる可能性があります。

- ExecutionClass - UTF-8 文字列。16 バイト長以下 (有効値: FLEX="" | STANDARD="")。

ジョブが標準実行クラスまたは柔軟な実行クラスのどちらで実行されるのかを示します。標準の実行クラスは、素早くジョブを起動する必要があり、専用のリソースが必要な時間的な制約のあるワークロードに最適です。

柔軟な実行クラスは、開始時刻と完了時刻が変化する時間的制約が厳しいジョブに適しています。

AWS Glue バージョン 3.0 以降およびコマンドタイプのジョブのみを ExecutionClass glueet1 に設定できます FLEX。柔軟な実行クラスは Spark ジョブで使用できます。

- MaintenanceWindow – UTF-8 文字列、「[Custom string pattern #30](#)」に一致。

このフィールドは、ストリーミングジョブのメンテナンスウィンドウの曜日と時間を指定します。AWS Glue は定期的にメンテナンスアクティビティを実行します。これらのメンテナンスウィンドウ中に、AWS Glue はストリーミングジョブを再起動する必要があります。

AWS Glue は、指定されたメンテナンスウィンドウから 3 時間以内にジョブを再起動します。たとえば、月曜日の午前 10 時 (GMT) にメンテナンスウィンドウを設定すると、ジョブは午前 10 時 (GMT) から午後 1 時 (GMT) までの間に再起動されます。

- ProfileName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ジョブ実行に関連付けられた AWS Glue 使用プロファイルの名前。

## Predecessor 構造

このジョブ実行をトリガーした条件トリガーの述語に使用されたジョブ実行。

フィールド

- JobName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

先行するジョブ実行で使用したジョブ定義の名前。

- RunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

先行するジョブ実行のジョブ実行 ID。

## JobBookmarkEntry 構造

ジョブの処理を再開できるポイントを定義します。

フィールド

- JobName – UTF-8 文字列。

該当するジョブの名前。

- Version – 数値 (整数)。

ジョブのバージョン。

- Run – 数値 (整数)。

実行 ID 番号。

- Attempt – 数値 (整数)。

試行 ID 番号。

- PreviousRunId – UTF-8 文字列。

前回のジョブ実行に関連付けられた一意の実行識別子。

- RunId – UTF-8 文字列。

実行 ID 番号。

- JobBookmark – UTF-8 文字列。

ブックマーク自体。

## BatchStopJobRunSuccessfulSubmission 構造

指定された JobRun を停止するリクエストの成功を記録します。

フィールド

- JobName – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。  
停止したジョブ実行で使ったジョブ定義の名前。
- JobRunId – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。  
停止したジョブ実行の JobRunId。

## BatchStopJobRunError 構造

指定したジョブ実行を停止しようとして発生したエラーを記録します。

フィールド

- JobName – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。  
該当するジョブ実行で使ったジョブ定義の名前。
- JobRunId – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。  
該当するジョブ実行の JobRunId。
- ErrorDetail – [ErrorDetail](#) オブジェクト。  
発生したエラーに関する詳細を指定します。

## NotificationProperty 構造

通知の構成プロパティを指定します。

フィールド

- NotifyDelayAfter - 数値 (整数)。1 以上。

ジョブの実行が開始された後、ジョブ実行遅延通知を送信するまでの待機時間 (分単位)。

## 操作

- [StartJobRun アクション \(Python: start\\_job\\_run\)](#)
- [BatchStopJobRun アクション \(Python: batch\\_stop\\_job\\_run\)](#)
- [GetJobRun アクション \(Python: get\\_job\\_run\)](#)
- [GetJobRuns アクション \(Python: get\\_job\\_runs\)](#)
- [GetJobBookmark アクション \(Python: get\\_job\\_bookmark\)](#)
- [GetJobBookmarks アクション \(Python: get\\_job\\_bookmarks\)](#)
- [ResetJobBookmark アクション \(Python: reset\\_job\\_bookmark\)](#)

## StartJobRun アクション (Python: start\_job\_run)

ジョブ定義を使用してジョブ実行を開始します。

### リクエスト

- JobName – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

使用するジョブ定義の名前。

- JobRunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

再試行する以前の JobRun の ID。

- Arguments – キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

この実行に関連付けられているジョブの引数。このジョブ実行では、ジョブ定義自体に設定されているデフォルト引数を置き換えます。

ここで独自のジョブ実行スクリプトが消費する引数と、AWS Glue それ自体が消費する引数を指定できます。

ジョブ引数はログに記録される場合があります。プレーンテキストのシークレットを引数として渡さないでください。ジョブ内に保持する場合は、AWS Glue 接続 AWS Secrets Manager、またはその他のシークレット管理メカニズムからシークレットを取得します。

独自のジョブ引数を指定および使用方法については、デベロッパーガイドのトピック「[Calling AWS Glue APIs in Python](#)」を参照してください。

Spark ジョブの設定時にこのフィールドに指定できる引数については、デベロッパーガイドのトピック「[Special Parameters Used by AWS Glue](#)」を参照してください。

Ray ジョブの設定時にこのフィールドに指定できる引数については、デベロッパーガイドの「[Ray ジョブでジョブパラメータを使用する](#)」を参照してください。

- AllocatedCapacity - 数値 (整数)。

このフィールドは廃止されました。代わりに MaxCapacity を使用します。

この に割り当てる AWS Glue データ処理ユニット (DPUs) の数 JobRun。最低 2 つの DPU を割り当てることができます。デフォルトは 10 です。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

- Timeout - 数値 (整数)。1 以上。

JobRun のタイムアウト (分)。ジョブ実行が終了済みになって TIMEOUT ステータスに入るまでに、ジョブ実行でリソースを消費できる最大時間です。この値は、親ジョブで設定したタイムアウト値を上書きします。

ストリーミングジョブのタイムアウト値は 7 日または 10080 分未満である必要があります。値を空白のままにすると、メンテナンスウィンドウを設定していない場合、ジョブは 7 日後に再起動されます。メンテナンスウィンドウを設定している場合、7 日後のメンテナンスウィンドウ中に再起動されます。

- MaxCapacity - 数値 (double)。

Glue バージョン 1.0 以前のジョブの場合、標準ワーカータイプを使用して、このジョブの実行時に割り当てることができる AWS Glue データ処理ユニット (DPUs) の数。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

一方、Glue バージョン 2.0 以降のジョブでは、Maximum capacity を指定できません。代わりに、Worker type と Number of workers を指定する必要があります。

WorkerType および NumberOfWorkers を使用している場合は MaxCapacity を設定しないでください。

MaxCapacity に割り当てることができる値は、Python シェルジョブ、Apache Spark ETL ジョブ、Apache Spark ストリーミング ETL ジョブのいずれを実行しているかによって異なります。

- Python シェルジョブを指定すると (JobCommand.Name="pythonshell")、0.0625 または 1 DPU のいずれかを割り当てることができます。デフォルトは 0.0625 DPU です。
- Apache Spark ETL ジョブ (JobCommand.Name="glueetl") または Apache Spark ストリーミング ETL ジョブ (JobCommand.Name="gluestreaming") を指定した場合は、2~100 の DPU を割り当てることができます。デフォルトでは 10 DPU になっています。このジョブタイプには、小数の DPU 割り当てを指定できません。
- SecurityConfiguration – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このジョブ実行で使用される SecurityConfiguration 構造の名前。

- NotificationProperty – [NotificationProperty](#) オブジェクト。

ジョブ実行通知の構成プロパティを指定します。

- WorkerType – UTF-8 文字列 (有効な値: Standard="" | G.1X="" | G.2X="" | G.025X="" | G.4X="" | G.8X="" | Z.2X="")。

ジョブの実行時に割り当てられる事前定義済みのワーカーの種類。Spark ジョブに使用できる値は G.1X、G.2X、G.4X、G.8X、または G.025X です。Ray ジョブに使用できる値は Z.2X です。

- G.1X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 1 DPU (4 vCPU、16 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケーラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- G.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 77 GB の空き容量) を備えた 2 DPU (8 vCPU、32 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケーラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。

- G.4X ワーカータイプでは、各ワーカーは 256 GB のディスク (約 235 GB の空き容量) を備えた 4 DPU (16 vCPU、64 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、米国東部 (オハイオ)、米国東部 (バージニア北部)、米国西部 (オレゴン)、アジアパシフィック (シンガポール)、アジアパシフィック (シドニー)、アジアパシフィック (東京)、カナダ (中部)、欧州 (フランクフルト)、欧州 (アイルランド)、欧州 (ストックホルム) の各 AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL ジョブでのみ使用できます。
- G.8X ワーカータイプでは、各ワーカーは 512 GB のディスク (約 487 GB の空き容量) を備えた 8 DPU (32 vCPU、128 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、ワーカータイプでサポートされているのと同じ AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL G.4X ジョブでのみ使用できます。
- G.025X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 0.25 DPU (2 vCPU、4 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。少量のストリーミングジョブには、このワーカータイプをお勧めします。このワーカータイプは、AWS Glue バージョン 3.0 ストリーミングジョブでのみ使用できます。
- Z.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 120 GB の空き容量) を備えた 2 M-DPU (8 vCPU、64 GB のメモリ) にマッピングされており、オートスケーラーに基づき最大 8 個の Ray ワーカーを提供します。
- `NumberOfWorkers` - 数値 (整数)。

ジョブの実行時に割り当てられた、定義済みの `workerType` ワーカー数。

- `ExecutionClass` - UTF-8 文字列。16 バイト長以下 (有効値: `FLEX=""` | `STANDARD=""`)。

ジョブが標準実行クラスまたは柔軟な実行クラスのどちらで実行されるのかを示します。標準の実行クラスは、素早くジョブを起動する必要があり、専用のリソースが必要な時間的な制約のあるワークロードに最適です。

柔軟な実行クラスは、開始時刻と完了時刻が変化する時間的な制約が厳しいジョブに適しています。

AWS Glue バージョン 3.0 以降およびコマンドタイプのジョブのみを `ExecutionClass glueetl` に設定できます FLEX。柔軟な実行クラスは Spark ジョブで使用できます。

- `ProfileName` - UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ジョブ実行に関連付けられた AWS Glue 使用プロファイルの名前。

## レスポンス

- JobRunId – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

このジョブ実行に割り当てられた ID。

## エラー

- InvalidInputException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException
- ResourceNumberLimitExceededException
- ConcurrentRunsExceededException

## BatchStopJobRun アクション (Python: batch\_stop\_job\_run)

指定したジョブ定義の 1 つ以上のジョブ実行を停止します。

### リクエスト

- JobName – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

ジョブ実行を停止するジョブ定義の名前。

- JobRunIds – 必須: UTF-8 文字列の配列。1 ~ 25 個の文字列。

このジョブ定義で停止する JobRunIds のリスト。

### 応答

- SuccessfulSubmissions – [BatchStopJobRunSuccessfulSubmission](#) オブジェクトの配列。

停止のために JobRuns 正常に送信された のリスト。

- Errors – [BatchStopJobRunError](#) オブジェクトの配列。

JobRuns を停止しようとして発生したエラーのリスト。各エラーで発生した JobRunId とエラーの詳細が含まれます。

## エラー

- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`

## GetJobRun アクション (Python: `get_job_run`)

指定されたジョブ実行のメタデータを取得します。ジョブ実行履歴には、ワークフローとジョブ実行のために 90 日間アクセスできます。

### リクエスト

- `JobName` – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

実行中のジョブ定義の名前。

- `RunId` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ジョブ実行の ID。

- `PredecessorsIncluded` – ブール。

以前の実行のリストが返される場合は、`true` です。

### レスポンス

- `JobRun` – [JobRun](#) オブジェクト。

リクエストされたジョブ実行のメタデータ。

## エラー

- `InvalidInputException`
- `EntityNotFoundException`
- `InternalServiceException`

- `OperationTimeoutException`

## GetJobRuns アクション (Python: `get_job_runs`)

特定のジョブ定義に該当するすべての実行のメタデータを取得します。

### リクエスト

- `JobName` – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
すべてのジョブ実行を取得するジョブ定義の名前。
- `NextToken` – UTF-8 文字列。  
継続トークン (これが継続呼び出しの場合)。
- `MaxResults` – 1 未満または 200 を超えない数値 (整数)。  
応答の最大サイズ。

### 応答

- `JobRuns` – [JobRun](#) オブジェクトの配列。  
ジョブ実行のメタデータオブジェクトのリスト。
- `NextToken` – UTF-8 文字列。  
継続トークン (一部のリクエストされたジョブがまだ返されていない場合)。

### エラー

- `InvalidInputException`
- `EntityNotFoundException`
- `InternalServiceException`
- `OperationTimeoutException`

## GetJobBookmark アクション (Python: `get_job_bookmark`)

ジョブブックマークエントリに関する情報を返します。

ジョブのブックマークの有効化と使用の詳細については、以下を参照してください。

- [ジョブのブックマークを使用した処理済みデータの追跡](#)
- [で使用されるジョブパラメータ AWS Glue](#)
- [Job 構造](#)

## リクエスト

- JobName – 必須: UTF-8 文字列。

該当するジョブの名前。

- Version – 数値 (整数)。

ジョブのバージョン。

- RunId – UTF-8 文字列。

このジョブの実行に関連付けられた一意の実行識別子。

## レスポンス

- JobBookmarkEntry – [JobBookmarkEntry](#) オブジェクト。

ジョブの処理を再開できるポイントを定義する構造。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- ValidationException

## GetJobBookmarks アクション (Python: get\_job\_bookmarks)

ジョブブックマークエントリに関する情報を返します。リストは、バージョン番号の小さい順に並べられます。

ジョブのブックマークの有効化と使用の詳細については、以下を参照してください。

- [ジョブのブックマークを使用した処理済みデータの追跡](#)
- [で使用されるジョブパラメータ AWS Glue](#)
- [Job 構造](#)

## リクエスト

- JobName – 必須: UTF-8 文字列。

該当するジョブの名前。

- MaxResults – 数値 (整数)。

応答の最大サイズ。

- NextToken – 数値 (整数)。

継続トークン (これが継続呼び出しの場合)。

## 応答

- JobBookmarkEntries – [JobBookmarkEntry](#) オブジェクトの配列。

ジョブの処理を再開できるポイントを定義するジョブブックマークエントリのリスト。

- NextToken – 数値 (整数)。

継続トークン。すべてのエントリが返される場合は 1、要求されたすべてのジョブ実行が返されなかった場合は 1 より大きい値を持ちます。

## エラー

- InvalidInputException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException

## ResetJobBookmark アクション (Python: reset\_job\_bookmark)

ブックマークエントリをリセットします。

ジョブのブックマークの有効化と使用の詳細については、以下を参照してください。

- [ジョブのブックマークを使用した処理済みデータの追跡](#)
- [で使用されるジョブパラメータ AWS Glue](#)
- [Job 構造](#)

### リクエスト

- JobName – 必須: UTF-8 文字列。

該当するジョブの名前。

- RunId – UTF-8 文字列。

このジョブの実行に関連付けられた一意の実行識別子。

### レスポンス

- JobBookmarkEntry – [JobBookmarkEntry](#) オブジェクト。

ブックマークエントリのリセット。

### エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## トリガー

Triggers API では、AWS Glue におけるジョブトリガーの作成、更新、削除と、開始および停止に関連するデータ型と API について説明します。

## データ型

- [トリガー構造](#)
- [TriggerUpdate 構造](#)
- [述語構造](#)
- [条件の構造](#)
- [アクション構造](#)
- [EventBatchingCondition 構造](#)

### トリガー構造

特定のトリガーに関する情報です。

#### フィールド

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

トリガーの名前。

- WorkflowName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

トリガーに関連付けられているワークフローの名前。

- Id – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

将来の利用のために予約されています。

- Type – UTF-8 文字列 (有効な値: SCHEDULED | CONDITIONAL | ON\_DEMAND | EVENT)。

これがトリガーのタイプです。

- State – UTF-8 文字列 (有効な値: CREATING | CREATED | ACTIVATING | ACTIVATED | DEACTIVATING | DEACTIVATED | DELETING | UPDATING)。

現在のトリガーの状態。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

このトリガーの説明。

- Schedule – UTF-8 文字列。

スケジュールを指定するために使用される cron 式 ([ジョブとクローラーの時間ベースのスケジュール](#)を参照してください。)たとえば、毎日 12:15 UTC に何かを実行するには、cron(15 12 \* \* ? \*) を指定します。

- Actions – [アクション](#) オブジェクトの配列。

このトリガーによって開始されるアクション。

- Predicate – [述語](#) オブジェクト。

このトリガーの述語は、いつトリガーを起動するかを定義します。

- EventBatchingCondition – [EventBatchingCondition](#) オブジェクト。

EventBridge イベントトリガーが起動する前に、満たす必要があるBatch 条件 (指定された受信イベント数またはバッチ時間ウィンドウの期限切れ)。

## TriggerUpdate 構造

トリガーの更新に使用する情報を提供するために使用される構造です。このオブジェクトは、前のトリガー定義を完全に上書きして更新します。

### フィールド

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

将来の利用のために予約されています。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

このトリガーの説明。

- Schedule – UTF-8 文字列。

スケジュールを指定するために使用される cron 式 ([ジョブとクローラーの時間ベースのスケジュール](#)を参照してください。)たとえば、毎日 12:15 UTC に何かを実行するには、cron(15 12 \* \* ? \*) を指定します。

- Actions – [アクション](#) オブジェクトの配列。

このトリガーによって開始されるアクション。

- Predicate – [述語](#) オブジェクト。

このトリガーの述語は、いつトリガーを起動するかを定義します。

- `EventBatchingCondition` – [EventBatchingCondition](#) オブジェクト。

EventBridge イベントトリガーが起動する前に、満たす必要があるBatch 条件 (指定された受信イベント数またはバッチ時間ウィンドウの期限切れ)。

## 述語構造

トリガーがいつ起動するかを決定するトリガーの述語を定義します。

### フィールド

- `Logical` – UTF-8 文字列 (有効な値: AND | ANY)。

1 つの条件のみが表示されている場合のオプションフィールドです。複数の条件が表示されている場合、このフィールドは必須です。

- `Conditions` – [条件](#) オブジェクトの配列。

トリガーがいつ起動するかを決定する条件のリスト。

## 条件の構造

トリガーが起動する条件を定義します。

### フィールド

- `LogicalOperator` – UTF-8 文字列 (有効な値: EQUALS)。

論理演算子。

- `JobName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

JobRuns のこの条件が適用され、このトリガーが待機するジョブの名前。

- `State` – UTF-8 文字列 (有効な値: STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT | ERROR | WAITING | EXPIRED)。

条件の状態。現在、トリガーがリスンできるジョブ状態は SUCCEEDED、STOPPED、FAILED、および TIMEOUT のみです。トリガーがリスンできるクローラー状態は SUCCEEDED、FAILED、および CANCELLED のみです。

- `CrawlerName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

この条件が適用されるクローラーの名前。

- `CrawlState` - UTF-8 文字列 (有効値: `RUNNING` | `CANCELLING` | `CANCELLED` | `SUCCEEDED` | `FAILED` | `ERROR`)。

この条件が適用されるクローラーの状態。

## アクション構造

トリガーによって開始されるアクションを定義します。

### フィールド

- `JobName` - UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

実行されるジョブの名前。

- `Arguments` - キーバリューペアのマップ配列。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

このトリガーが起動するとき使用されるジョブ引数。このジョブ実行では、ジョブ定義自体に設定されているデフォルト引数を置き換えます。

独自のジョブ実行スクリプトが消費する引数だけでなく、AWS Glue が消費する引数もここで指定できます。

独自のジョブ引数を指定および使用する方法については、デベロッパーガイドのトピック「[Calling AWS Glue APIs in Python](#)」を参照してください。

AWS Glue がジョブを設定するために使用するキーと値のペアについては、デベロッパーガイドのトピック「[Special Parameters Used by AWS Glue](#)」を参照してください。

- `Timeout` - 数値 (整数)。1 以上。

`JobRun` のタイムアウト (分)。ジョブ実行が終了済みになって `TIMEOUT` ステータスに入るまでに、ジョブ実行でリソースを消費できる最大時間です。デフォルト値は 2,880 分 (48 時間) です。これにより、親ジョブで設定したタイムアウト値が上書きされます。

- `SecurityConfiguration` - UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このジョブで使用される SecurityConfiguration 構造の名前。

- NotificationProperty – [NotificationProperty](#) オブジェクト。

ジョブ実行通知の構成プロパティを指定します。

- CrawlerName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このアクションで使用されるクローラーの名前。

## EventBatchingCondition 構造

EventBridge イベントトリガーが起動する前に、満たす必要がある Batch 条件 (指定された受信イベント数またはバッチ時間ウィンドウの期限切れ)。

### フィールド

- BatchSize – 必須: 数値 (integer)。1 ~ 100。

EventBridge イベントトリガーが発生する前に、Amazon EventBridge から受信する必要があるイベントの数。

- BatchWindow – 数値 (integer)。1 ~ 900。

EventBridge イベントトリガーが起動するまでの時間 (秒単位)。最初のイベントを受信したときからの時間です。

## 操作

- [CreateTrigger アクション \(Python: create\\_trigger\)](#)
- [StartTrigger アクション \(Python: start\\_trigger\)](#)
- [GetTrigger アクション \(Python: get\\_trigger\)](#)
- [GetTriggers アクション \(Python: get\\_triggers\)](#)
- [UpdateTrigger アクション \(Python: update\\_trigger\)](#)
- [StopTrigger アクション \(Python: stop\\_trigger\)](#)
- [DeleteTrigger アクション \(Python: delete\\_trigger\)](#)
- [ListTriggers アクション \(Python: list\\_triggers\)](#)
- [BatchGetTriggers アクション \(Python: batch\\_get\\_triggers\)](#)

## CreateTrigger アクション (Python: create\_trigger)

新しいトリガーを作成します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

トリガーの名前。

- WorkflowName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

トリガーに関連付けられているワークフローの名前。

- Type – 必須: UTF-8 文字列 (有効な値: SCHEDULED | CONDITIONAL | ON\_DEMAND | EVENT)。

新しいトリガーのタイプ。

- Schedule – UTF-8 文字列。

スケジュールを指定するために使用される cron 式 ([ジョブとクローラーの時間ベースのスケジュール](#)を参照してください。)たとえば、毎日 12:15 UTC に何かを実行するには、cron(15 12 \* \* ? \*) を指定します。

このフィールドは、トリガータイプが SCHEDULED の場合に必要です。

- Predicate – [述語](#) オブジェクト。

新しいトリガーがいつ起動するかを指定する述語です。

このフィールドは、トリガータイプが CONDITIONAL の場合に必要です。

- Actions – 必須: [アクション](#) オブジェクトの配列。

このトリガーが起動したときに開始されるアクション。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

新しいトリガーの説明。

- StartOnCreation – ブール。

作成時に SCHEDULED および CONDITIONAL トリガーを起動するには、true に設定します。True は ON\_DEMAND トリガーではサポートされません。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

このトリガーで使用するタグ。トリガーへのアクセスを制限するためにタグを使用することができません。AWS Glue のタグの詳細については、デベロッパーガイドの「[AWS Tags in AWS Glue](#)」を参照してください。

- EventBatchingCondition – [EventBatchingCondition](#) オブジェクト。

EventBridge イベントトリガーが起動する前に、満たす必要がある Batch 条件 (指定された受信イベント数またはバッチ時間ウィンドウの期限切れ)。

## 応答

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

トリガーの名前。

## エラー

- AlreadyExistsException
- EntityNotFoundException
- InvalidInputException
- IdempotentParameterMismatchException
- InternalServiceException
- OperationTimeoutException
- ResourceNumberLimitExceededException
- ConcurrentModificationException

## StartTrigger アクション (Python: start\_trigger)

既存のトリガーを開始します。さまざまなタイプのトリガーの開始方法については、「[ジョブのトリガー](#)」を参照してください。

## リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

開始するトリガーの名前。

応答

- Name – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

開始されたトリガーの名前。

エラー

- `InvalidInputException`
- `InternalServiceException`
- `EntityNotFoundException`
- `OperationTimeoutException`
- `ResourceNumberLimitExceededException`
- `ConcurrentRunsExceededException`

GetTrigger アクション (Python: `get_trigger`)

トリガーの定義を取得します。

リクエスト

- Name – 必須: UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

取得するトリガーの名前。

応答

- `Trigger` – [Trigger トリガー](#) オブジェクト。

リクエストされたトリガー定義。

エラー

- `EntityNotFoundException`

- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`

## GetTriggers アクション (Python: `get_triggers`)

ジョブに関連付けられているすべてのトリガーを取得します。

### リクエスト

- `NextToken` – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

- `DependentJobName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

トリガーを取得するジョブの名前。このジョブを開始できるトリガーが返されます。このようなトリガーがない場合、すべてのトリガーが返されます。

- `MaxResults` – 1 未満または 200 を超えない数値 (整数)。

応答の最大サイズ。

### 応答

- `Triggers` – [Trigger トリガー](#) オブジェクトの配列。

指定されたジョブのトリガーのリスト。

- `NextToken` – UTF-8 文字列。

リクエストされたトリガーのすべてがまだ返されていない場合は、継続トークン。

### エラー

- `EntityNotFoundException`
- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`

## UpdateTrigger アクション (Python: update\_trigger)

トリガー定義を更新します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

更新するトリガーの名前。

- TriggerUpdate – 必須: [TriggerUpdate](#) オブジェクト。

トリガーの更新に使用する新しい値。

### 応答

- Trigger – [Trigger トリガー](#)) オブジェクト。

結果として生じるトリガー定義。

### エラー

- InvalidInputException
- InternalServiceException
- EntityNotFoundException
- OperationTimeoutException
- ConcurrentModificationException

## StopTrigger アクション (Python: stop\_trigger)

指定されたトリガーを停止します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

停止するトリガーの名前。

## 応答

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

停止したトリガーの名前。

## エラー

- `InvalidInputException`
- `InternalServiceException`
- `EntityNotFoundException`
- `OperationTimeoutException`
- `ConcurrentModificationException`

## DeleteTrigger アクション (Python: delete\_trigger)

指定されたトリガーを削除します。トリガーが見つからない場合、例外はスローされません。

## リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除するトリガーの名前。

## 応答

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除されたトリガーの名前。

## エラー

- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`
- `ConcurrentModificationException`

## ListTriggers アクション (Python: list\_triggers)

この AWS アカウントのすべてのトリガーリソース、または指定されたタグを持つリソースの名前を取得します。このオペレーションにより、アカウントで利用可能なリソースとその名前を確認できます。

このオペレーションはオプションの Tags フィールドを受け取ります。このフィールドを応答のフィルターとして使用すると、タグ付きリソースをグループとして取得できます。タグフィルタリングの使用を選択した場合は、タグが付いたリソースのみが取得されます。

### リクエスト

- NextToken – UTF-8 文字列。

継続トークン (これが継続リクエストの場合)。

- DependentJobName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

トリガーを取得するジョブの名前。このジョブを開始できるトリガーが返されます。このようなトリガーがない場合は、すべてのトリガーが返ります。

- MaxResults – 1 未満または 200 を超えない数値 (整数)。

返されるリストの最大サイズ。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

これらのタグ付きリソースのみを返すように指定します。

### 応答

- TriggerNames – UTF-8 文字列の配列。

アカウント内のすべてのトリガーの名前、または指定されたタグを持つトリガーの名前。

- NextToken – UTF-8 文字列。

継続トークン (戻されたリストに最後に使用可能なメトリクスが含まれていない場合)。

## エラー

- EntityNotFoundException
- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## BatchGetTriggers アクション (Python: batch\_get\_triggers)

指定されたトリガー名のリストのリソースメタデータのリストを返します。ListTriggers オペレーションを呼び出した後で、このオペレーションを呼び出すことで、アクセス許可が付与されているデータにアクセスできます。このオペレーションは、タグを使用するアクセス許可条件を含め、すべての IAM のアクセス許可をサポートします。

### リクエスト

- TriggerNames – 必須: UTF-8 文字列の配列。

トリガー名のリスト。これは ListTriggers オペレーションから返される名前です。

### 応答

- Triggers – [Trigger トリガー](#)) オブジェクトの配列。

トリガー定義のリスト。

- TriggersNotFound – UTF-8 文字列の配列。

トリガーの名前のリストが見つかりません。

## エラー

- InternalServiceException
- OperationTimeoutException
- InvalidInputException

# インタラクティブセッション API

インタラクティブセッション API は、AWS Glue インタラクティブセッションを使用してデータ統合用の抽出、変換、ロード (ETL) スクリプトを構築およびテストすることに関連する AWS Glue API について説明します。

## データ型

- [セッション構造](#)
- [SessionCommand 構造](#)
- [Statement 構造](#)
- [StatementOutput 構造](#)
- [StatementOutputData 構造](#)
- [ConnectionsList 構造](#)

## セッション構造

リモート Spark ランタイム環境が実行されている期間。

### フィールド

- Id – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

セッションの ID

- CreatedOn – タイムスタンプ。

セッションが作成された日時。

- Status - UTF-8 文字列 (有効値: PROVISIONING | READY | FAILED | TIMEOUT | STOPPING | STOPPED)。

セッションのステータスです。

- ErrorMessage – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

セッション中に表示されるエラーメッセージです。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

セッションの説明です。

- Role - UTF-8 文字列。20 ~ 2,048 バイト長。 [Custom string pattern #26](#) に一致。

セッションに関連付けられている IAM ロールの名前または Amazon リソースネーム (ARN)。

- Command - [SessionCommand](#) オブジェクト。

コマンドオブジェクト。「」を参照してください SessionCommand。

- DefaultArguments - キーバリューペアのマップ配列。75 ペア以下。

各キーは UTF-8 文字列、1 ~ 128 バイト長で、 [Custom string pattern #27](#) に一致します。

各値は UTF-8 文字列、4096 バイト長で、 [URI address multi-line string pattern](#) に一致します。

キーバリューペアのマップ配列。最大ペア数は 75 です。

- Connections - [ConnectionsList](#) オブジェクト。

セッションに使用される接続の数。

- Progress - 数値 (double)。

セッションのコード実行の進行状況。

- MaxCapacity - 数値 (double)。

ジョブの実行時に割り当てることができる AWS Glue データ処理ユニット (DPU) の数。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。

- SecurityConfiguration - UTF-8 文字列、1 ~ 255 バイト長、 [Single-line string pattern](#) に一致。

セッションで使用する SecurityConfiguration 構造の名前。

- GlueVersion - UTF-8 文字列、1 ~ 255 バイト長、 [Custom string pattern #20](#) に一致。

AWS Glue バージョンによって、 が AWS Glue サポートする Apache Spark および Python のバージョンが決まります。は 2.0 より大きい GlueVersion 必要があります。

- DataAccessId - UTF-8 文字列、1 ~ 36 バイト。

セッションのデータアクセス ID。

- PartitionId - UTF-8 文字列、1 ~ 36 バイト。

セッションのパーティション ID。

- `NumberOfWorkers` – 数値 (整数)。

セッションに使用する、定義済み `WorkerType` のワーカー数。

- `WorkerType` – UTF-8 文字列 (有効な値: `Standard=""` | `G.1X=""` | `G.2X=""` | `G.025X=""` | `G.4X=""` | `G.8X=""` | `Z.2X=""`)。

セッションの実行時に割り当てられる事前定義済みのワーカーの種類。Spark セッション用に `G.1X`、`G.2X`、`G.4X`、または `G.8X` の値を使用できます。Ray セッション用に値 `Z.2X` を使用できます。

- `CompletedOn` – タイムスタンプ。

このセッションが完了した日付と時刻。

- `ExecutionTime` – 数値 (double)。

セッションの合計実行時間。

- `DPUSeconds` – 数値 (double)。

セッションによって消費される DPUs (式: `ExecutionTime * MaxCapacity`)。

- `IdleTimeout` – 数値 (整数)。

セッションがタイムアウトするまでのアイドル時の分数。

- `ProfileName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

セッションに関連付けられた AWS Glue 使用プロファイルの名前。

## SessionCommand 構造

ジョブを実行する `SessionCommand`。

フィールド

- `Name` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

の名前を指定します `SessionCommand`。「`glueetl`」または「`gluestreaming`」を使用できます。

- `PythonVersion` – UTF-8 文字列、「[Custom string pattern #21](#)」に一致。

Python バージョンを指定します。Python バージョンは、Spark タイプのジョブでサポートされるバージョンを示します。

## Statement 構造

セッションで特定のアクションを発生させるためのステートメントまたはリクエスト。

フィールド

- Id – 数値 (整数)。

ステートメントの ID。

- Code – UTF-8 文字列。

ステートメントの実行コード。

- State - UTF-8 文字列 (有効値: WAITING | RUNNING | AVAILABLE | CANCELLING | CANCELLED | ERROR)。

リクエストがアクションされている間の状態。

- Output – [StatementOutput](#) オブジェクト。

JSON 形式での出力。

- Progress – 数値 (double)。

コード実行の進行状況。

- StartedOn – 数値 (long 型)。

ジョブ定義が開始された UNIX の日時。

- CompletedOn – 数値 (long 型)。

ジョブ定義が完了した UNIX の日時。

## StatementOutput 構造

JSON 形式のコード実行の出力。

フィールド

- Data – [StatementOutputData](#) オブジェクト。

コード実行の出力。

- ExecutionCount – 数値 (整数)。

出力の実行回数。

- Status - UTF-8 文字列 (有効値: WAITING | RUNNING | AVAILABLE | CANCELLING | CANCELLED | ERROR)。

コード実行の出力状態。

- ErrorName - UTF-8 文字列。

出力内エラーの名前。

- ErrorValue - UTF-8 文字列。

出力のエラー値。

- Traceback - UTF-8 文字列の配列。

出力のトレースバック。

## StatementOutputData 構造

JSON 形式のコード実行の出力。

フィールド

- TextPlain - UTF-8 文字列。

テキスト形式のコード実行の出力。

## ConnectionsList 構造

ジョブが使用する接続を指定します。

フィールド

- Connections - UTF-8 文字列の配列。

ジョブが使用する接続のリスト。

## 操作

- [CreateSession アクション \(Python: create\\_session\)](#)

- [StopSession アクション \(Python: stop\\_session\)](#)
- [DeleteSession アクション \(Python: delete\\_session\)](#)
- [GetSession アクション \(Python: get\\_session\)](#)
- [ListSessions アクション \(Python: list\\_sessions\)](#)
- [RunStatement アクション \(Python: run\\_statement\)](#)
- [CancelStatement アクション \(Python: cancel\\_statement\)](#)
- [GetStatement アクション \(Python: get\\_statement\)](#)
- [ListStatements アクション \(Python: list\\_statements\)](#)

## CreateSession アクション (Python: create\_session)

新しいセッションを作成します。

### リクエスト

新しいセッションの作成をリクエストします。

- Id – 必須: UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。

セッションリクエストの ID。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

セッションの説明です。

- Role – 必須: UTF-8 文字列。20 ~ 2048 バイト長。[Custom string pattern #26](#) に一致。

IAM ロールの ARN

- Command – 必須: [SessionCommand](#) オブジェクト。

ジョブを実行する SessionCommand。

- Timeout - 数値 (整数)。1 以上。

セッションがタイムアウトするまでの時間 (分)。Spark ETLジョブのデフォルトは48時間 (2,880 分) であり、このジョブタイプの最大セッション持続時間です。他のジョブタイプについては、ドキュメントを参照してください。

- IdleTimeout - 数値 (整数)。1 以上。

セッションがタイムアウトするまでのアイドル時の分数。Spark ETL ジョブのデフォルトは、タイムアウト値です。他のジョブタイプについては、ドキュメントを参照してください。

- `DefaultArguments` – キーバリューペアのマップ配列。75 ペア以下。

各キーは UTF-8 文字列、1~128 バイト長で、[Custom string pattern #27](#) に一致します。

各値は UTF-8 文字列、4096 バイト長で、[URI address multi-line string pattern](#) に一致します。

キーバリューペアのマップ配列。最大ペア数は 75 です。

- `Connections` – [ConnectionsList](#) オブジェクト。

セッションに使用する接続の数。

- `MaxCapacity` – 数値 (double)。

ジョブの実行時に割り当てることができる AWS Glue データ処理ユニット (DPUs) の数。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。

- `NumberOfWorkers` – 数値 (整数)。

セッションに使用する、定義済み WorkerType のワーカー数。

- `WorkerType` – UTF-8 文字列 (有効な値: `Standard=""` | `G.1X=""` | `G.2X=""` | `G.025X=""` | `G.4X=""` | `G.8X=""` | `Z.2X=""`)。

ジョブの実行時に割り当てられる事前定義済みのワーカーの種類。Spark ジョブに使用できる値は G.1X、G.2X、G.4X、または G.8X です。Ray ノートブックに使用できる値は Z.2X です。

- G.1X ワーカータイプでは、各ワーカーは 84 GB のディスク (約 34 GB の空き容量) を備えた 1 DPU (4 vCPU、16 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケーラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。
- G.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 77 GB の空き容量) を備えた 2 DPU (8 vCPU、32 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。データ変換、結合、クエリなどのワークロードには、ほとんどのジョブを実行するためのスケーラブルで費用対効果の高い方法として、このワーカータイプをお勧めします。

- G.4X ワーカータイプでは、各ワーカーは 256 GB のディスク (約 235 GB の空き容量) を備えた 4 DPU (16 vCPU、64 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、米国東部 (オハイオ)、米国東部 (バージニア北部)、米国西部 (オレゴン)、アジアパシフィック (シンガポール)、アジアパシフィック (シドニー)、アジアパシフィック (東京)、カナダ (中部)、欧州 (フランクフルト)、欧州 (アイルランド)、欧州 (ストックホルム) の各 AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL ジョブでのみ使用できます。
- G.8X ワーカータイプでは、各ワーカーは 512 GB のディスク (約 487 GB の空き容量) を備えた 8 DPU (32 vCPU、128 GB のメモリ) にマッピングされており、ワーカーごとに 1 つのエグゼキューターを提供します。ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブには、このワーカータイプをお勧めします。このワーカータイプは、ワーカータイプでサポートされているのと同じ AWS リージョンで、AWS Glue バージョン 3.0 以降の Spark ETL G.4X ジョブでのみ使用できます。
- Z.2X ワーカータイプでは、各ワーカーは 128 GB のディスク (約 120 GB の空き容量) を備えた 2 M-DPU (8 vCPU、64 GB のメモリ) にマッピングされており、オートスケーラーに基づき最大 8 個の Ray ワーカーを提供します。
- SecurityConfiguration – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

セッションで使用する SecurityConfiguration 構造の名前

- GlueVersion – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

AWS Glue バージョンによって、が AWS Glue サポートする Apache Spark および Python のバージョンが決まります。は 2.0 より大きい GlueVersion 必要があります。

- DataAccessId – UTF-8 文字列、1 ~ 36 バイト。

セッションのデータアクセス ID。

- PartitionId – UTF-8 文字列、1 ~ 36 バイト。

セッションのパーティション ID。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

セッションに属するキーバリューペアのマップ (タグ)。

- RequestOrigin – UTF-8 文字列、1~128 バイト長、[「Custom string pattern #27」](#) に一致。

リクエストの送信元。

- ProfileName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

セッションに関連付けられた AWS Glue 使用プロファイルの名前。

## レスポンス

- Session – [セッション](#) オブジェクト。

レスポンス内のセッションオブジェクトを返します。

## エラー

- AccessDeniedException
- IdempotentParameterMismatchException
- InternalServiceException
- OperationTimeoutException
- InvalidInputException
- ValidationException
- AlreadyExistsException
- ResourceNumberLimitExceededException

## StopSession アクション (Python: stop\_session)

セッションを停止します。

### リクエスト

- Id – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

停止するセッションの ID。

- RequestOrigin – UTF-8 文字列、1~128 バイト長、[「Custom string pattern #27」](#) に一致。

リクエストの送信元。

## レスポンス

- Id – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

停止したセッションの ID を返します。

## エラー

- AccessDeniedException
- InternalServiceException
- OperationTimeoutException
- InvalidInputException
- IllegalSessionStateException
- ConcurrentModificationException

## DeleteSession アクション (Python: delete\_session)

セッションを削除します。

### リクエスト

- Id – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

削除するセッションの ID。

- RequestOrigin – UTF-8 文字列、1~128 バイト長、「[Custom string pattern #27](#)」に一致。

セッション削除リクエスト送信元の名前。

### レスポンス

- Id – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

削除したセッションの ID を返します。

## エラー

- `AccessDeniedException`
- `InternalServiceException`
- `OperationTimeoutException`
- `InvalidInputException`
- `IllegalSessionStateException`
- `ConcurrentModificationException`

## GetSession アクション (Python: `get_session`)

セッションを取得します。

### リクエスト

- `Id` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

セッションの ID

- `RequestOrigin` – UTF-8 文字列、1~128 バイト長、[「Custom string pattern #27」](#) に一致。

リクエストの送信元。

### レスポンス

- `Session` – [セッション](#) オブジェクト。

セッションオブジェクトはレスポンスで返されます。

## エラー

- `AccessDeniedException`
- `EntityNotFoundException`
- `InternalServiceException`
- `OperationTimeoutException`
- `InvalidInputException`

## ListSessions アクション (Python: list\_sessions)

セッションのリストを取得します。

### リクエスト

- NextToken - UTF-8 文字列。400000 バイト長以下。

次の結果セットのトークン、または追加の結果がない場合は null。

- MaxResults - 1 ~ 1000 の数値 (整数)。

結果の最大数。

- Tags - キーと値のペアのマッピング配列。50 ペア以下。

各キーは UTF-8 文字列で、1 ~ 128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

セッションに属するタグ。

- RequestOrigin - UTF-8 文字列、1 ~ 128 バイト長、「[Custom string pattern #27](#)」に一致。

リクエストの送信元。

### レスポンス

- Ids - UTF-8 文字列の配列。

セッションの ID を返します。

- Sessions - [セッション](#) オブジェクトの配列。

セッションオブジェクトを返します。

- NextToken - UTF-8 文字列。400000 バイト長以下。

次の結果セットのトークン、または追加の結果がない場合は null。

### エラー

- AccessDeniedException
- InvalidInputException

- `InternalServiceException`
- `OperationTimeoutException`

## RunStatement アクション (Python: `run_statement`)

ステートメントを実行します。

### リクエスト

- `SessionId` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

実行するステートメントのセッション ID。

- `Code` – 必須: UTF-8 文字列。68,000 バイト長未満。

実行するステートメントのコード。

- `RequestOrigin` – UTF-8 文字列、1~128 バイト長、[「Custom string pattern #27」](#) に一致。

リクエストの送信元。

### レスポンス

- `Id` – 数値 (整数)。

実行されたステートメントの ID を返します。

### エラー

- `EntityNotFoundException`
- `AccessDeniedException`
- `InternalServiceException`
- `OperationTimeoutException`
- `InvalidInputException`
- `ValidationException`
- `ResourceNumberLimitExceededException`
- `IllegalSessionStateException`

## CancelStatement アクション (Python: cancel\_statement)

ステートメントをキャンセルします。

### リクエスト

- `SessionId` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
キャンセルするステートメントのセッション ID。
- `Id` – 必須: 数値 (integer)。  
キャンセルするステートメントの ID。
- `RequestOrigin` – UTF-8 文字列、1~128 バイト長、[「Custom string pattern #27」](#) に一致。  
ステートメントのキャンセルをリクエストした送信元。

### レスポンス

- 応答パラメータはありません。

### エラー

- `AccessDeniedException`
- `EntityNotFoundException`
- `InternalServiceException`
- `OperationTimeoutException`
- `InvalidInputException`
- `IllegalSessionStateException`

## GetStatement アクション (Python: get\_statement)

ステートメントを取得します。

### リクエスト

- `SessionId` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
ステートメントのセッション ID。

- Id – 必須: 数値 (integer)。  
ステートメントの ID。
- RequestOrigin – UTF-8 文字列、1~128 バイト長、「[Custom string pattern #27](#)」に一致。  
リクエストの送信元。

## レスポンス

- Statement – [Statement](#) オブジェクト。  
ステートメントを返します。

## エラー

- AccessDeniedException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException
- InvalidInputException
- IllegalSessionStateException

## ListStatements アクション (Python: list\_statements)

セッションのステートメントをリスト表示します。

### リクエスト

- SessionId – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
ステートメントのセッション ID。
- RequestOrigin – UTF-8 文字列、1~128 バイト長、「[Custom string pattern #27](#)」に一致。  
ステートメントのリスト表示をリクエストした送信元。
- NextToken - UTF-8 文字列。400000 バイト長以下。  
継続トークン (これが継続呼び出しの場合)。

## 応答

- Statements – [Statement](#) オブジェクトの配列。  
ステートメントのリストを返します。
- NextToken - UTF-8 文字列。400000 バイト長以下。  
一部のステートメントがまだ返されていない場合は、継続トークン。

## エラー

- AccessDeniedException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException
- InvalidInputException
- IllegalSessionStateException

## 開発エンドポイント API

開発エンドポイント API では、カスタム DevEndpoint を使用したテストに関連した AWS Glue API について説明します。

## データ型

- [DevEndpoint 構造](#)
- [DevEndpointCustomLibraries 構造](#)

## DevEndpoint 構造

開発者が (ETL) スクリプトをリモートでデバッグ、抽出、変換、ロードする開発エンドポイント。

### フィールド

- EndpointName – UTF-8 文字列。  
DevEndpoint の名前。

- RoleArn – UTF-8 文字列、「[AWS IAM ARN string pattern](#)」に一致。

この DevEndpoint で使用される IAM ロールの Amazon リソースネーム (ARN)。

- SecurityGroupIds – UTF-8 文字列の配列。

この DevEndpoint で使用されるセキュリティグループ識別子のリスト。

- SubnetId – UTF-8 文字列。

この DevEndpoint のサブネット ID。

- YarnEndpointAddress – UTF-8 文字列。

この DevEndpoint で使用される YARN エンドポイントアドレス。

- PrivateAddress – UTF-8 文字列。

DevEndpoint が VPC 内に作成されている場合は、VPC 内の DevEndpoint にアクセスするためのプライベート IP アドレス。PrivateAddress フィールドは、VPC 内に DevEndpoint を作成する場合にのみ表示されます。

- ZeppelinRemoteSparkInterpreterPort – 数値 (整数)。

リモート Apache Spark インタープリタの Apache Zeppelin ポート。

- PublicAddress – UTF-8 文字列。

この DevEndpoint で使用されるパブリック IP アドレス。PublicAddress フィールドは、非 Virtual Private Cloud (VPC) の DevEndpoint を作成する場合にのみ存在します。

- Status – UTF-8 文字列。

DevEndpoint の現在の状態。

- WorkerType – UTF-8 文字列 (有効な値: Standard="" | G.1X="" | G.2X="" | G.025X="" | G.4X="" | G.8X="" | Z.2X="")。

開発エンドポイントに割り当てられる事前定義されたワーカーのタイプ。使用できる値は、Standard、G.1X、または G.2X です。

- Standard ワーカータイプでは、各ワーカーは 4 vCPU、16 GB のメモリ、50 GB のディスク、ワーカーあたり 2 個のエグゼキューターを提供します。
- G.1X ワーカータイプでは、各ワーカーは 1 DPU (4 vCPU、16 GB のメモリ、64 GB のディスク) にマッピングされており、ワーカーごとに 1 個のエグゼキューターを提供します。メモリを大量に消費するジョブには、このワーカータイプをお勧めします。

- G.2X ワーカータイプでは、各ワーカーは 2 DPU (8 vCPU、32 GB のメモリ、128 GB のディスク) にマッピングされており、ワーカーごとに 1 個のエグゼキュターを提供します。メモリを大量に消費するジョブには、このワーカータイプをお勧めします。

既知の問題: G.2X WorkerType 設定で開発エンドポイントを作成すると、開発エンドポイントの Spark ドライバーは 4 vCPU、16 GB のメモリ、および 64 GB のディスクで実行されます。

- GlueVersion – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

Glue バージョンは、AWS Glue がサポートする Apache Spark と Python のバージョンを決定します。Python バージョンは、開発エンドポイントで ETL スクリプトを実行するためにサポートされているバージョンを示します。

利用可能な AWS Glue のバージョン、および対応する Spark および Python のバージョンの詳細については、デベロッパーガイドの「[Glue version](#)」を参照してください。

Glue バージョンを指定せずに作成された開発エンドポイントは、デフォルトで Glue 0.9 に設定されます。

CreateDevEndpoint または UpdateDevEndpoint API で Arguments パラメータを使用して、開発エンドポイントの Python サポートのバージョンを指定できます。引数を指定しない場合、バージョンはデフォルトで Python 2 になります。

- NumberOfWorkers – 数値 (整数)。

開発エンドポイントに割り当てられた、定義された workerType のワーカーの数。

定義可能なワーカーの最大数は、299 (G.1X) または 149 (G.2X) です。

- NumberOfNodes – 数値 (整数)。

この DevEndpoint に割り当てられた AWS Glue データ処理ユニット (DPU) の数。

- AvailabilityZone – UTF-8 文字列。

この DevEndpoint が配置されている AWS アベイラビリティーゾーン。

- VpcId – UTF-8 文字列。

この DevEndpoint によって使用される Virtual Private Cloud (VPC) の ID。

- ExtraPythonLibsS3Path – UTF-8 文字列。

DevEndpoint にロードする Amazon S3 バケットの 1 つ以上の Python ライブラリへのパス。複数の値はコンマで区切られた完全なパスでなければなりません。

**Note**

DevEndpoint を持つ純粋な Python ライブラリのみを使用することができます。[pandas](#) Python データ解析ライブラリなど、C の拡張に依存するライブラリは現在サポートされていません。

- ExtraJarsS3Path – UTF-8 文字列。

DevEndpoint にロードする S3 バケットの 1 つ以上の Java .jar へのパス。

**Note**

DevEndpoint を持つ純粋な Java/Scala ライブラリのみを使用することができます。

- FailureReason – UTF-8 文字列。

この DevEndpoint で現在障害が発生している原因。

- LastUpdateStatus – UTF-8 文字列。

最終更新のステータス。

- CreatedTimestamp – タイムスタンプ。

この DevEndpoint が作成された時点。

- LastModifiedTimestamp – タイムスタンプ。

この DevEndpoint が最後に変更された時点。

- PublicKey – UTF-8 文字列。

この DevEndpoint が認証に使用するパブリックキー。この属性は、使用が推奨される属性がパブリックキーであるため、下位互換性のために提供されています。

- PublicKeys - UTF-8 文字列の配列。文字列 5 個以下。

認証のために DevEndpoints により使用されるパブリックキーのリスト。パブリックキーを使用するとクライアントごとに異なるプライベートキーを設定できるため、この属性は単一のパブリックキーで使用するものが推奨されます。

**Note**

以前にパブリックキーを使用してエンドポイントを作成した場合は、パブリックキーのリストを設定できるようにそのキーを削除する必要があります。パブリックキーのコンテンツを `deletePublicKeys` 属性で、新しいキーのリストを `addPublicKeys` 属性で `UpdateDevEndpoint` API オペレーションを呼び出します。

- `SecurityConfiguration` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

この `DevEndpoint` で使用される `SecurityConfiguration` 構造の名前。

- `Arguments` – キーと値のペアのマップ配列。100 ペア以下。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

`DevEndpoint` を設定するために使用される引数のマップ。

有効な引数は次のとおりです。

- `"--enable-glue-datacatalog": ""`

`CreateDevEndpoint` または `UpdateDevEndpoint` API で `Arguments` パラメータを使用して、開発エンドポイントの Python サポートのバージョンを指定できます。引数を指定しない場合、バージョンはデフォルトで Python 2 になります。

## DevEndpointCustomLibraries 構造

開発エンドポイントにロードされるカスタムライブラリ。

フィールド

- `ExtraPythonLibsS3Path` – UTF-8 文字列。

`DevEndpoint` にロードする Amazon Simple Storage Service (Amazon S3) バケットの 1 つ以上の Python ライブラリへのパス。複数の値はコンマで区切られた完全なパスでなければなりません。

**Note**

DevEndpoint を持つ純粋な Python ライブラリのみを使用することができます。[pandas](#) Python データ解析ライブラリなど、C の拡張に依存するライブラリは現在サポートされていません。

- ExtraJarsS3Path – UTF-8 文字列。

DevEndpoint にロードする S3 バケットの 1 つ以上の Java .jar へのパス。

**Note**

DevEndpoint を持つ純粋な Java/Scala ライブラリのみを使用することができます。

## 操作

- [CreateDevEndpoint アクション \(Python: create\\_dev\\_endpoint\)](#)
- [UpdateDevEndpoint アクション \(Python: update\\_dev\\_endpoint\)](#)
- [DeleteDevEndpoint アクション \(Python: delete\\_dev\\_endpoint\)](#)
- [GetDevEndpoint アクション \(Python: get\\_dev\\_endpoint\)](#)
- [GetDevEndpoints アクション \(Python: get\\_dev\\_endpoints\)](#)
- [BatchGetDevEndpoints アクション \(Python: batch\\_get\\_dev\\_endpoints\)](#)
- [ListDevEndpoints アクション \(Python: list\\_dev\\_endpoints\)](#)

## CreateDevEndpoint アクション (Python: create\_dev\_endpoint)

新しい開発エンドポイントを作成します。

### リクエスト

- EndpointName – 必須: UTF-8 文字列。

新しい DevEndpoint に割り当てられる名前。

- RoleArn – 必須: UTF-8 文字列。[AWS IAM ARN string pattern](#) に一致。

DevEndpoint の IAM ロール。

- SecurityGroupIds – UTF-8 文字列の配列。

新しい DevEndpoint によって使用されるセキュリティグループのセキュリティグループ ID。

- SubnetId – UTF-8 文字列。

使用する新しい DevEndpoint のサブネット ID。

- PublicKey – UTF-8 文字列。

この DevEndpoint が認証に使用するパブリックキー。この属性は、使用が推奨される属性がパブリックキーであるため、下位互換性のために提供されています。

- PublicKeys - UTF-8 文字列の配列。文字列 5 個以下。

認証のために開発エンドポイントで使用されるパブリックキーのリスト。パブリックキーを使用するとクライアントごとに異なるプライベートキーを設定できるため、この属性は単一のパブリックキーで使用するものが推奨されます。

#### Note

以前にパブリックキーを使用してエンドポイントを作成した場合は、パブリックキーのリストを設定できるようにそのキーを削除する必要があります。パブリックキーのコンテンツを deletePublicKeys 属性で、新しいキーのリストを addPublicKeys 属性で指定して、UpdateDevEndpoint API オペレーションを呼び出します。

- NumberOfNodes – 数値 (整数)。

この DevEndpoint に割り当てる AWS Glue データ処理ユニット (DPU) の数。

- WorkerType – UTF-8 文字列 (有効な値: Standard="" | G.1X="" | G.2X="" | G.025X="" | G.4X="" | G.8X="" | Z.2X="")。

開発エンドポイントに割り当てられる事前定義されたワーカーのタイプ。使用できる値は、Standard、G.1X、または G.2X です。

- Standard ワーカータイプでは、各ワーカーは 4 vCPU、16 GB のメモリ、50 GB のディスク、ワーカーあたり 2 個のエグゼキュターを提供します。
- G.1X ワーカータイプでは、各ワーカーは 1 DPU (4 vCPU、16 GB のメモリ、64 GB のディスク) にマッピングされており、ワーカーごとに 1 個のエグゼキュターを提供します。メモリを大量に消費するジョブには、このワーカータイプをお勧めします。

- G.2X ワーカータイプでは、各ワーカーは 2 DPU (8 vCPU、32 GB のメモリ、128 GB のディスク) にマッピングされており、ワーカーごとに 1 個のエグゼキューターを提供します。メモリを大量に消費するジョブには、このワーカータイプをお勧めします。

既知の問題: G.2X WorkerType 設定で開発エンドポイントを作成すると、開発エンドポイントの Spark ドライバーは 4 vCPU、16 GB のメモリ、および 64 GB のディスクで実行されます。

- GlueVersion – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

Glue バージョンは、AWS Glue がサポートする Apache Spark と Python のバージョンを決定します。Python バージョンは、開発エンドポイントで ETL スクリプトを実行するためにサポートされているバージョンを示します。

利用可能な AWS Glue のバージョン、および対応する Spark および Python のバージョンの詳細については、デベロッパーガイドの「[Glue version](#)」を参照してください。

Glue バージョンを指定せずに作成された開発エンドポイントは、デフォルトで Glue 0.9 に設定されます。

CreateDevEndpoint または UpdateDevEndpoint API で Arguments パラメータを使用して、開発エンドポイントの Python サポートのバージョンを指定できます。引数を指定しない場合、バージョンはデフォルトで Python 2 になります。

- NumberOfWorkers – 数値 (整数)。

開発エンドポイントに割り当てられた、定義された workerType のワーカーの数。

定義可能なワーカーの最大数は、299 (G.1X) または 149 (G.2X) です。

- ExtraPythonLibsS3Path – UTF-8 文字列。

DevEndpoint にロードする Amazon S3 バケットの 1 つ以上の Python ライブラリへのパス。複数の値はコンマで区切られた完全なパスでなければなりません。

#### Note

DevEndpoint を持つ純粋な Python ライブラリのみを使用することができます。[pandas](#) Python データ解析ライブラリなど、C の拡張機能に依存するライブラリはまだサポートされていません。

- ExtraJarsS3Path – UTF-8 文字列。

DevEndpoint にロードする S3 バケットの 1 つ以上の Java .jar へのパス。

- SecurityConfiguration – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

この DevEndpoint で使用される SecurityConfiguration 構造の名前。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

この DevEndpoint で使用するタグ。DevEndpoint へのアクセスを制限するためにタグを使用することができます。AWS Glue のタグの詳細については、デベロッパーガイドの「[AWS Tags in AWS Glue](#)」を参照してください。

- Arguments – キーと値のペアのマップ配列。100 ペア以下。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

DevEndpoint を設定するために使用される引数のマップ。

## レスポンス

- EndpointName – UTF-8 文字列。

新しい DevEndpoint に割り当てられた名前。

- Status – UTF-8 文字列。

新しい DevEndpoint の現在の状態。

- SecurityGroupIds – UTF-8 文字列の配列。

新しい DevEndpoint に割り当てられたセキュリティグループ。

- SubnetId – UTF-8 文字列。

新しい DevEndpoint に割り当てられたサブネット ID。

- RoleArn – UTF-8 文字列、「[AWS IAM ARN string pattern](#)」に一致。

新しい DevEndpoint に割り当てられたロールの Amazon リソースネーム (ARN)。

- YarnEndpointAddress – UTF-8 文字列。

この DevEndpoint で使用される YARN エンドポイントのアドレス。

- ZeppelinRemoteSparkInterpreterPort – 数値 (整数)。

リモート Apache Spark インタープリタの Apache Zeppelin ポート。

- NumberOfNodes – 数値 (整数)。

この DevEndpoint に割り当てられた AWS Glue データ処理ユニット (DPU) の数。

- WorkerType – UTF-8 文字列 (有効な値: Standard="" | G.1X="" | G.2X="" | G.025X="" | G.4X="" | G.8X="" | Z.2X="")。

開発エンドポイントに割り当てられる事前定義されたワーカーのタイプ。使用できる値は、Standard、G.1X、または G.2X です。

- GlueVersion – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

Glue バージョンは、AWS Glue がサポートする Apache Spark と Python のバージョンを決定します。Python バージョンは、開発エンドポイントで ETL スクリプトを実行するためにサポートされているバージョンを示します。

利用可能な AWS Glue のバージョン、および対応する Spark および Python のバージョンの詳細については、デベロッパーガイドの「[Glue version](#)」を参照してください。

- NumberOfWorkers – 数値 (整数)。

開発エンドポイントに割り当てられた、定義された workerType のワーカーの数。

- AvailabilityZone – UTF-8 文字列。

この DevEndpoint が配置されている AWS アベイラビリティーゾーン。

- VpcId – UTF-8 文字列。

この DevEndpoint によって使用される Virtual Private Cloud (VPC) の ID。

- ExtraPythonLibsS3Path – UTF-8 文字列。

DevEndpoint にロードする S3 バケットの 1 つ以上の Python ライブラリへのパス。

- ExtraJarsS3Path – UTF-8 文字列。

DevEndpoint にロードする S3 バケットの 1 つ以上の Java .jar ファイルへのパス。

- FailureReason – UTF-8 文字列。

この DevEndpoint で現在障害が発生している原因。

- SecurityConfiguration – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

この DevEndpoint で使用されている SecurityConfiguration 構造の名前。

- CreatedTimestamp – タイムスタンプ。

この DevEndpoint が最後に変更された時点。

- Arguments – キーと値のペアのマップ配列。100 ペア以下。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

この DevEndpoint を設定するために使用される引数のマップ。

有効な引数は次のとおりです。

- "--enable-glue-datacatalog": ""

CreateDevEndpoint または UpdateDevEndpoint API で Arguments パラメータを使用して、開発エンドポイントの Python サポートのバージョンを指定できます。引数を指定しない場合、バージョンはデフォルトで Python 2 になります。

## エラー

- AccessDeniedException
- AlreadyExistsException
- IdempotentParameterMismatchException
- InternalServiceException
- OperationTimeoutException
- InvalidInputException
- ValidationException
- ResourceNumberLimitExceededException

## UpdateDevEndpoint アクション (Python: update\_dev\_endpoint)

指定された開発エンドポイントを更新します。

### リクエスト

- `EndpointName` – 必須: UTF-8 文字列。

更新される `DevEndpoint` の名前。

- `PublicKey` – UTF-8 文字列。

使用する `DevEndpoint` のパブリックキー。

- `AddPublicKeys` - UTF-8 文字列の配列。文字列 5 個以下。

使用する `DevEndpoint` のパブリックキーのリスト。

- `DeletePublicKeys` - UTF-8 文字列の配列。文字列 5 個以下。

`DevEndpoint` から削除するパブリックキーのリスト。

- `CustomLibraries` – [DevEndpointCustomLibraries](#) オブジェクト。

`DevEndpoint` でロードされるカスタム Python または Java ライブラリ。

- `UpdateEtlLibraries` – ブール。

開発エンドポイントでロードされるカスタムライブラリのリストを更新する必要がある場合は `True`、それ以外の場合は `False`。

- `DeleteArguments` – UTF-8 文字列の配列。

`DevEndpoint` の設定に使用される引数のマップから削除される引数キーのリスト。

- `AddArguments` – キーと値のペアのマップ配列。100 ペア以下。

各キーは UTF-8 文字列。

各値は UTF-8 文字列。

`DevEndpoint` の設定に使用される引数のマップを追加する引数のマップ。

有効な引数は次のとおりです。

- `"--enable-glue-datacatalog": ""`

CreateDevEndpoint または UpdateDevEndpoint API で Arguments パラメータを使用して、開発エンドポイントの Python サポートのバージョンを指定できます。引数を指定しない場合、バージョンはデフォルトで Python 2 になります。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException
- InvalidInputException
- ValidationException

## DeleteDevEndpoint アクション (Python: delete\_dev\_endpoint)

指定された開発エンドポイントを削除します。

## リクエスト

- EndpointName – 必須: UTF-8 文字列。

DevEndpoint の名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException

- `InvalidInputException`

## GetDevEndpoint アクション (Python: `get_dev_endpoint`)

指定した開発エンドポイントに関する情報を取得します。

### Note

Virtual Private Cloud (VPC) で開発エンドポイントを作成する場合、AWS Glue はプライベート IP アドレスのみを返します。パブリック IP アドレスフィールドは入力されません。VPC 以外の開発エンドポイントを作成する場合、AWS Glue はパブリック IP アドレスのみを返します。

### リクエスト

- `EndpointName` – 必須: UTF-8 文字列。

情報を取得する `DevEndpoint` の名前。

### レスポンス

- `DevEndpoint` – [DevEndpoint](#) オブジェクト。

`DevEndpoint` の定義。

### エラー

- `EntityNotFoundException`
- `InternalServiceException`
- `OperationTimeoutException`
- `InvalidInputException`

## GetDevEndpoints アクション (Python: `get_dev_endpoints`)

この AWS アカウントのすべての開発エンドポイントを取得します。

**Note**

Virtual Private Cloud (VPC) で開発エンドポイントを作成する場合、AWS Glue はプライベート IP アドレスのみを返します。パブリック IP アドレスフィールドには入力されません。VPC 以外の開発エンドポイントを作成する場合、AWS Glue はパブリック IP アドレスのみを返します。

## リクエスト

- `MaxResults` – 1 ~ 1000 の数値 (整数)。

返される情報の最大サイズ。

- `NextToken` – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

## 応答

- `DevEndpoints` – [DevEndpoint](#) オブジェクトの配列。

`DevEndpoint` の定義のリスト。

- `NextToken` – UTF-8 文字列。

一部の `DevEndpoint` 定義がまだ返されていない場合は、継続トークン。

## エラー

- `EntityNotFoundException`
- `InternalServiceException`
- `OperationTimeoutException`
- `InvalidInputException`

## BatchGetDevEndpoints アクション (Python: `batch_get_dev_endpoints`)

指定された開発エンドポイントのリストのリソースメタデータのリストを返します。 `ListDevEndpoints` オペレーションを呼び出した後で、このオペレーションを呼び出すこと

で、アクセス許可が付与されているデータにアクセスできます。このオペレーションは、タグを使用するアクセス許可条件を含め、すべての IAM のアクセス許可をサポートします。

## リクエスト

- `customerAccountId` – UTF-8 文字列。

AWS アカウント ID。

- `DevEndpointNames` – 必須: UTF-8 文字列の配列。1 ~ 25 個の文字列。

`DevEndpoint` 名のリスト。これは `ListDevEndpoint` オペレーションから返される名前です。

## 応答

- `DevEndpoints` – [DevEndpoint](#) オブジェクトの配列。

`DevEndpoint` の定義のリスト。

- `DevEndpointsNotFound` – UTF-8 文字列の配列。1 ~ 25 個の文字列。

`DevEndpoints` のリストが見つかりません。

## エラー

- `AccessDeniedException`
- `InternalServiceException`
- `OperationTimeoutException`
- `InvalidInputException`

## ListDevEndpoints アクション (Python: `list_dev_endpoints`)

この AWS アカウントのすべての `DevEndpoint` リソースまたは指定されたタグを持つリソースの名前を取得します。このオペレーションにより、アカウントで利用可能なリソースとその名前を確認できます。

このオペレーションはオプションの `Tags` フィールドを受け取ります。このフィールドを応答のフィルターとして使用すると、タグ付きリソースをグループとして取得できます。タグフィルタリングの使用を選択した場合は、タグが付いたリソースのみが取得されます。

## リクエスト

- NextToken – UTF-8 文字列。

継続トークン (これが継続リクエストの場合)。

- MaxResults – 1 ~ 1000 の数値 (整数)。

返されるリストの最大サイズ。

- Tags – キーと値のペアのマッピング配列。50 ペア以下。

各キーは UTF-8 文字列で、1 ~ 128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

これらのタグ付きリソースのみを返すように指定します。

## レスポンス

- DevEndpointNames – UTF-8 文字列の配列。

アカウント内のすべての DevEndpoint の名前、または指定されたタグを持つ DevEndpoint の名前。

- NextToken – UTF-8 文字列。

継続トークン (戻されたリストに最後に使用可能なメトリクスが含まれていない場合)。

## エラー

- InvalidInputException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException

## スキーマレジストリ

スキーマレジストリ API は、でのスキーマの操作に関連するデータ型と API について説明します AWS Glue。

## データ型

- [RegistryId 構造](#)
- [RegistryListItem 構造](#)
- [MetadataInfo 構造](#)
- [OtherMetadataValueListItem 構造](#)
- [SchemaListItem 構造](#)
- [SchemaVersionListItem 構造](#)
- [MetadataKeyValuePair 構造](#)
- [SchemaVersionErrorItem 構造](#)
- [ErrorDetails 構造](#)
- [SchemaVersionNumber 構造](#)
- [Schemald 構造](#)

### RegistryId 構造

レジストリ名と Amazon リソースネーム (ARN) を含むことができるラッパー構造体。

#### フィールド

- RegistryName – UTF-8 文字列、1~255 バイト長、[Custom string pattern #18](#) に一致。

レジストリの名前。ルックアップにのみ使用されます。RegistryArn または RegistryName のいずれかを指定する必要があります。

- RegistryArn – UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。

更新するレジストリの ARN。RegistryArn または RegistryName のいずれかを指定する必要があります。

### RegistryListItem 構造

レジストリの詳細を含む構造。

#### フィールド

- RegistryName – UTF-8 文字列、1~255 バイト長、[Custom string pattern #18](#) に一致。

レジストリの名前。

- RegistryArn – UTF-8 文字列。1～10240 バイト長。[Custom string pattern #22](#) に一致。

レジストリの Amazon リソースネーム (ARN)

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

レジストリの説明。

- Status – UTF-8 文字列 (有効な値: AVAILABLE | DELETING)。

レジストリのステータス。

- CreatedTime – UTF-8 文字列。

レジストリが作成されたデータ。

- UpdatedTime – UTF-8 文字列。

レジストリが更新された日付。

## MetadataInfo 構造

スキーマバージョンのメタデータ情報を含む構造。

フィールド

- MetadataValue – UTF-8 文字列。1～256 バイト長。[Custom string pattern #33](#) に一致。

メタデータキーの対応する値。

- CreatedTime – UTF-8 文字列。

エントリの作成時刻。

- OtherMetadataValueList – [OtherMetadataValueListItem](#) オブジェクトの配列。

同じメタデータキーに属する他のメタデータ。

## OtherMetadataValueListItem 構造

同じメタデータキーに属するスキーマバージョンの他のメタデータが含まれる構造。

## フィールド

- `MetadataValue` – UTF-8 文字列。1~256 バイト長。 [Custom string pattern #33](#) に一致。  
同じメタデータキーに属する他のメタデータ用の、それらに対応したメタデータキーの値。
- `CreatedTime` – UTF-8 文字列。  
エントリの作成時刻。

## SchemaListItem 構造

スキーマに対する最小限の詳細が含まれるオブジェクト。

### フィールド

- `RegistryName` – UTF-8 文字列、1~255 バイト長、 [Custom string pattern #18](#) に一致。  
スキーマが存在するレジストリの名前。
- `SchemaName` – UTF-8 文字列、1~255 バイト長、 [Custom string pattern #18](#) に一致。  
スキーマの名前。
- `SchemaArn` – UTF-8 文字列。1~10240 バイト長。 [Custom string pattern #22](#) に一致。  
スキーマの Amazon リソースネーム (ARN)。
- `Description` – 説明文字列、2048 バイト長以下、 [URI address multi-line string pattern](#) に一致。  
スキーマの説明。
- `SchemaStatus` – UTF-8 文字列 (有効な値: AVAILABLE | PENDING | DELETING)。  
スキーマのステータス。
- `CreatedTime` – UTF-8 文字列。  
スキーマが作成された日時。
- `UpdatedTime` – UTF-8 文字列。  
スキーマが更新された日時。

## SchemaVersionListItem 構造

スキーマのバージョンに関する詳細を含むオブジェクト。

フィールド

- `SchemaArn` – UTF-8 文字列。1 ~ 10240 バイト長。 [Custom string pattern #22](#) に一致。  
スキーマの Amazon リソースネーム (ARN)。
- `SchemaVersionId` – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。  
スキーマバージョンの一意の識別子。
- `VersionNumber` – 数値 (long)。1 ~ 100000。  
スキーマのバージョン番号。
- `Status` – UTF-8 文字列 (有効な値: AVAILABLE | PENDING | FAILURE | DELETING)。  
スキーマバージョンのステータス。
- `CreatedTime` – UTF-8 文字列。  
スキーマバージョンが作成された日時。

## MetadataKeyValuePair 構造

メタデータのキーバリューペアが含まれる構造。

フィールド

- `MetadataKey` – UTF-8 文字列。1 ~ 128 バイト長。 [Custom string pattern #33](#) に一致。  
メタデータキー。
- `MetadataValue` – UTF-8 文字列。1 ~ 256 バイト長。 [Custom string pattern #33](#) に一致。  
メタデータキーの対応する値。

## SchemaVersionErrorItem 構造

スキーマバージョンに対するオペレーションに対するエラーの詳細が含まれるオブジェクト。

## フィールド

- `VersionNumber` – 数値 (long)。1 ~ 100000。  
スキーマのバージョン番号。
- `ErrorDetails` – [ErrorDetails](#) オブジェクト。  
スキーマバージョンのエラーの詳細。

## ErrorDetails 構造

エラーの詳細を含むオブジェクト。

### フィールド

- `ErrorCode` – UTF-8 文字列。  
エラーのエラーコード。
- `ErrorMessage` – UTF-8 文字列。  
エラーのエラーメッセージ。

## SchemaVersionNumber 構造

スキーマのバージョン情報を含む構造。

### フィールド

- `LatestVersion` – ブール。  
スキーマで利用可能な最新バージョン。
- `VersionNumber` – 数値 (long)。1 ~ 100000。  
スキーマのバージョン番号。

## Schemald 構造

スキーマレジストリ内の AWS Glue スキーマの一意的 ID。

## フィールド

- SchemaArn – UTF-8 文字列。1～10240 バイト長。[Custom string pattern #22](#) に一致。

スキーマの Amazon リソースネーム (ARN)。SchemaArn または SchemaName のいずれかを指定する必要があります。

- SchemaName – UTF-8 文字列、1～255 バイト長、[Custom string pattern #18](#) に一致。

スキーマの名前。SchemaArn または SchemaName のいずれかを指定する必要があります。

- RegistryName – UTF-8 文字列、1～255 バイト長、[Custom string pattern #18](#) に一致。

スキーマを含むスキーマレジストリの名前。

## 操作

- [CreateRegistry アクション](#) (Python: `create_registry`)
- [CreateSchema アクション](#) (Python: `create_schema`)
- [GetSchema アクション](#) (Python: `get_schema`)
- [ListSchemaVersions アクション](#) (Python: `list_schema_versions`)
- [GetSchemaVersion アクション](#) (Python: `get_schema_version`)
- [GetSchemaVersionsDiff アクション](#) (Python: `get_schema_versions_diff`)
- [ListRegistries アクション](#) (Python: `list_registries`)
- [ListSchemas アクション](#) (Python: `list_schemas`)
- [RegisterSchemaVersion アクション](#) (Python: `register_schema_version`)
- [UpdateSchema アクション](#) (Python: `update_schema`)
- [CheckSchemaVersionValidity アクション](#) (Python: `check_schema_version_validity`)
- [UpdateRegistry アクション](#) (Python: `update_registry`)
- [GetSchemaByDefinition アクション](#) (Python: `get_schema_by_definition`)
- [GetRegistry アクション](#) (Python: `get_registry`)
- [PutSchemaVersionMetadata アクション](#) (Python: `put_schema_version_metadata`)
- [QuerySchemaVersionMetadata アクション](#) (Python: `query_schema_version_metadata`)
- [RemoveSchemaVersionMetadata アクション](#) (Python: `remove_schema_version_metadata`)

- [DeleteRegistry アクション \(Python: delete\\_registry\)](#)
- [DeleteSchema アクション \(Python: delete\\_schema\)](#)
- [DeleteSchemaVersions アクション \(Python: delete\\_schema\\_versions\)](#)

## CreateRegistry アクション (Python: create\_registry)

スキーマのコレクションを保持するために使用できる新しいレジストリを作成します。

### リクエスト

- RegistryName – 必須: UTF-8 文字列、1~255 バイト長、「[Custom string pattern #18](#)」に一致。  
作成するスキーマの名前の最大長は 255 で、英字、数字、ハイフン、アンダースコア、ドル記号、ハッシュ記号のみを使用できます。空白はありません。
- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
レジストリの説明。説明がない場合、これに対するデフォルト値はありません。
- Tags – キーと値のペアのマップ配列。50 ペア以下。  
各キーは UTF-8 文字列で、1~128 バイト長です。  
各値は UTF-8 文字列で、256 バイト長以下です。  
AWS キーと値のペアを含み、コンソール、コマンドライン、または API で検索できるタグ。

### レスポンス

- RegistryArn – UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。  
新たに作成されたレジストリの Amazon リソースネーム (ARN)。
- RegistryName – UTF-8 文字列、1~255 バイト長、[Custom string pattern #18](#) に一致。  
レジストリの名前。
- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
レジストリの説明。
- Tags – キーと値のペアのマップ配列。50 ペア以下。  
各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

レジストリのタグ。

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `AlreadyExistsException`
- `ResourceNumberLimitExceededException`
- `ConcurrentModificationException`
- `InternalServiceException`

## CreateSchema アクション (Python: `create_schema`)

新しいスキーマセットを作成し、スキーマ定義を登録します。スキーマセットがすでに存在する場合は、実際にバージョンを登録せずにエラーを返します。

スキーマセットが作成されると、バージョンチェックポイントが最初のバージョンに設定されます。互換モードが「DISABLED」であると、最初のスキーマバージョンの後に追加のスキーマバージョンが追加されないように制限されます。他のすべての互換モードでは、`RegisterSchemaVersion` API が使用された場合、2 番目以降のバージョンでのみ互換性設定の検証が適用されます。

この API が `RegistryId` なしで呼び出されたとき、レジストリデータベーステーブルに「default-registry」のエントリがまだ存在しなければ、エントリが作成されます。

## リクエスト

- `RegistryId` – [RegistryId](#) オブジェクト。

これは、レジストリ ID フィールドを含むラッパー図形です。これを指定しない場合、デフォルトのレジストリが使用されます。同じものの ARN 形式は次のようになります: `arn:aws:glue:us-east-2:<customer id>:registry/default-registry:random-5-letter-id`。

- `SchemaName` – 必須: UTF-8 文字列、1~255 バイト長、[Custom string pattern #18](#) に一致。

作成するスキーマの名前の最大長は 255 で、文字、数字、ハイフン、アンダースコア、ドル記号、またはハッシュ記号のみを使用できます。空白はありません。

- `DataFormat` – 必須: UTF-8 文字列 (有効な値: AVRO | JSON | PROTOBUF)。

スキーマ定義のデータ形式。現在サポートされている形式は、AVRO、JSON、および PROTOBUF です。

- `Compatibility` – UTF-8 文字列 (有効な値: NONE | DISABLED | BACKWARD | BACKWARD\_ALL | FORWARD | FORWARD\_ALL | FULL | FULL\_ALL)。

スキーマの互換モード。指定できる値は以下のとおりです。

- `NONE`: 互換モードは適用されません。この選択は、開発シナリオで使用するか、スキーマに適用する互換モードがわからない場合に使用できます。新しいバージョンが追加されたとき、互換性チェックを受けずに受け入れられます。
- `DISABLED`: 互換性についてこの選択をすると、特定のスキーマのバージョン管理が不要になります。この選択を使用して、スキーマの今後のバージョン管理を不要にできます。
- `BACKWARD`: データ受信者が現在のスキーマバージョンと直前のスキーマバージョンの両方を読み取ることができるため、この互換性の選択をお勧めします。つまり、例えば、以前のバージョンを使用して読み取れなくなならないように、新しいスキーマバージョンでデータフィールドを削除したり、これらのフィールドのタイプを変更したりできなくなります。
- `BACKWARD_ALL`: この互換性の選択により、データ受信者は、現在および以前のすべてのスキーマバージョンの両方を読み取ることができます。この選択は、フィールドを削除したり、オプションのフィールドを追加したり、以前のすべてのスキーマバージョンとの互換性を確認する必要がある場合に使用できます。
- `FORWARD`: この互換性の選択により、データ受信者は、現在のスキーマバージョンとすぐ次のスキーマバージョンの両方を読み取ることができますが、それより後のバージョンは必ずしも読み取れません。この選択は、フィールドを追加したり、オプションフィールドを削除したりする必要があり、最後のスキーマバージョンとの互換性だけをチェックする必要がある場合に使用できます。
- `FORWARD_ALL`: この互換性の選択により、データ受信者は、新しく登録されるすべてのスキーマのプロデューサーによって書き込まれたデータを読み取ることができます。この選択は、フィールドを追加したり、オプションフィールドを削除したり、以前のすべてのスキーマバージョンとの互換性を確認したりする必要がある場合に使用できます。
- `FULL`: この互換性の選択により、データ受信者は、スキーマの直前または直後のバージョンを使用してプロデューサーによって書き込まれたデータを読み取ることができますが、それより前または後のバージョンは必ずしも読み取れません。この選択は、オプションフィールドを追加または削除する必要があり、最後のスキーマバージョンとの互換性だけをチェックする必要がある場合に使用できます。

- FULL\_ALL: この互換性の選択により、データ受信者は、すべての以前のスキーマバージョンを使用してプロデューサーによって書き込まれたデータを読み取ることができます。この選択は、オプションのフィールドを追加または削除したり、以前のすべてのスキーマバージョンとの互換性を確認したりする必要がある場合に使用できます。
- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
スキーマの説明 (省略可能)。説明がない場合、これに自動的に適用されるデフォルト値はありません。
- Tags – キーと値のペアのマップ配列。50 ペア以下。  
各キーは UTF-8 文字列で、1~128 バイト長です。  
各値は UTF-8 文字列で、256 バイト長以下です。  
AWS キーと値のペアを含み、コンソール、コマンドライン、または API で検索できるタグ。指定した場合、AWS tags-on-create パターンに従います。
- SchemaDefinition – UTF-8 文字列。1~170000 バイト長。[Custom string pattern #32](#) に一致。  
SchemaName の DataFormat 設定を使用したスキーマ定義。

## レスポンス

- RegistryName – UTF-8 文字列、1~255 バイト長、「[Custom string pattern #18](#)」に一致。  
レジストリの名前。
- RegistryArn – UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。  
レジストリの Amazon リソースネーム (ARN)
- SchemaName – UTF-8 文字列、1~255 バイト長、[Custom string pattern #18](#) に一致。  
スキーマの名前。
- SchemaArn – UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。  
スキーマの Amazon リソースネーム (ARN)。
- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
スキーマの説明 (作成時に指定されている場合)。
- DataFormat – UTF-8 文字列 (有効な値: AVRO | JSON | PROTOBUF)。

スキーマ定義のデータ形式。現在サポートされている形式は、AVRO、JSON、および PROTOBUF です。

- `Compatibility` – UTF-8 文字列 (有効な値: NONE | DISABLED | BACKWARD | BACKWARD\_ALL | FORWARD | FORWARD\_ALL | FULL | FULL\_ALL)。

スキーマの互換モード。

- `SchemaCheckpoint` – 数値 (long)。1 ~ 100000。

チェックポイントのバージョン番号 (互換モードが最後に変更されたときの変更回数)。

- `LatestSchemaVersion` – 数値 (long)。1 ~ 100000。

返されるスキーマ定義に関連付けられたスキーマの最新バージョン。

- `NextSchemaVersion` – 数値 (long)。1 ~ 100000。

返されたスキーマ定義に関連付けられたスキーマの次のバージョン。

- `SchemaStatus` – UTF-8 文字列 (有効な値: AVAILABLE | PENDING | DELETING)。

スキーマのステータス。

- `Tags` – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1 ~ 128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

スキーマのタグ。

- `SchemaVersionId` – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。

最初のスキーマバージョンの一意の識別子。

- `SchemaVersionStatus` – UTF-8 文字列 (有効な値: AVAILABLE | PENDING | FAILURE | DELETING)。

最初に作成されたスキーマバージョンのステータス。

## エラー

- `InvalidInputException`
- `AccessDeniedException`

- EntityNotFoundException
- AlreadyExistsException
- ResourceNumberLimitExceededException
- ConcurrentModificationException
- InternalServiceException

## GetSchema アクション (Python: get\_schema)

指定されたスキーマについて詳細に説明します。

### リクエスト

- SchemaId – 必須: [Schemald](#) オブジェクト。

スキーマ ID フィールドを含むラッパー構造体です。構造体には以下が含まれています。

- Schemald\$SchemaArn: スキーマの Amazon リソースネーム (ARN)。SchemaArn または SchemaName、および RegistryName を指定する必要があります。
- Schemald\$SchemaName: スキーマの名前。SchemaArn または SchemaName、および RegistryName を指定する必要があります。

### レスポンス

- RegistryName – UTF-8 文字列、1~255 バイト長、「[Custom string pattern #18](#)」に一致。

レジストリの名前。

- RegistryArn – UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。

レジストリの Amazon リソースネーム (ARN)

- SchemaName – UTF-8 文字列、1~255 バイト長、[Custom string pattern #18](#) に一致。

スキーマの名前。

- SchemaArn – UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。

スキーマの Amazon リソースネーム (ARN)。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

スキーマの説明 (作成時に指定されている場合)。

- `DataFormat` – UTF-8 文字列 (有効な値: AVRO | JSON | PROTOBUF)。

スキーマ定義のデータ形式。現在サポートされている形式は、AVRO、JSON、および PROTOBUF です。

- `Compatibility` – UTF-8 文字列 (有効な値: NONE | DISABLED | BACKWARD | BACKWARD\_ALL | FORWARD | FORWARD\_ALL | FULL | FULL\_ALL)。

スキーマの互換モード。

- `SchemaCheckpoint` – 数値 (long)。1 ~ 100000。

チェックポイントのバージョン番号 (互換モードが最後に変更されたときの変更回数)。

- `LatestSchemaVersion` – 数値 (long)。1 ~ 100000。

返されるスキーマ定義に関連付けられたスキーマの最新バージョン。

- `NextSchemaVersion` – 数値 (long)。1 ~ 100000。

返されたスキーマ定義に関連付けられたスキーマの次のバージョン。

- `SchemaStatus` – UTF-8 文字列 (有効な値: AVAILABLE | PENDING | DELETING)。

スキーマのステータス。

- `CreatedTime` – UTF-8 文字列。

スキーマが作成された日時。

- `UpdatedTime` – UTF-8 文字列。

スキーマが更新された日時。

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `EntityNotFoundException`
- `InternalServiceException`

## ListSchemaVersions アクション (Python: list\_schema\_versions)

作成したスキーマバージョンのリストを、最小限の情報とともに返します。Deleted ステータスのスキーマのバージョンは、結果に含まれません。使用可能なスキーマバージョンがない場合、空の結果が返されます。

### リクエスト

- SchemaId – 必須: [Schemald](#) オブジェクト。

スキーマ ID フィールドを含むラッパー構造体です。構造体には以下が含まれています。

- Schemald\$SchemaArn: スキーマの Amazon リソースネーム (ARN)。SchemaArn または SchemaName、および RegistryName を指定する必要があります。
- Schemald\$SchemaName: スキーマの名前。SchemaArn または SchemaName、および RegistryName を指定する必要があります。
- MaxResults – 数値 (integer)。1 ~ 100。

1 ページあたりに必要な結果の最大数。値が指定されていない場合、これはデフォルトで 1 ページあたり 25 に設定されます。

- NextToken – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

### 応答

- Schemas – [SchemaVersionListItem](#) オブジェクトの配列。

各スキーマのバージョンの詳細を含む SchemaVersionList オブジェクトの配列。

- NextToken – UTF-8 文字列。

返されたトークンのリストをページ分割するための継続トークン。リストの現在のセグメントが最後のセグメントではない場合に返されます。

### エラー

- InvalidInputException
- AccessDeniedException
- EntityNotFoundException

- `InternalServiceException`

## GetSchemaVersion アクション (Python: `get_schema_version`)

スキーマのバージョンの作成または登録時に割り当てられた一意の ID によって指定されたスキーマを取得します。Deleted ステータスのスキーマのバージョンは、結果に含まれません。

### リクエスト

- `SchemaId` – [Schemald](#) オブジェクト。

スキーマ ID フィールドを含むラッパー構造体です。構造体には以下が含まれています。

- `Schemald$SchemaArn`: スキーマの Amazon リソースネーム (ARN)。SchemaArn または SchemaName、および RegistryName を指定する必要があります。
- `Schemald$SchemaName`: スキーマの名前。SchemaArn または SchemaName、および RegistryName を指定する必要があります。
- `SchemaVersionId` – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。

スキーマバージョンの SchemaVersionId。このフィールドは、スキーマ ID による取得に必須です。このスキーマ ID が SchemaId ラッパーのいずれかを指定する必要があります。

- `SchemaVersionNumber` – [SchemaVersionNumber](#) オブジェクト。

スキーマのバージョン番号。

### レスポンス

- `SchemaVersionId` – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。

スキーマバージョンの SchemaVersionId。

- `SchemaDefinition` – UTF-8 文字列。1~170000 バイト長。 [Custom string pattern #32](#) に一致。

スキーマ ID のスキーマ定義。

- `DataFormat` – UTF-8 文字列 (有効な値: AVRO | JSON | PROTOBUF)。

スキーマ定義のデータ形式。現在サポートされている形式は、AVRO、JSON、および PROTOBUF です。

- `SchemaArn` – UTF-8 文字列。1~10240 バイト長。 [Custom string pattern #22](#) に一致。

スキーマの Amazon リソースネーム (ARN)。

- `VersionNumber` – 数値 (long)。1 ~ 100000。

スキーマのバージョン番号。

- `Status` – UTF-8 文字列 (有効な値: AVAILABLE | PENDING | FAILURE | DELETING)。

スキーマバージョンのステータス。

- `CreatedTime` – UTF-8 文字列。

スキーマバージョンが作成された日時。

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `EntityNotFoundException`
- `InternalServiceException`

## GetSchemaVersionsDiff アクション (Python: `get_schema_versions_diff`)

スキーマレジストリに格納された 2 つのスキーマバージョン間で、指定された差分タイプによるスキーマバージョンの差分を取得します。

この API を使用すると、同じスキーマの下にある 2 つのスキーマ定義間で 2 つのスキーマバージョンを比較できます。

### リクエスト

- `SchemaId` – 必須: [Schemald](#) オブジェクト。

スキーマ ID フィールドを含むラッパー構造体です。構造体には以下が含まれています。

- `Schemald$SchemaArn`: スキーマの Amazon リソースネーム (ARN)。SchemaArn または SchemaName のいずれかを指定する必要があります。
- `Schemald$SchemaName`: スキーマの名前。SchemaArn または SchemaName のいずれかを指定する必要があります。
- `FirstSchemaVersionNumber` – 必須: [SchemaVersionNumber](#) オブジェクト。

比較する 2 つのスキーマバージョンのうち、1 番目のスキーマバージョンです。

- SecondSchemaVersionNumber – 必須: [SchemaVersionNumber](#) オブジェクト。

比較する 2 つのスキーマバージョンのうち、2 番目のスキーマバージョンです。

- SchemaDiffType – 必須: UTF-8 文字列 (有効な値: SYNTAX\_DIFF)。

現在サポートされている差分タイプである SYNTAX\_DIFF が参照されます。

## レスポンス

- Diff – UTF-8 文字列。1 ~ 340000 バイト長。 [Custom string pattern #32](#) に一致。

JsonPatch 形式の文字列としてのスキーマの違い。

## エラー

- InvalidInputException
- EntityNotFoundException
- AccessDeniedException
- InternalServiceException

## ListRegistries アクション (Python: list\_registries)

作成したレジストリの一覧を、最小限のレジストリ情報とともに返します。Deleting ステータスのレジストリは結果に含まれません。利用可能なレジストリがない場合、空の結果が返されます。

### リクエスト

- MaxResults – 数値 (integer)。1 ~ 100。

1 ページあたりに必要な結果の最大数。値が指定されていない場合、これはデフォルトで 1 ページあたり 25 に設定されます。

- NextToken – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

## 応答

- Registries – [RegistryListItem](#) オブジェクトの配列。

各レジストリの最小の詳細を含む RegistryDetailedListItem オブジェクトの配列。

- NextToken – UTF-8 文字列。

返されたトークンのリストをページ分割するための継続トークン。リストの現在のセグメントが最後のセグメントではない場合に返されます。

## エラー

- InvalidInputException
- AccessDeniedException
- InternalServiceException

## ListSchemas アクション (Python: list\_schemas)

スキーマのリストを最小限の詳細とともに返します。Deleting ステータスのスキーマは結果に含まれません。使用可能なスキーマがない場合、空の結果が返されます。

RegistryId が指定されていない場合、レジストリ間のすべてのスキーマは API レスポンスの一部になります。

## リクエスト

- RegistryId – [RegistryId](#) オブジェクト。

レジストリ名と Amazon リソースネーム (ARN) を含むことができるラッパー構造体。

- MaxResults – 数値 (integer)。1 ~ 100。

1 ページあたりに必要な結果の最大数。値が指定されていない場合、これはデフォルトで 1 ページあたり 25 に設定されます。

- NextToken – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

## 応答

- Schemas – [SchemaListItem](#) オブジェクトの配列。

各スキーマの詳細を含む SchemaListItem オブジェクトの配列。

- NextToken – UTF-8 文字列。

返されたトークンのリストをページ分割するための継続トークン。リストの現在のセグメントが最後のセグメントではない場合に返されます。

## エラー

- InvalidInputException
- AccessDeniedException
- EntityNotFoundException
- InternalServiceException

## RegisterSchemaVersion アクション (Python: register\_schema\_version)

既存のスキーマに新しいバージョンを追加します。スキーマの新しいバージョンがスキーマセットの互換性要件を満たしていない場合、エラーを返します。この API では、スキーマセットがまだスキーマレジストリに存在しない場合、新しいスキーマセットは作成されず、404 エラーが返されません。

これがスキーマレジストリに登録される最初のスキーマ定義である場合、この API はスキーマバージョンを保存し、すぐに戻ります。それ以外の場合、この呼び出しは、互換モードのために他のオペレーションよりも長時間実行される可能性があります。SchemaVersionId を使用して GetSchemaVersion API を呼び出して、互換モードをチェックできます。

同じスキーマ定義がすでにバージョンとして Schema Registry に格納されている場合は、既存のスキーマのスキーマ ID が呼び出し元に返されます。

## リクエスト

- SchemaId – 必須: [Schemald](#) オブジェクト。

スキーマ ID フィールドを含むラッパー構造体です。構造体には以下が含まれています。

- `SchemaId$SchemaArn`: スキーマの Amazon リソースネーム (ARN)。SchemaArn または SchemaName、および RegistryName を指定する必要があります。
- `SchemaId$SchemaName`: スキーマの名前。SchemaArn または SchemaName、および RegistryName を指定する必要があります。
- `SchemaDefinition` – 必須: UTF-8 文字列。1~170000 バイト長。 [Custom string pattern #32](#) に一致。

SchemaName の DataFormat 設定を使用したスキーマ定義。

## レスポンス

- `SchemaVersionId` – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。

このスキーマのバージョンを表す一意の ID。

- `VersionNumber` – 数値 (long)。1~100000。

このスキーマのバージョン (これが最初のバージョンの場合、同期フローのみ)。

- `Status` – UTF-8 文字列 (有効な値: AVAILABLE | PENDING | FAILURE | DELETING)。

スキーマバージョンのステータス。

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `EntityNotFoundException`
- `ResourceNumberLimitExceededException`
- `ConcurrentModificationException`
- `InternalServiceException`

## UpdateSchema アクション (Python: update\_schema)

スキーマセットの説明、互換性設定、またはバージョンチェックポイントを更新します。

互換性設定を更新する場合、その呼び出しでは、新しい互換性設定を持つスキーマバージョンのセット全体の互換性は検証されません。Compatibility の値が指定されている場合、VersionNumber

(チェックポイント) も指定する必要があります。APIは、整合性のためにチェックポイントのバージョン番号を検証します。

VersionNumber (チェックポイント) の値が指定されている場合、Compatibility はオプションであり、スキーマのチェックポイントを設定/リセットするために使用できます。

この更新は、スキーマが AVAILABLE 状態である場合にのみ発生します。

## リクエスト

- SchemaId – 必須: [SchemaId](#) オブジェクト。

スキーマ ID フィールドを含むラッパー構造体です。構造体には以下が含まれています。

- SchemaId\$SchemaArn: スキーマの Amazon リソースネーム (ARN)。SchemaArn または SchemaName のいずれかを指定する必要があります。
- SchemaId\$SchemaName: スキーマの名前。SchemaArn または SchemaName のいずれかを指定する必要があります。
- SchemaVersionNumber – [SchemaVersionNumber](#) オブジェクト。

チェックポイントに必要なバージョン番号。VersionNumber または Compatibility のいずれかを指定する必要があります。

- Compatibility – UTF-8 文字列 (有効な値: NONE | DISABLED | BACKWARD | BACKWARD\_ALL | FORWARD | FORWARD\_ALL | FULL | FULL\_ALL)。

スキーマの新しい互換性設定。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

スキーマの新しい説明。

## レスポンス

- SchemaArn – UTF-8 文字列。1 ~ 10240 バイト長。[Custom string pattern #22](#) に一致。

スキーマの Amazon リソースネーム (ARN)。

- SchemaName – UTF-8 文字列、1 ~ 255 バイト長、[Custom string pattern #18](#) に一致。

スキーマの名前。

- RegistryName – UTF-8 文字列、1 ~ 255 バイト長、[Custom string pattern #18](#) に一致。

スキーマを含むレジストリの名前。

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `EntityNotFoundException`
- `ConcurrentModificationException`
- `InternalServiceException`

## CheckSchemaVersionValidity アクション (Python: `check_schema_version_validity`)

指定されたスキーマを検証します。この呼び出しには副作用はなく、指定されたスキーマを使用して `DataFormat` を形式として使用して単に検証します。スキーマセット名を取らないため、互換性チェックは実行されません。

### リクエスト

- `DataFormat` – 必須: UTF-8 文字列 (有効な値: AVRO | JSON | PROTOBUF)。

スキーマ定義のデータ形式。現在サポートされている形式は、AVRO、JSON、および PROTOBUF です。

- `SchemaDefinition` – 必須: UTF-8 文字列。1~170000 バイト長。 [Custom string pattern #32](#) に一致。

検証する必要があるスキーマの定義。

### レスポンス

- `Valid` – ブール。

スキーマが有効な場合は `true` を返し、そうでない場合は `false` を返します。

- `Error` – UTF-8 文字列。1~5000 バイト長。

検証失敗のエラーメッセージ。

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `InternalServiceException`

## UpdateRegistry アクション (Python: `update_registry`)

スキーマのコレクションを保持するために使用される既存のレジストリを更新します。更新されたプロパティはレジストリに関連づけられており、レジストリ内のスキーマは変更されません。

### リクエスト

- `RegistryId` – 必須: [RegistryId](#) オブジェクト。

これは、レジストリ名と Amazon リソースネーム (ARN) を含む可能性があるラッパー構造です。

- `Description` – 必須: 説明文字列。2048 バイト長以下。 [URI address multi-line string pattern](#) に一致。

レジストリの説明。説明が指定されていない場合、このフィールドは更新されません。

### レスポンス

- `RegistryName` – UTF-8 文字列、1~255 バイト長、「[Custom string pattern #18](#)」に一致。

更新されたレジストリの名前。

- `RegistryArn` – UTF-8 文字列。1~10240 バイト長。 [Custom string pattern #22](#) に一致。

更新されたレジストリの Amazon リソースネーム (ARN)

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `EntityNotFoundException`
- `ConcurrentModificationException`
- `InternalServiceException`

## GetSchemaByDefinition アクション (Python: get\_schema\_by\_definition)

SchemaDefinition によってスキーマを取得します。スキーマ定義は、スキーマレジストリに送信され、正規化され、ハッシュされます。ハッシュが SchemaName または ARN (指定されていない場合はデフォルトのレジストリ) のスコープ内で一致する場合、そのスキーマのメタデータが返されます。それ以外の場合は、404 または NotFound エラーが返されます。Deleted ステータスのスキーマバージョンは結果に含まれません。

### リクエスト

- SchemaId – 必須: [Schemald](#) オブジェクト。

スキーマ ID フィールドを含むラッパー構造体です。構造体には以下が含まれています。

- Schemald\$SchemaArn: スキーマの Amazon リソースネーム (ARN)。SchemaArn または SchemaName のいずれかを指定する必要があります。
- Schemald\$SchemaName: スキーマの名前。SchemaArn または SchemaName のいずれかを指定する必要があります。
- SchemaDefinition – 必須: UTF-8 文字列。1~170000 バイト長。 [Custom string pattern #32](#) に一致。

スキーマの詳細が必要なスキーマの定義。

### レスポンス

- SchemaVersionId – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。

スキーマバージョンのスキーマ ID。

- SchemaArn – UTF-8 文字列。1~10240 バイト長。 [Custom string pattern #22](#) に一致。

スキーマの Amazon リソースネーム (ARN)。

- DataFormat – UTF-8 文字列 (有効な値: AVRO | JSON | PROTOBUF)。

スキーマ定義のデータ形式。現在サポートされている形式は、AVRO、JSON、および PROTOBUF です。

- Status – UTF-8 文字列 (有効な値: AVAILABLE | PENDING | FAILURE | DELETING)。

スキーマバージョンのステータス。

- CreatedTime – UTF-8 文字列。

スキーマが作成された日時。

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `EntityNotFoundException`
- `InternalServiceException`

## GetRegistry アクション (Python: `get_registry`)

指定されたレジストリの詳細について説明します。

### リクエスト

- `RegistryId` – 必須: [RegistryId](#) オブジェクト。

これは、レジストリ名と Amazon リソースネーム (ARN) を含む可能性があるラッパー構造です。

### レスポンス

- `RegistryName` – UTF-8 文字列、1~255 バイト長、[「Custom string pattern #18」](#) に一致。

レジストリの名前。

- `RegistryArn` – UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。

レジストリの Amazon リソースネーム (ARN)

- `Description` – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

レジストリの説明。

- `Status` – UTF-8 文字列 (有効な値: AVAILABLE | DELETING)。

レジストリのステータス。

- `CreatedTime` – UTF-8 文字列。

レジストリが作成された日時。

- `UpdatedTime` – UTF-8 文字列。

レジストリが更新された日時。

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `EntityNotFoundException`
- `InternalServiceException`

## PutSchemaVersionMetadata アクション (Python: `put_schema_version_metadata`)

指定されたスキーマバージョン ID のメタデータのキーバリューペアを配置します。各スキーマバージョンに対して最大 10 個のキーバリューペアを持つことができます。これらは、1 つ以上の呼び出しで追加できます。

### リクエスト

- `SchemaId` – [SchemaId](#) オブジェクト。

スキーマの一意の ID。

- `SchemaVersionNumber` – [SchemaVersionNumber](#) オブジェクト。

スキーマのバージョン番号。

- `SchemaVersionId` – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。

スキーマバージョンの一意のバージョン ID。

- `MetadataKeyValue` – 必須: [MetadataKeyValuePair](#) オブジェクト。

メタデータキーの対応する値。

### レスポンス

- `SchemaArn` – UTF-8 文字列。1 ~ 10240 バイト長。 [Custom string pattern #22](#) に一致。

スキーマの Amazon リソースネーム (ARN)。

- `SchemaName` – UTF-8 文字列、1~255 バイト長、[Custom string pattern #18](#) に一致。  
スキーマの名前。
- `RegistryName` – UTF-8 文字列、1~255 バイト長、[Custom string pattern #18](#) に一致。  
レジストリの名前。
- `LatestVersion` – ブール。  
スキーマの最新バージョン。
- `VersionNumber` – 数値 (long)。1~100000。  
スキーマのバージョン番号。
- `SchemaVersionId` – UTF-8 文字列。36 バイト長ちょうど。[Custom string pattern #17](#) に一致。  
スキーマバージョンの一意のバージョン ID。
- `MetadataKey` – UTF-8 文字列。1~128 バイト長。[Custom string pattern #33](#) に一致。  
メタデータキー。
- `MetadataValue` – UTF-8 文字列。1~256 バイト長。[Custom string pattern #33](#) に一致。  
メタデータキーの値。

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `AlreadyExistsException`
- `EntityNotFoundException`
- `ResourceNumberLimitExceededException`

## QuerySchemaVersionMetadata アクション (Python: `query_schema_version_metadata`)

スキーマバージョンのメタデータ情報をクエリします。

## リクエスト

- SchemaId – [SchemaId](#) オブジェクト。

スキーマ名と Amazon リソースネーム (ARN) を含む可能性のあるラッパー構造体。

- SchemaVersionNumber – [SchemaVersionNumber](#) オブジェクト。

スキーマのバージョン番号。

- SchemaVersionId – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。

スキーマバージョンの一意のバージョン ID。

- MetadataList – [MetadataKeyValuePair](#) オブジェクトの配列。

キーバリューペアからメタデータを検索します。キーバリューペアが指定されていない場合は、すべてのメタデータ情報がフェッチされます。

- MaxResults – 数値 (整数)。1 ~ 50。

1 ページあたりに必要な結果の最大数。値が指定されていない場合、これはデフォルトで 1 ページあたり 25 に設定されます。

- NextToken – UTF-8 文字列。

継続トークン (これが継続呼び出しの場合)。

## 応答

- MetadataInfoMap – キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1 ~ 128 バイト長で、 [Custom string pattern #33](#) に一致します。

各値は [MetadataInfo](#) オブジェクトです。

メタデータキーおよび関連する値のマップ。

- SchemaVersionId – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。

スキーマバージョンの一意のバージョン ID。

- NextToken – UTF-8 文字列。

返されたトークンのリストをページ分割するための継続トークン。リストの現在のセグメントが最後のセグメントではない場合に返されます。

## エラー

- `InvalidInputException`
- `AccessDeniedException`
- `EntityNotFoundException`

## RemoveSchemaVersionMetadata アクション (Python: `remove_schema_version_metadata`)

指定されたスキーマバージョン ID のスキーマバージョンメタデータからキーバリューペアを削除します。

### リクエスト

- `SchemaId` – [Schemald](#) オブジェクト。  
スキーマ名と Amazon リソースネーム (ARN) を含む可能性のあるラッパー構造体。
- `SchemaVersionNumber` – [SchemaVersionNumber](#) オブジェクト。  
スキーマのバージョン番号。
- `SchemaVersionId` – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。  
スキーマバージョンの一意のバージョン ID。
- `MetadataKeyValue` – 必須: [MetadataKeyValuePair](#) オブジェクト。  
メタデータキーの値。

### レスポンス

- `SchemaArn` – UTF-8 文字列。1 ~ 10240 バイト長。 [Custom string pattern #22](#) に一致。  
スキーマの Amazon リソースネーム (ARN)。
- `SchemaName` – UTF-8 文字列、1 ~ 255 バイト長、 [Custom string pattern #18](#) に一致。  
スキーマの名前。
- `RegistryName` – UTF-8 文字列、1 ~ 255 バイト長、 [Custom string pattern #18](#) に一致。  
レジストリの名前。

- LatestVersion – ブール。  
スキーマの最新バージョン。
- VersionNumber – 数値 (long)。1 ~ 100000。  
スキーマのバージョン番号。
- SchemaVersionId – UTF-8 文字列。36 バイト長ちょうど。 [Custom string pattern #17](#) に一致。  
スキーマバージョンのバージョン ID。
- MetadataKey – UTF-8 文字列。1 ~ 128 バイト長。 [Custom string pattern #33](#) に一致。  
メタデータキー。
- MetadataValue – UTF-8 文字列。1 ~ 256 バイト長。 [Custom string pattern #33](#) に一致。  
メタデータキーの値。

## エラー

- InvalidInputException
- AccessDeniedException
- EntityNotFoundException

## DeleteRegistry アクション (Python: delete\_registry)

スキーマとそのすべてのバージョンを含め、レジストリ全体を削除します。削除オペレーションのステータスを取得するには、非同期呼び出し後に GetRegistry API を呼び出します。レジストリを削除すると、レジストリに対する UpdateRegistry、CreateSchema、UpdateSchema、および RegisterSchemaVersion API など、すべてのオンラインオペレーションが非アクティブ化されません。

## リクエスト

- RegistryId – 必須: [RegistryId](#) オブジェクト。

これは、レジストリ名と Amazon リソースネーム (ARN) を含む可能性があるラッパー構造です。

## レスポンス

- RegistryName – UTF-8 文字列、1～255 バイト長、「[Custom string pattern #18](#)」に一致。

削除されるレジストリの名前。

- RegistryArn – UTF-8 文字列。1～10240 バイト長。[Custom string pattern #22](#) に一致。

削除されるレジストリの Amazon リソースネーム (ARN)。

- Status – UTF-8 文字列 (有効な値: AVAILABLE | DELETING)。

レジストリのステータス。オペレーションが正常終了すると、Deleting ステータスを返しません。

## エラー

- InvalidInputException
- EntityNotFoundException
- AccessDeniedException
- ConcurrentModificationException

## DeleteSchema アクション (Python: delete\_schema)

スキーマセットとそのすべてのバージョンを含め、スキーマセット全体を削除します。削除オペレーションのステータスを取得するには、非同期呼び出し後に GetSchema API を呼び出します。レジストリを削除すると、スキーマに対する GetSchemaByDefinition および RegisterSchemaVersion API など、すべてのオンライン操作が非アクティブ化されます。

## リクエスト

- SchemaId – 必須: [Schemald](#) オブジェクト。

これは、スキーマ名と Amazon リソースネーム (ARN) を含む可能性のあるラッパー構造です。

## レスポンス

- SchemaArn – UTF-8 文字列。1～10240 バイト長。[Custom string pattern #22](#) に一致。

削除されるスキーマの Amazon リソースネーム (ARN)。

- SchemaName – UTF-8 文字列、1~255 バイト長、[Custom string pattern #18](#) に一致。

削除されるスキーマの名前。

- Status – UTF-8 文字列 (有効な値: AVAILABLE | PENDING | DELETING)。

スキーマのステータス。

## エラー

- InvalidInputException
- EntityNotFoundException
- AccessDeniedException
- ConcurrentModificationException

## DeleteSchemaVersions アクション (Python: delete\_schema\_versions)

指定したスキーマからバージョンを削除します。バージョン番号または範囲を指定できます。互換モードが BACKWARDS\_FULL など、必要なバージョンの削除を禁止している場合、エラーが返されます。この呼び出しの後に GetSchemaVersions API を呼び出すと、削除されたバージョンのステータスが一覧表示されます。

バージョン番号の範囲にチェックポイントが付加されたバージョンが含まれている場合、APIは 409 の競合を返し、削除を続行しません。このAPIを使用する前に、まず、DeleteSchemaCheckpoint APIを使用してチェックポイントを削除する必要があります。

DeleteSchemaVersions API を使用して、スキーマセット内の最初のスキーマバージョンを削除することはできません。最初のスキーマバージョンは、DeleteSchema API でのみ削除できます。この操作により、スキーマバージョンの下にアタッチされた SchemaVersionMetadata も削除されます。データベースでは必ず物理削除 (ハード削除) となります。

互換モードが BACKWARDS\_FULL など、必要なバージョンの削除を禁止している場合、エラーが返されます。

## リクエスト

- SchemaId – 必須: [Schemald](#) オブジェクト。

これは、スキーマ名と Amazon リソースネーム (ARN) を含む可能性のあるラッパー構造です。

- Versions – 必須: UTF-8 文字列。1 ~ 100000 バイト長。 [Custom string pattern #34](#) に一致。

バージョン範囲は、次のような形式で指定することができます。

- 単一のバージョン番号、5
- 範囲、5-8: deletes versions 5, 6, 7, 8

## 応答

- SchemaVersionErrors – [SchemaVersionErrorItem](#) オブジェクトの配列。

SchemaVersionErrorItem オブジェクトのリスト。それぞれにはエラーとスキーマのバージョンが含まれています。

## エラー

- InvalidInputException
- EntityNotFoundException
- AccessDeniedException
- ConcurrentModificationException

## ワークフロー

ワークフロー API では、AWS Glue でのワークフローの作成、更新、表示に関連するデータ型と API について説明します。ジョブ実行履歴には、ワークフローとジョブ実行のために 90 日間アクセスできます。

## データ型

- [JobNodeDetails の構造](#)
- [CrawlerNodeDetails 構造](#)
- [TriggerNodeDetails 構造](#)
- [Crawl 構造](#)
- [Node 構造](#)
- [Edge 構造](#)
- [Workflow 構造](#)

- [WorkflowGraph 構造](#)
- [WorkflowRun 構造](#)
- [WorkflowRunStatistics 構造](#)
- [StartingEventBatchCondition 構造](#)
- [Blueprint 構造](#)
- [BlueprintDetails 構造](#)
- [LastActiveDefinition 構造](#)
- [BlueprintRun 構造](#)

## JobNodeDetails の構造

ワークフロー内のジョブノードの詳細。

フィールド

- JobRuns – [JobRun](#) オブジェクトの配列。  
ジョブノードが表すジョブ実行に関する情報。

## CrawlerNodeDetails 構造

ワークフロー内に存在するクローラノードの詳細。

フィールド

- Crawls – [Crawl](#) オブジェクトの配列。  
クローラノードが表すクローラのリスト。

## TriggerNodeDetails 構造

ワークフロー内に存在するトリガーノードの詳細。

フィールド

- Trigger – [Trigger トリガー](#) オブジェクト。

トリガーノードが表すトリガーの情報。

## Crawl 構造

ワークフロー内のクロールの詳細。

フィールド

- State - UTF-8 文字列 (有効値: RUNNING | CANCELLING | CANCELLED | SUCCEEDED | FAILED | ERROR)。

クローラの状態。

- StartedOn - タイムスタンプ。

クロールが開始された日時。

- CompletedOn - タイムスタンプ。

クロールが完了した日時。

- ErrorMessage - 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

クロールに関連付けられたエラーメッセージ。

- LogGroup - UTF-8 文字列。1 ~ 512 バイト長。[Log group string pattern](#) に一致。

クロールに関連付けられたロググループ。

- LogStream - UTF-8 文字列。1 ~ 512 バイト長。[Log-stream string pattern](#) に一致。

クロールに関連付けられたログストリーム。

## Node 構造

ノードは、ワークフローグラフ上で、AWS Glue コンポーネント (トリガー、クローラ、またはジョブ) を表しています。

フィールド

- Type - UTF-8 文字列 (有効な値: CRAWLER | JOB | TRIGGER)。

ノードが表す AWS Glue コンポーネントのタイプ。

- Name – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。  
ノードが表す AWS Glue コンポーネントの名前。
- UniqueId – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。  
ワークフロー内のノードに割り当てられた一意の ID。
- TriggerDetails – [TriggerNodeDetails](#) オブジェクト。  
ノードがトリガーを表すときのトリガーの詳細。
- JobDetails – [JobNodeDetails](#) オブジェクト。  
ノードがジョブを表すときのジョブの詳細。
- CrawlerDetails – [CrawlerNodeDetails](#) オブジェクト。  
ノードがクローラを表すときのクローラの詳細。

## Edge 構造

エッジは、エッジが属するワークフローの一部である 2 つの AWS Glue コンポーネント間の有向接続を表します。

### フィールド

- SourceId – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。  
エッジの始点となるワークフロー内のノードの一意の ID。
- DestinationId – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。  
エッジの終点となるワークフロー内のノードの一意の ID。

## Workflow 構造

ワークフローは、複雑な ETL タスクを完了するために実行される複数の関連する AWS Glue ジョブとクローラのコレクションです。ワークフローは、それに含まれるすべてのジョブとクローラの実行とモニタリングを管理します。

### フィールド

- Name – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

ワークフローの名前。

- `Description` – UTF-8 文字列。

ワークフローの説明。

- `DefaultRunProperties` – キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。

各値は UTF-8 文字列。

ワークフローの各実行の一部として使用されるプロパティのコレクション。実行プロパティは、ワークフローの各ジョブで使用できるようになっています。ジョブは、フロー内の次のジョブのためにプロパティを変更できます。

- `CreatedOn` – タイムスタンプ。

ワークフローが作成された日時。

- `LastModifiedOn` – タイムスタンプ。

ワークフローが最終に変更された日時。

- `LastRun` – [WorkflowRun](#) オブジェクト。

ワークフローの最終の実行に関する情報。

- `Graph` – [WorkflowGraph](#) オブジェクト。

ワークフローに属するすべての AWS Glue コンポーネントをノードとして表し、これらのコンポーネント間の有向接続をエッジとして表すグラフ。

- `CreationStatus` – UTF-8 文字列 (有効な値: `CREATING` | `CREATED` | `CREATION_FAILED`)。

ワークフローの作成ステータス。

- `MaxConcurrentRuns` – 数値 (整数)。

このパラメータを使用すると、データが不必要に複数回更新されないようにしたり、コストを管理したり、場合によってはコンポーネントジョブの同時実行の最大数を超えないようにしたりできます。このパラメータを空白のままにした場合、ワークフロー実行の数に制限はありません。

- `BlueprintDetails` – [BlueprintDetails](#) オブジェクト。

この構造は、この特定のワークフローの作成元となるブループリントの詳細を示します。

## WorkflowGraph 構造

ワークフローグラフは、ワークフロー内に存在するすべての AWS Glue コンポーネントと、これらのコンポーネント間のすべての有向接続を含むワークフロー全体を表します。

フィールド

- Nodes – [ノード](#) オブジェクトの配列。  
ワークフロー内でノードとして表される AWS Glue コンポーネントのリスト。
- Edges – [Edge](#) オブジェクトの配列。  
ワークフローに属するノード間のすべての有向接続のリスト。

## WorkflowRun 構造

ワークフロー実行は、すべてのランタイム情報を提供するワークフローの実行です。

フィールド

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
実行されたワークフローの名前。
- WorkflowRunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
このワークフロー実行の ID。
- PreviousRunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
前回のワークフロー実行の ID。
- WorkflowRunProperties – キーバリューペアのマップ配列。  
各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。  
各値は UTF-8 文字列。  
実行中に設定されたワークフロー実行のプロパティ。
- StartedOn – タイムスタンプ。  
ワークフロー実行の開始日時。
- CompletedOn – タイムスタンプ。

ワークフロー実行の完了日時。

- Status – UTF-8 文字列 (有効な値: RUNNING | COMPLETED | STOPPING | STOPPED | ERROR)。

ワークフロー実行のステータス。

- ErrorMessage – UTF-8 文字列。

このエラーメッセージは、ワークフローの実行の開始時に発生した可能性のあるエラーを示します。現在のところ、唯一のエラーメッセージは「Concurrent runs exceeded for workflow: foo」(ワークフローの同時実行数を超過しました:) のみです。

- Statistics – [WorkflowRunStatistics](#) オブジェクト。

実行の統計情報。

- Graph – [WorkflowGraph](#) オブジェクト。

ワークフローに属するすべての AWS Glue コンポーネントをノードとして表し、これらのコンポーネント間の有向接続をエッジとして表すグラフ。

- StartingEventBatchCondition – [StartingEventBatchCondition](#) オブジェクト。

ワークフロー実行を開始したバッチ条件。

## WorkflowRunStatistics 構造

ワークフロー実行の統計は、ワークフロー実行に関する統計情報を提供します。

フィールド

- TotalActions – 数値 (整数)。

ワークフロー実行のアクションの合計数。

- TimeoutActions – 数値 (整数)。

タイムアウトしたアクションの合計数。

- FailedActions – 数値 (整数)。

失敗したアクションの合計数。

- StoppedActions – 数値 (整数)。

停止したアクションの合計数。

- SucceededActions – 数値 (整数)。

成功したアクションの合計数。

- RunningActions – 数値 (整数)。

実行中状態のアクションの合計数。

- ErroredActions – 数値 (整数)。

ワークフロー実行で ERROR 状態のジョブ実行数を示します。

- WaitingActions – 数値 (整数)。

ワークフロー実行で WAITING 状態のジョブ実行数を示します。

## StartingEventBatchCondition 構造

ワークフロー実行を開始したバッチ条件。到着したイベント数がバッチサイズに到達したか (BatchSize メンバーがゼロ以外の場合)、バッチウィンドウの時間が経過したか (BatchWindow メンバーがゼロ以外の場合) のいずれかです。

フィールド

- BatchSize – 数値 (整数)。

バッチ内のイベント数。

- BatchWindow – 数値 (整数)。

バッチウィンドウの持続時間 (秒単位)。

## Blueprint 構造

ブループリントの詳細。

フィールド

- Name – UTF-8 文字列。1 ~ 128 バイト長。 [Custom string pattern #27](#) に一致。

ブループリントの名前。

- Description – UTF-8 文字列。1 ~ 512 バイト長。

ブループリントの説明。

- CreatedOn – タイムスタンプ。

ブループリントが登録された日時。

- LastModifiedOn – タイムスタンプ。

ブループリントの最終変更日時。

- ParameterSpec – UTF-8 文字列。1 ~ 131072 バイト長。

ブループリントのパラメータ仕様のリストを示す JSON 文字列。

- BlueprintLocation – UTF-8 文字列。

ブループリントを公開する Simple Storage Service (Amazon S3) のパスを指定します。

- BlueprintServiceLocation – UTF-8 文字列。

CreateBlueprint/UpdateBlueprint を呼び出してブループリントを AWS Glue に登録するときに、ブループリントがコピーされる Simple Storage Service (Amazon S3) のパスを指定します。

- Status – UTF-8 文字列 (有効な値: CREATING | ACTIVE | UPDATING | FAILED)。

ブループリント登録のステータス。

- Creating – ブループリントの登録が進行中です。
  - Active – ブループリントは正常に登録されました。
  - Updating – ブループリント登録の更新が進行中です。
  - Failed – ブループリントの登録に失敗しました。
- ErrorMessage – UTF-8 文字列。

エラーメッセージ。

- LastActiveDefinition – [LastActiveDefinition](#) オブジェクト。

ブループリントに複数のバージョンがあり、最新バージョンに何らかのエラーがある場合、この属性は、サービスで使用できる正常な最後のブループリント定義を示します。

## BlueprintDetails 構造

ブループリントの詳細。

## フィールド

- `BlueprintName` – UTF-8 文字列。1～128 バイト長。 [Custom string pattern #27](#) に一致。

ブループリントの名前。

- `RunId` – UTF-8 文字列、1～255 バイト長、 [Single-line string pattern](#) に一致。

このブループリントの実行 ID。

## LastActiveDefinition 構造

ブループリントに複数のバージョンがあり、最新バージョンに何らかのエラーがある場合、この属性は、サービスで使用できる正常な最後のブループリント定義を示します。

### フィールド

- `Description` – UTF-8 文字列。1～512 バイト長。

ブループリントの説明。

- `LastModifiedOn` – タイムスタンプ。

ブループリントの最終変更日時。

- `ParameterSpec` – UTF-8 文字列。1～131072 バイト長。

ブループリントのパラメータを指定する JSON 文字列。

- `BlueprintLocation` – UTF-8 文字列。

AWS Glue デベロッパーによってブループリントが公開される Simple Storage Service (Amazon S3) 内のパスを指定します。

- `BlueprintServiceLocation` – UTF-8 文字列。

ブループリントの作成または更新時にコピーされる Simple Storage Service (Amazon S3) 内のパスを指定します。

## BlueprintRun 構造

ブループリントの実行の詳細。

## フィールド

- `BlueprintName` – UTF-8 文字列。1～128 バイト長。 [Custom string pattern #27](#) に一致。

ブループリントの名前。

- `RunId` – UTF-8 文字列、1～255 バイト長、 [Single-line string pattern](#) に一致。

このブループリント実行の実行 ID。

- `WorkflowName` – UTF-8 文字列、1～255 バイト長、 [Single-line string pattern](#) に一致。

ブループリントが正常に実行された結果として作成されるワークフローの名前。ブループリントの実行にエラーがある場合、ワークフローは作成されません。

- `State` – UTF-8 文字列 (有効な値: `RUNNING` | `SUCCEEDED` | `FAILED` | `ROLLING_BACK`)。

ブループリントの実行の状態。可能な値は以下のとおりです。

- `Running` – ブループリントの実行が進行中です。
- `Succeeded` – ブループリントの実行が正常に完了しました。
- `Failed` – ブループリントの実行が失敗し、ロールバックが完了しました。
- `Rolling Back` – ブループリントの実行が失敗し、ロールバックが進行中です。
- `StartedOn` – タイムスタンプ。

ブループリントの実行が開始された日時。

- `CompletedOn` – タイムスタンプ。

ブループリントの実行が完了した日時。

- `ErrorMessage` – UTF-8 文字列。

ブループリントの実行中に発生したエラーを示します。

- `RollbackErrorMessage` – UTF-8 文字列。

ワークフローのエンティティの作成中にエラーが発生した場合は、作成されたエンティティをその時点までロールバックして削除を試みます。この属性は、作成されたエンティティを削除しようとしたときにエラーが発生したことを示します。

- `Parameters` – UTF-8 文字列。1～131072 バイト長。

文字列のブループリントパラメータです。 `Blueprint$ParameterSpec` で定義されているパラメータ仕様から必要な各キーの値を指定する必要があります。

- RoleArn – UTF-8 文字列。1 ~ 1024 バイト長。 [Custom string pattern #26](#) に一致。

ロール ARN。このロールは、AWS Glue サービスが継承し、ワークフローとワークフローの他のエンティティを作成するために使用されます。

## 操作

- [CreateWorkflow アクション \(Python: create\\_workflow\)](#)
- [UpdateWorkflow アクション \(Python: update\\_workflow\)](#)
- [DeleteWorkflow アクション \(Python: delete\\_workflow\)](#)
- [GetWorkflow アクション \(Python: get\\_workflow\)](#)
- [ListWorkflows アクション \(Python: list\\_workflows\)](#)
- [BatchGetWorkflows アクション \(Python: batch\\_get\\_workflows\)](#)
- [GetWorkflowRun アクション \(Python: get\\_workflow\\_run\)](#)
- [GetWorkflowRuns アクション \(Python: get\\_workflow\\_runs\)](#)
- [GetWorkflowRunProperties アクション \(Python: get\\_workflow\\_run\\_properties\)](#)
- [PutWorkflowRunProperties アクション \(Python: put\\_workflow\\_run\\_properties\)](#)
- [CreateBlueprint アクション \(Python: create\\_blueprint\)](#)
- [UpdateBlueprint アクション \(Python: update\\_blueprint\)](#)
- [DeleteBlueprint アクション \(Python: delete\\_blueprint\)](#)
- [ListBlueprints アクション \(Python: list\\_blueprints\)](#)
- [BatchGetBlueprints アクション \(Python: batch\\_get\\_blueprints\)](#)
- [StartBlueprintRun アクション \(Python: start\\_blueprint\\_run\)](#)
- [GetBlueprintRun アクション \(Python: get\\_blueprint\\_run\)](#)
- [GetBlueprintRuns アクション \(Python: get\\_blueprint\\_runs\)](#)
- [StartWorkflowRun アクション \(Python: start\\_workflow\\_run\)](#)
- [StopWorkflowRun アクション \(Python: stop\\_workflow\\_run\)](#)
- [ResumeWorkflowRun アクション \(Python: resume\\_workflow\\_run\)](#)

## CreateWorkflow アクション (Python: create\_workflow)

新しいワークフローを作成します。

## リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ワークフローに割り当てる名前。アカウント内で一意である必要があります。

- Description – UTF-8 文字列。

ワークフローの説明。

- DefaultRunProperties – キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。

各値は UTF-8 文字列。

ワークフローの各実行の一部として使用されるプロパティのコレクション。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

このワークフローで使用するタグ。

- MaxConcurrentRuns – 数値 (整数)。

このパラメータを使用すると、データが不必要に複数回更新されないようにしたり、コストを管理したり、場合によってはコンポーネントジョブの同時実行の最大数を超えないようにしたりできます。このパラメータを空白のままにした場合、ワークフロー実行の数に制限はありません。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

リクエストの一部として指定されたワークフローの名前。

## エラー

- AlreadyExistsException
- InvalidInputException
- InternalServiceException

- `OperationTimeoutException`
- `ResourceNumberLimitExceededException`
- `ConcurrentModificationException`

## UpdateWorkflow アクション (Python: `update_workflow`)

既存のワークフローを更新します。

### リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
更新するワークフローの名前。
- `Description` – UTF-8 文字列。  
ワークフローの説明。
- `DefaultRunProperties` – キーバリューペアのマップ配列。  
各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。  
各値は UTF-8 文字列。  
ワークフローの各実行の一部として使用されるプロパティのコレクション。
- `MaxConcurrentRuns` – 数値 (整数)。  
このパラメータを使用すると、データが不必要に複数回更新されないようにしたり、コストを管理したり、場合によってはコンポーネントジョブの同時実行の最大数を超えないようにしたりできます。このパラメータを空白のままにした場合、ワークフロー実行の数に制限はありません。

### レスポンス

- `Name` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
入力で指定されたワークフローの名前。

### エラー

- `InvalidInputException`

- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException
- ConcurrentModificationException

## DeleteWorkflow アクション (Python: delete\_workflow)

ワークフローを削除します。

リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
削除するワークフローの名前。

レスポンス

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
入力で指定されたワークフローの名前。

エラー

- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- ConcurrentModificationException

## GetWorkflow アクション (Python: get\_workflow)

ワークフローのリソースメタデータを取得します。

リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
取得するワークフローの名前。

- IncludeGraph – ブール。

ワークフローのリソースメタデータを返すときに、グラフを含めるかどうかを指定します。

## レスポンス

- Workflow – [ワークフロー](#) オブジェクト。

ワークフローのリソースメタデータ。

## エラー

- InvalidInputException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException

## ListWorkflows アクション (Python: list\_workflows)

アカウントで作成されたワークフローの名前を一覧表示します。

### リクエスト

- NextToken – UTF-8 文字列。  
継続トークン (これが継続リクエストの場合)。
- MaxResults – 1 未満または 25 を超えない数値 (整数)。

返されるリストの最大サイズ。

### レスポンス

- Workflows – UTF-8 文字列の配列。1 ~ 25 個の文字列。  
アカウント内のワークフローの名前のリスト。
- NextToken – UTF-8 文字列。

継続トークン (一部のワークフロー名が返されていない場合)。

## エラー

- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`

## BatchGetWorkflows アクション (Python: `batch_get_workflows`)

特定のワークフロー名のリストに対応するリソースメタデータのリストを返します。ListWorkflows オペレーションを呼び出した後で、このオペレーションを呼び出すことで、アクセス許可が付与されているデータにアクセスできます。このオペレーションは、タグを使用するアクセス許可条件を含め、すべての IAM のアクセス許可をサポートします。

### リクエスト

- `Names` – 必須: UTF-8 文字列の配列。1 ~ 25 個の文字列。

ワークフロー名のリスト。これらの名前は ListWorkflows オペレーションから返される場合があります。

- `IncludeGraph` – ブール。

ワークフローのリソースメタデータを返すときに、グラフを含めるかどうかを指定します。

### レスポンス

- `Workflows` – [ワークフロー](#) オブジェクトの配列。1~25 個の構造。

ワークフローのリソースメタデータのリスト。

- `MissingWorkflows` – UTF-8 文字列の配列。1~25 個の文字列。

検出されなかったワークフロー名のリスト。

## エラー

- `InternalServiceException`
- `OperationTimeoutException`
- `InvalidInputException`

## GetWorkflowRun アクション (Python: get\_workflow\_run)

特定のワークフロー実行のメタデータを取得します。ジョブ実行履歴には、ワークフローとジョブ実行のために 90 日間アクセスできます。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

実行中のワークフローの名前。

- RunId – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ワークフロー実行の ID。

- IncludeGraph – ブール。

レスポンスにワークフローグラフを含めるかどうかを指定します。

### レスポンス

- Run – [WorkflowRun](#) オブジェクト。

リクエストされたワークフロー実行のメタデータ。

### エラー

- `InvalidInputException`
- `EntityNotFoundException`
- `InternalServiceException`
- `OperationTimeoutException`

## GetWorkflowRuns アクション (Python: get\_workflow\_runs)

特定のワークフローのすべての実行のメタデータを取得します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

返すべき実行のメタデータが属するワークフローの名前。

- IncludeGraph – ブール。

レスポンスにワークフローグラフを含めるかどうかを指定します。

- NextToken – UTF-8 文字列。

応答の最大サイズ。

- MaxResults – 1~1000 の数値 (整数)。

レスポンスに含めるワークフロー実行の最大数。

## レスポンス

- Runs – [WorkflowRun](#) オブジェクトの配列。1~1000 個の構造体。

ワークフロー実行のメタデータオブジェクトのリスト。

- NextToken – UTF-8 文字列。

継続トークン (リクエストしたワークフロー実行の一部が返されていない場合)。

## エラー

- InvalidInputException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException

## GetWorkflowRunProperties アクション (Python: get\_workflow\_run\_properties)

実行中に設定されたワークフロー実行のプロパティを取得します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

実行されたワークフローの名前。

- RunId – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

返す必要がある実行プロパティが属するワークフロー実行の ID。

## レスポンス

- `RunProperties` – キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。

各値は UTF-8 文字列。

指定した実行の進行中に設定されたワークフロー実行のプロパティ。

## エラー

- `InvalidInputException`
- `EntityNotFoundException`
- `InternalServiceException`
- `OperationTimeoutException`

## PutWorkflowRunProperties アクション (Python: `put_workflow_run_properties`)

特定のワークフロー実行に対して、指定されたワークフロー実行プロパティを挿入します。指定された実行内に該当するプロパティが既に存在する場合は、既存の値が上書きされます。それ以外の場合は、既存のプロパティにプロパティが追加されます。

## リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

実行されたワークフローの名前。

- `RunId` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

実行プロパティを更新する対象のワークフロー実行の ID。

- `RunProperties` – 必須: キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。

各値は UTF-8 文字列。

指定された実行に対して挿入するプロパティ。

## レスポンス

- 応答パラメータはありません。

## エラー

- `AlreadyExistsException`
- `EntityNotFoundException`
- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`
- `ResourceNumberLimitExceededException`
- `ConcurrentModificationException`

## CreateBlueprint アクション (Python: `create_blueprint`)

ブループリントを AWS Glue に登録します。

### リクエスト

- `Name` – 必須: UTF-8 文字列。1~128 バイト長。 [Custom string pattern #27](#) に一致。

ブループリントの名前。

- `Description` – UTF-8 文字列。1~512 バイト長。

ブループリントの説明。

- `BlueprintLocation` – 必須: UTF-8 文字列。1~8192 バイト長。 [Custom string pattern #28](#) に一致。

ブループリントを公開する Amazon S3 のパスを指定します。

- `Tags` – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1～128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

このブループリントに適用されるタグ。

## レスポンス

- Name – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

登録されたブループリントの名前を返します。

## エラー

- AlreadyExistsException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException
- ResourceNumberLimitExceededException

## UpdateBlueprint アクション (Python: update\_blueprint)

登録済みのブループリントを更新します。

### リクエスト

- Name – 必須: UTF-8 文字列。1～128 バイト長。[Custom string pattern #27](#) に一致。

ブループリントの名前。

- Description – UTF-8 文字列。1～512 バイト長。

ブループリントの説明。

- BlueprintLocation – 必須: UTF-8 文字列。1～8192 バイト長。[Custom string pattern #28](#) に一致。

ブループリントを公開する Amazon S3 のパスを指定します。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

更新されたブループリントの名前を返します。

## エラー

- EntityNotFoundException
- ConcurrentModificationException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException
- IllegalBlueprintStateException

## DeleteBlueprint アクション (Python: delete\_blueprint)

既存のブループリントを削除します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除するブループリントの名前。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

削除されたブループリントの名前を返します。

## エラー

- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## ListBlueprints アクション (Python: list\_blueprints)

アカウント内のすべてのブループリント名をリスト表示します。

### リクエスト

- NextToken – UTF-8 文字列。  
継続トークン (これが継続リクエストの場合)。
- MaxResults – 1 未満または 25 を超えない数値 (整数)。  
返されるリストの最大サイズ。
- Tags – キーと値のペアのマップ配列。50 ペア以下。  
各キーは UTF-8 文字列で、1~128 バイト長です。  
各値は UTF-8 文字列で、256 バイト長以下です。  
リストを AWS リソースタグでフィルタリングします。

### レスポンス

- Blueprints – UTF-8 文字列の配列。  
アカウント内のブループリントの名前のリスト。
- NextToken – UTF-8 文字列。  
継続トークン (返されていないワークフロー名がある場合)。

### エラー

- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## BatchGetBlueprints アクション (Python: batch\_get\_blueprints)

ブループリントのリストに関する情報を取得します。

## リクエスト

- **Names** – 必須: UTF-8 文字列の配列。1 ~ 25 個の文字列。

ブループリント名のリスト。

- **IncludeBlueprint** – ブール。

レスポンスにブループリントを含めるかどうかを指定します。

- **IncludeParameterSpec** – ブール。

レスポンスにブループリントのパラメータを JSON 文字列として含めるかどうかを指定します。

## レスポンス

- **Blueprints** – [Blueprint](#) オブジェクトの配列。

ブループリントのリストを **Blueprints** オブジェクトとして返します。

- **MissingBlueprints** – UTF-8 文字列の配列。

検出されなかった **BlueprintNames** のリストを返します。

## エラー

- **InternalServiceException**
- **OperationTimeoutException**
- **InvalidInputException**

## StartBlueprintRun アクション (Python: start\_blueprint\_run)

指定したブループリントの新しい実行を開始します。

### リクエスト

- **BlueprintName** – 必須: UTF-8 文字列。1 ~ 128 バイト長。 [Custom string pattern #27](#) に一致。

ブループリントの名前。

- **Parameters** – UTF-8 文字列。1 ~ 131072 バイト長。

パラメータを `BlueprintParameters` オブジェクトとして指定します。

- `RoleArn` – 必須: UTF-8 文字列。1 ~ 1024 バイト長。 [Custom string pattern #26](#) に一致。

ワークフローの作成に使用される IAM ロールを指定します。

## レスポンス

- `RunId` – UTF-8 文字列、1 ~ 255 バイト長、 [Single-line string pattern](#) に一致。

このブループリント実行の実行 ID。

## エラー

- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`
- `ResourceNumberLimitExceededException`
- `EntityNotFoundException`
- `IllegalBlueprintStateException`

## GetBlueprintRun アクション (Python: `get_blueprint_run`)

ブループリント実行の詳細を取得します。

### リクエスト

- `BlueprintName` – 必須: UTF-8 文字列。1 ~ 128 バイト長。 [Custom string pattern #27](#) に一致。

ブループリントの名前。

- `RunId` – 必須: UTF-8 文字列、1 ~ 255 バイト長、 [Single-line string pattern](#) に一致。

取得するブループリント実行の実行 ID。

## レスポンス

- `BlueprintRun` – [BlueprintRun](#) オブジェクト。

BlueprintRun オブジェクトを返します。

## エラー

- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException

## GetBlueprintRuns アクション (Python: get\_blueprint\_runs)

指定したブループリントでのブループリント実行の詳細を取得します。

### リクエスト

- BlueprintName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
ブループリントの名前。
- NextToken – UTF-8 文字列。  
継続トークン (これが継続リクエストの場合)。
- MaxResults – 1~1000 の数値 (整数)。  
返されるリストの最大サイズ。

### レスポンス

- BlueprintRuns – [BlueprintRun](#) オブジェクトの配列。  
BlueprintRun オブジェクトのリストを返します。
- NextToken – UTF-8 文字列。  
継続トークン (一部のブループリント実行が返されていない場合)。

## エラー

- EntityNotFoundException
- InternalServiceException

- `OperationTimeoutException`
- `InvalidInputException`

## StartWorkflowRun アクション (Python: `start_workflow_run`)

指定したワークフローの新しい実行を開始します。

### リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

開始するワークフローの名前。

- `RunProperties` – キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。

各値は UTF-8 文字列。

新しいワークフロー実行のためのワークフロー実行プロパティ。

### レスポンス

- `RunId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

新しい実行の ID。

### エラー

- `InvalidInputException`
- `EntityNotFoundException`
- `InternalServiceException`
- `OperationTimeoutException`
- `ResourceNumberLimitExceededException`
- `ConcurrentRunsExceededException`

## StopWorkflowRun アクション (Python: stop\_workflow\_run)

指定したワークフロー実行の実行を停止します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
停止するワークフローの名前。
- RunId – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
停止するワークフロー実行の ID。

### レスポンス

- 応答パラメータはありません。

### エラー

- InvalidInputException
- EntityNotFoundException
- InternalServiceException
- OperationTimeoutException
- IllegalWorkflowStateException

## ResumeWorkflowRun アクション (Python: resume\_workflow\_run)

以前に部分的に完了したワークフロー実行から一部のノードを再起動し、ワークフロー実行を再開します。選択したノードと、選択したノードの下流にあるすべてのノードが実行されます。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
再開するワークフローの名前。
- RunId – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
再開するワークフロー実行の ID。

- NodeIds – 必須: UTF-8 文字列の配列。

再起動するノードのノード ID のリスト。再起動するノードは、元の実行で実行を試行したことがなければなりません。

## レスポンス

- RunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

再開されたワークフロー実行に割り当てられた新しい ID。ワークフロー実行を再開するたびに、新しい実行 ID が割り当てられます。

- NodeIds – UTF-8 文字列の配列。

実際に再起動されたノードのノード ID のリスト。

## エラー

- `InvalidInputException`
- `EntityNotFoundException`
- `InternalServiceException`
- `OperationTimeoutException`
- `ConcurrentRunsExceededException`
- `IllegalWorkflowStateException`

## 使用プロファイル

使用状況プロファイル API は、での使用プロファイルの作成、更新、または表示に関連するデータ型と API について説明します AWS Glue。

## データ型

- [ProfileConfiguration](#) 構造
- [ConfigurationObject](#) 構造
- [UsageProfileDefinition](#) 構造

## ProfileConfiguration 構造

管理者が AWS Glue 使用プロファイルで設定するジョブとセッションの値を指定します。

フィールド

- SessionConfiguration – キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。

各値は [ConfigurationObject](#) オブジェクトです。

AWS Glue セッションの設定パラメータのキーと値のマップ。

- JobConfiguration – キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。

各値は [ConfigurationObject](#) オブジェクトです。

AWS Glue ジョブの設定パラメータのキーと値のマップ。

## ConfigurationObject 構造

AWS Glue 管理者が使用プロファイルで設定された各ジョブまたはセッションパラメータに設定する値を指定します。

フィールド

- DefaultValue – UTF-8 文字列、1~128 バイト長、「[Custom string pattern #31](#)」に一致。

パラメータのデフォルト値。

- AllowedValues – UTF-8 文字列の配列。

パラメータに許可される値のリスト。

- MinValue – UTF-8 文字列、1~128 バイト長、「[Custom string pattern #31](#)」に一致。

パラメータに許可される最小値。

- MaxValue – UTF-8 文字列、1~128 バイト長、「[Custom string pattern #31](#)」に一致。

パラメータに許容される最大値。

## UsageProfileDefinition 構造

AWS Glue 使用状況プロファイルについて説明します。

### フィールド

- Name – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。

使用状況プロファイルの名前。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

使用プロファイルの説明。

- CreatedOn – タイムスタンプ。

使用状況プロファイルが作成された日時。

- LastModifiedOn – タイムスタンプ。

使用状況プロファイルが最後に変更された日時。

## 操作

- [CreateUsageProfile アクション \(Python: create\\_usage\\_profile\)](#)
- [GetUsageProfile アクション \(Python: get\\_usage\\_profile\)](#)
- [UpdateUsageProfile アクション \(Python: update\\_usage\\_profile\)](#)
- [DeleteUsageProfile アクション \(Python: delete\\_usage\\_profile\)](#)
- [ListUsageProfiles アクション \(Python: list\\_usage\\_profiles\)](#)

## CreateUsageProfile アクション (Python: create\_usage\_profile)

AWS Glue 使用プロファイルを作成します。

### リクエスト

- Name – 必須: UTF-8 文字列、1～255 バイト長、「[Single-line string pattern](#)」に一致。

使用状況プロファイルの名前。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

使用プロファイルの説明。

- Configuration – 必須: [ProfileConfiguration](#) オブジェクト。

プロファイルのジョブ値とセッション値を指定するProfileConfigurationオブジェクト。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

使用プロファイルに適用されるタグのリスト。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

作成された使用プロファイルの名前。

## エラー

- InvalidInputException
- InternalServiceException
- AlreadyExistsException
- OperationTimeoutException
- ResourceNumberLimitExceededException
- OperationNotSupportedException

## GetUsageProfile アクション (Python: get\_usage\_profile)

指定された AWS Glue 使用プロファイルに関する情報を取得します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

取得する使用プロファイルの名前。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
使用状況プロファイルの名前。
- Description – 説明文字列、2048 バイト長以下、「[URI address multi-line string pattern](#)」に一致。  
使用プロファイルの説明。
- Configuration – [ProfileConfiguration](#) オブジェクト。  
プロファイルのジョブ値とセッション値を指定するProfileConfigurationオブジェクト。
- CreatedOn – タイムスタンプ。  
使用状況プロファイルが作成された日時。
- LastModifiedOn – タイムスタンプ。  
使用状況プロファイルが最後に変更された日時。

## エラー

- InvalidInputException
- InternalServiceException
- EntityNotFoundException
- OperationTimeoutException
- OperationNotSupportedException

## UpdateUsageProfile アクション (Python: update\_usage\_profile)

AWS Glue 使用状況プロファイルを更新します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
使用状況プロファイルの名前。
- Description – 説明文字列、2048 バイト長以下、「[URI address multi-line string pattern](#)」に一致。  
使用プロファイルの説明。

- Configuration – 必須: [ProfileConfiguration](#) オブジェクト。

プロファイルのジョブ値とセッション値を指定するProfileConfigurationオブジェクト。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

更新した使用状況プロファイルの名前。

## エラー

- InvalidInputException
- InternalServiceException
- EntityNotFoundException
- OperationTimeoutException
- OperationNotSupportedException
- ConcurrentModificationException

## DeleteUsageProfile アクション (Python: delete\_usage\_profile)

AWS Glue 指定された使用プロファイルを削除します。

## リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

削除する使用プロファイルの名前。

## レスポンス

- 応答パラメータはありません。

## エラー

- InvalidInputException
- InternalServiceException

- `OperationTimeoutException`
- `OperationNotSupportedException`

## ListUsageProfiles アクション (Python: `list_usage_profiles`)

すべての AWS Glue 使用プロファイルを一覧表示します。

### リクエスト

- `NextToken` - UTF-8 文字列。400000 バイト長以下。  
継続トークン (これが継続呼び出しの場合)。
- `MaxResults` - 1 未満または 200 を超えない数値 (整数)。  
1 回のレスポンスで返される使用プロファイルの最大数。

### 応答

- `Profiles` - [UsageProfileDefinition](#) オブジェクトの配列。  
使用状況プロファイル (`UsageProfileDefinition`) オブジェクトのリスト。
- `NextToken` - UTF-8 文字列。400000 バイト長以下。  
継続トークン (現在のリストセグメントが最後ではない場合に存在)。

### エラー

- `InternalServiceException`
- `OperationTimeoutException`
- `InvalidInputException`
- `OperationNotSupportedException`

## 機械学習 API

機械学習 API では、機械学習のデータ型について説明します。これには、変換を作成、削除、更新したり、機械学習のタスク実行を開始したりするための API が含まれます。

## データ型

- [TransformParameters 構造](#)
- [EvaluationMetrics 構造](#)
- [MLTransform 構造](#)
- [FindMatchesParameters 構造](#)
- [FindMatchesMetrics 構造](#)
- [ConfusionMatrix 構造](#)
- [GlueTable 構造](#)
- [TaskRun 構造](#)
- [TransformFilterCriteria 構造](#)
- [TransformSortCriteria 構造](#)
- [TaskRunFilterCriteria 構造](#)
- [TaskRunSortCriteria 構造](#)
- [TaskRunProperties 構造](#)
- [FindMatchesTaskRunProperties 構造](#)
- [ImportLabelsTaskRunProperties 構造](#)
- [ExportLabelsTaskRunProperties 構造](#)
- [LabelingSetGenerationTaskRunProperties 構造](#)
- [SchemaColumn 構造](#)
- [TransformEncryption 構造](#)
- [MLUserDataEncryption 構造](#)
- [ColumnImportance 構造](#)

### TransformParameters 構造

機械学習変換に関連付けられたアルゴリズム固有のパラメータ。

#### フィールド

- `TransformType` – 必須: UTF-8 文字列 (有効な値: `FIND_MATCHES`)。

機械学習変換のタイプ。

機械学習変換のタイプの詳細については、「[機械学習変換の作成](#)」を参照してください。

- FindMatchesParameters – [FindMatchesParameters](#) オブジェクト。  
一致検索アルゴリズムのパラメータ。

## EvaluationMetrics 構造

評価メトリクスは、機械学習変換の品質を推定します。

フィールド

- TransformType – 必須: UTF-8 文字列 (有効な値: FIND\_MATCHES)。  
機械学習変換のタイプ。
- FindMatchesMetrics – [FindMatchesMetrics](#) オブジェクト。  
「一致の検索」アルゴリズムの評価メトリクス。

## MLTransform 構造

機械学習変換の構造。

フィールド

- TransformId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
機械学習変換用に生成される一意の変換 ID。ID は一意であることが保証されており、変更されません。
- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
機械学習変換のユーザー定義名。名前は一意であることが保証されておらず、いつでも変更される可能性があります。
- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
機械学習変換に対する長い形式のユーザー定義による説明。説明は一意であることが保証されておらず、いつでも変更される可能性があります。
- Status – UTF-8 文字列 (有効な値: NOT\_READY | READY | DELETING)。

機械学習変換の現在のステータス。

- CreatedOn – タイムスタンプ。

タイムスタンプ。この機械学習変換の作成日時。

- LastModifiedOn – タイムスタンプ。

タイムスタンプ。この機械学習変換の最終変更日時。

- InputRecordTables – [GlueTable](#) オブジェクトの配列。構造体 10 個以下。

変換で使用される AWS Glue テーブル定義のリスト。

- Parameters – [TransformParameters](#) オブジェクト。

TransformParameters オブジェクト。パラメータを使用して機械学習変換の動作を調整 (カスタマイズ) できます。そのために、変換に学習させる要素と、さまざまなトレードオフ (適合率と再現率、精度とコストなど) の選好を指定します。

- EvaluationMetrics – [EvaluationMetrics](#) オブジェクト。

EvaluationMetrics オブジェクト。評価メトリクスは、機械学習変換の品質を推定します。

- LabelCount – 数値 (整数)。

この変換のために AWS Glue で生成されるラベリングファイルのカウント識別子。より適切な変換を作成する場合、ラベリングファイルのダウンロード、ラベル付け、アップロードを繰り返すことができます。

- Schema – [SchemaColumn](#) オブジェクトの配列。構造体 100 個以下。

この変換の実行対象となる列とデータ型を表すキーバリューペアのマッピング。100 列が上限です。

- Role – UTF-8 文字列。

必要なアクセス許可を持つ IAM ロールの Amazon リソースネーム (ARN) の名前。必要なアクセス許可には、AWS Glue リソースに対する AWS Glue サービスロールのアクセス許可と、変換に必要な Amazon S3 アクセス許可の両方が含まれます。

- このロールには、AWS Glue のリソースへのアクセスを許可する AWS Glue サービスロールのアクセス許可が必要です。「[Attach a Policy to IAM Users That Access AWS Glue](#)」を参照してください。

- このロールは、この変換のタスク実行で使用される Amazon Simple Storage Service (Amazon S3) のソース、ターゲット、一時ディレクトリ、スクリプト、およびライブラリへのアクセス許可を必要とします。
- `GlueVersion` – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

この値により、この機械学習変換と互換性がある AWS Glue のバージョンが決定します。ほとんどのお客様に、Glue 1.0 が推奨されます。値が設定されていない場合、Glue の互換性はデフォルトで Glue 0.9 に設定されます。詳細については、デベロッパーガイドの「[AWS Glue Versions](#)」を参照してください。

- `MaxCapacity` – 数値 (double)。

この変換のタスク実行に割り当てられる AWS Glue データ処理ユニット (DPU) の数。DPU は、2~100 の範囲で割り当てることができます。デフォルトは 10 です。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

`MaxCapacity` は、`NumberOfWorkers` と `WorkerType` との相互排他的なオプションです。

- `NumberOfWorkers` または `WorkerType` のいずれかが設定されている場合、`MaxCapacity` は設定できません。
- `MaxCapacity` が設定されている場合、`NumberOfWorkers` または `WorkerType` は設定できません。
- `WorkerType` が設定されている場合、`NumberOfWorkers` は必須です ( 逆も同様です ) 。
- `MaxCapacity` と `NumberOfWorkers` は両方とも少なくとも 1 である必要があります。

`WorkerType` フィールドを `Standard` 以外の値に設定すると、`MaxCapacity` フィールドが自動的に設定され、読み取り専用になります。

- `WorkerType` – UTF-8 文字列 (有効な値: `Standard=""` | `G.1X=""` | `G.2X=""` | `G.025X=""` | `G.4X=""` | `G.8X=""` | `Z.2X=""`) 。

この変換のタスク実行時に割り当てられる事前定義済みワーカーのタイプ。使用できる値は、`Standard`、`G.1X`、または `G.2X` です。

- `Standard` ワーカータイプでは、各ワーカーは 4 vCPU、16 GB のメモリ、50 GB のディスク、ワーカーあたり 2 個のエグゼキューターを提供します。
- `G.1X` ワーカータイプでは、各ワーカーは 4 vCPU、16 GB のメモリ、64 GB のディスク、ワーカーあたり 1 個のエグゼキューターを提供します。

- G.2X ワーカータイプでは、各ワーカーは 8 vCPU、32 GB のメモリ、128 GB のディスク、ワーカーあたり 1 個のエグゼキュターを提供します。

MaxCapacity は、NumberOfWorkers と WorkerType との相互排他的なオプションです。

- NumberOfWorkers または WorkerType のいずれかが設定されている場合、MaxCapacity は設定できません。
- MaxCapacity が設定されている場合、NumberOfWorkers または WorkerType は設定できません。
- WorkerType が設定されている場合、NumberOfWorkers は必須です ( 逆も同様です ) 。
- MaxCapacity と NumberOfWorkers は両方とも少なくとも 1 である必要があります。
- NumberOfWorkers – 数値 (整数)。

変換のタスク実行時に割り当てられる定義済み workerType のワーカー数。

WorkerType が設定されている場合、NumberOfWorkers は必須です ( 逆も同様です ) 。

- Timeout - 数値 (整数)。1 以上。

機械学習変換のタイムアウト (分単位)。

- MaxRetries – 数値 (整数)。

機械学習変換の MLTaskRun が失敗した後の最大再試行回数。

- TransformEncryption – [TransformEncryption](#) オブジェクト。

ユーザーデータへのアクセスに適用される変換の、保管時の暗号化設定。機械学習の変換では、KMS を使用して Amazon S3 で暗号化されたユーザーデータにアクセスできます。

## FindMatchesParameters 構造

「一致検索」変換を設定するためのパラメータ。

### フィールド

- PrimaryKeyColumnName – UTF-8 文字列。1~1024 バイト長。 [Single-line string pattern](#) に一致。

ソーステーブルの行を一意に識別する列の名前。一致するレコードを識別するために使用します。

- PrecisionRecallTradeoff – 数値 (double)。1.0 以下。

変換を調整して適合率と再現率のバランスを取るときに選択する値。値 0.5 は指定なし、値 1.0 は純粋な適合率寄りとなり、値 0.0 は再現率寄りであることを意味します。これはトレードオフの関係にあるため、1.0 に近い値を選択すると、再現率が非常に低くなります。0.0 に近い値を選択すると、適合率が非常に低くなります。

適合率メトリクスは、モデルが一致を推定して、その推定が正確である度合いを示します。

再現率メトリクスは、実際的一致件数に対して、モデルが一致を推定した件数を示します。

- AccuracyCostTradeoff – 数値 (double)。1.0 以下。

変換を調整して精度とコストのバランスを取るときに選択する値。値 0.5 は、システムが精度とコスト懸念のバランスを取っていることを意味します。値 1.0 は完全な精度を意味し、通常、コストが高くなります。著しく高くなる場合もあります。値 0.0 は、完全なコスト優先を意味し、FindMatches 変換の精度が低くなります。精度が許容できないレベルになる場合もあります。

精度は、変換が真陽性と真陰性を正確に検索する度合いを測定します。精度を高めるには、マシンリソースとコストを増やす必要があります。ただし、再現率も高くなります。

コストは、変換を実行するために消費されるコンピューティングリソースの数およびそれに伴う金額を測定します。

- EnforceProvidedLabels – ブール。

ユーザーから指定されたラベルに出力を一致させるためにオン/オフを切り替える値。値が True の場合、find matches 変換は指定されたラベルに出力を一致させます。この結果は通常の合成結果より優先されます。値が False の場合、find matches 変換は、すべての指定されたラベルを優先するとは限らず、結果はトレーニングされたモデルに依存します。

この値を true に設定すると、合成の実行時間が長くなる場合があります。

## FindMatchesMetrics 構造

「一致の検索」アルゴリズムの評価メトリクス。機械学習変換の品質を測定するには、変換によっていくつかの一致を推定し、その結果を同じデータセットの既知の一致と比較します。品質メトリクスは、データのサブセットに基づくものであるため、厳密ではありません。

## フィールド

- `AreaUnderPRCurve` – 数値 (double)。1.0 以下。

適合率/再現率曲線の下の部分の面積 (AUPRC) は、変換の品質全体を測定するための単一の数値です。適合率と再現率のために行われた選択とは関係ありません。値が高いほど、適合率と再現率のトレードオフが適切であることを示します。

詳細については、Wikipedia の「[適合率と再現率](#)」を参照してください。

- `Precision` – 数値 (double)。1.0 以下。

適合率メトリクスは、変換が一致を推定したときに、その推定が正しい度合いを示します。具体的には、真陽性の総数に対して、変換が実際に検索した真陽性の割合を測定します。

詳細については、Wikipedia の「[適合率と再現率](#)」を参照してください。

- `Recall` – 数値 (double)。1.0 以下。

再現率メトリクスは、実際的一致件数に対して、変換が推定した一致件数の割合です。具体的には、ソースデータのレコード総数に対して、変換が検索した真陽性の件数の割合です。

詳細については、Wikipedia の「[適合率と再現率](#)」を参照してください。

- `F1` – 数値 (double)。1.0 以下。

最大の F1 メトリクスは、変換の精度 (1~0) を示します。値 1 は最適な精度です。

詳細については、Wikipedia の「[F1 スコア](#)」を参照してください。

- `ConfusionMatrix` – [ConfusionMatrix](#) オブジェクト。

混同行列は、変換が何を正確に推定しているか、どのような種類のエラーを犯しているかを示します。

詳細については、Wikipedia の「[混同行列](#)」を参照してください。

- `ColumnImportances` – [ColumnImportance](#) オブジェクトの配列。構造体 100 個以下。

重要度の降順でソートされている列重要度メトリクスを含む `ColumnImportance` 構造のリスト。

## ConfusionMatrix 構造

混同行列は、変換が何を正確に推定しているか、どのような種類のエラーを犯しているかを示します。

詳細については、Wikipedia の「[混同行列](#)」を参照してください。

### フィールド

- NumTruePositives – 数値 (long 型)。

変換の混同行列において、変換が正確に検索したデータ内の一致の数。

- NumFalsePositives – 数値 (long 型)。

変換の混同行列において、変換が間違っ一致として分類したデータ内の不一致の数。

- NumTrueNegatives – 数値 (long 型)。

変換の混同行列において、変換が正しく拒否したデータ内の不一致の数。

- NumFalseNegatives – 数値 (long 型)。

変換の混同行列において、変換が検索しなかったデータ内の一致の数。

## GlueTable 構造

入力データまたは出力データとして使用される AWS Glue Data Catalog 内のデータベースとテーブル。

### フィールド

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

AWS Glue Data Catalog のデータベース名。

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

AWS Glue Data Catalog のテーブル名。

- CatalogId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

AWS Glue Data Catalog 用の一意の識別子。

- ConnectionName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

AWS Glue Data Catalog への接続の名前。

- `AdditionalOptions` – キーと値のペアのマップ配列。1~10 個のペア。

各キーは UTF-8 文字列、1~255 バイト長で、[Single-line string pattern](#) に一致します。

各値は説明文字列であり、2048 バイト長以下で、[URI address multi-line string pattern](#) に一致しません。

テーブルの追加オプション。現在、次の 2 つのキーがサポートされています。

- `pushDownPredicate`: データセットのすべてのファイルを一覧表示して読み込むことなく、パーティションをフィルタリングするためのキー。
- `catalogPartitionPredicate`: AWS Glue Data Catalog のパーティションインデックスを使用して、サーバー側のパーティションプルーニングを使用するためのキー。

## TaskRun 構造

機械学習変換に関連付けられているサンプリングパラメータ。

フィールド

- `TransformId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

変換に対する一意の識別子。

- `TaskRunId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このタスク実行に対する一意の識別子。

- `Status` – UTF-8 文字列 (有効な値: STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT)。

リクエストされたタスク実行の現在のステータス。

- `LogGroupName` – UTF-8 文字列。

このタスク実行に関連付けられている、安全なログ記録用のロググループの名前。

- `Properties` – [TaskRunProperties](#) オブジェクト。

このタスク実行に関連付けられている設定プロパティを指定します。

- `ErrorString` – UTF-8 文字列。

このタスク実行に関連付けられているエラー文字列のリスト。

- StartedOn – タイムスタンプ。

このタスク実行の開始日時。

- LastModifiedOn – タイムスタンプ。

リクエストされたタスク実行の最終更新日時。

- CompletedOn – タイムスタンプ。

リクエストされたタスク実行の最終完了日時。

- ExecutionTime – 数値 (整数)。

タスク実行でリソースを消費した時間 (秒)。

## TransformFilterCriteria 構造

機械学習変換のフィルタリング基準。

フィールド

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

機械学習変換のフィルタリングに使用する一意の変換名。

- TransformType – UTF-8 文字列 (有効な値: FIND\_MATCHES)。

機械学習変換のフィルタリングに使用する機械学習変換のタイプ。

- Status – UTF-8 文字列 (有効な値: NOT\_READY | READY | DELETING)。

変換の直近の既知のステータスを使用して機械学習変換のリストをフィルタリングします (変換が使用可能かどうかを示します)。有効な値は「NOT\_READY」、「READY」、または「DELETING」です。

- GlueVersion – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

この値により、この機械学習変換と互換性がある AWS Glue のバージョンが決定します。ほとんどのお客様に、Glue 1.0 が推奨されます。値が設定されていない場合、Glue の互換性はデフォルトで Glue 0.9 に設定されます。詳細については、デベロッパーガイドの「[AWS Glue Versions](#)」を参照してください。

- `CreatedBefore` – タイムスタンプ。

この日時より前に作成された変換をフィルタリングします。

- `CreatedAfter` – タイムスタンプ。

この日時より後に作成された変換をフィルタリングします。

- `LastModifiedBefore` – タイムスタンプ。

この日付より前に最終変更された変換をフィルタリングします。

- `LastModifiedAfter` – タイムスタンプ。

この日付より後に最終更新された変換をフィルタリングします。

- `Schema` – [SchemaColumn](#) オブジェクトの配列。構造体 100 個以下。

特定のスキーマを持つデータセットをフィルタリングします。Map<Column, Type> オブジェクトは、この変換で利用できるスキーマを表すキーと値のペアの配列です。Column は列名、Type はデータ型 (整数や文字列など) です。100 列が上限です。

## TransformSortCriteria 構造

機械学習変換に関連付けられている並べ替え基準。

フィールド

- `Column` – 必須: UTF-8 文字列 (有効な値: NAME | TRANSFORM\_TYPE | STATUS | CREATED | LAST\_MODIFIED)。

機械学習変換に関連付けられている並べ替え基準で使用する列。

- `SortDirection` – 必須: UTF-8 文字列 (有効な値: DESCENDING | ASCENDING)。

機械学習変換に関連付けられている並べ替え基準で使用する並べ替え方向。

## TaskRunFilterCriteria 構造

機械学習変換のタスク実行のフィルタリング基準。

## フィールド

- `TaskRunType` – UTF-8 文字列 (有効な値: `EVALUATION` | `LABELING_SET_GENERATION` | `IMPORT_LABELS` | `EXPORT_LABELS` | `FIND_MATCHES`)。

タスク実行のタイプ。

- `Status` – UTF-8 文字列 (有効な値: `STARTING` | `RUNNING` | `STOPPING` | `STOPPED` | `SUCCEEDED` | `FAILED` | `TIMEOUT`)。

タスク実行の現在のステータス。

- `StartedBefore` – タイムスタンプ。

この日付より前に開始されたタスク実行をフィルタリングします。

- `StartedAfter` – タイムスタンプ。

この日付より後に開始されたタスク実行をフィルタリングします。

## TaskRunSortCriteria 構造

機械学習変換のタスク実行リストを並べ替えるために使用する並べ替え基準。

### フィールド

- `Column` – 必須: UTF-8 文字列 (有効な値: `TASK_RUN_TYPE` | `STATUS` | `STARTED`)。

機械学習変換のタスク実行リストを並べ替えるために使用する列。

- `SortDirection` – 必須: UTF-8 文字列 (有効な値: `DESCENDING` | `ASCENDING`)。

機械学習変換のタスク実行リストを並べ替えるために使用する並べ替え方向。

## TaskRunProperties 構造

タスク実行の設定プロパティ。

### フィールド

- `TaskType` – UTF-8 文字列 (有効な値: `EVALUATION` | `LABELING_SET_GENERATION` | `IMPORT_LABELS` | `EXPORT_LABELS` | `FIND_MATCHES`)。

タスク実行のタイプ。

- `ImportLabelsTaskRunProperties` – [ImportLabelsTaskRunProperties](#) オブジェクト。  
「ラベルのインポート」タスク実行の設定プロパティ。
- `ExportLabelsTaskRunProperties` – [ExportLabelsTaskRunProperties](#) オブジェクト。  
「ラベルのエクスポート」タスク実行の設定プロパティ。
- `LabelingSetGenerationTaskRunProperties` – [LabelingSetGenerationTaskRunProperties](#) オブジェクト。  
「ラベリングセットの生成」タスク実行の設定プロパティ。
- `FindMatchesTaskRunProperties` – [FindMatchesTaskRunProperties](#) オブジェクト。  
「一致検索」タスク実行の設定プロパティ。

## FindMatchesTaskRunProperties 構造

「一致検索」タスク実行の設定プロパティを指定します。

フィールド

- `JobId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
「一致検索」タスク実行のジョブ ID。
- `JobName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
「一致検索」タスク実行のジョブに割り当てられた名前。
- `JobRunId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
「一致検索」タスク実行のジョブ実行 ID。

## ImportLabelsTaskRunProperties 構造

「ラベルのインポート」タスク実行の設定プロパティを指定します。

フィールド

- `InputS3Path` – UTF-8 文字列。

ラベルをインポートする元の Amazon Simple Storage Service (Amazon S3) のパス。

- Replace – ブール。

既存のラベルを上書きするかどうかを示します。

## ExportLabelsTaskRunProperties 構造

「ラベルのエクスポート」タスク実行の設定プロパティを指定します。

フィールド

- OutputS3Path – UTF-8 文字列。

ラベルをエクスポートする先の Amazon Simple Storage Service (Amazon S3) のパス。

## LabelingSetGenerationTaskRunProperties 構造

「ラベリングセットの生成」タスク実行の設定プロパティを指定します。

フィールド

- OutputS3Path – UTF-8 文字列。

ラベリングセットを生成する Amazon Simple Storage Service (Amazon S3) のパス。

## SchemaColumn 構造

この変換の実行対象となる列とデータ型を表すキーバリューペア。MLTransform の Schema パラメータには、これらの構造を最大 100 個まで含めることができます。

フィールド

- Name – UTF-8 文字列。1 ~ 1024 バイト長。 [Single-line string pattern](#) に一致。

列の名前。

- DataType – UTF-8 文字列。131,072 バイト長以下。 [Single-line string pattern](#) に一致。

列のデータ型。

## TransformEncryption 構造

ユーザーデータへのアクセスに適用される変換の、保管時の暗号化設定。機械学習の変換では、KMS を使用して Amazon S3 で暗号化されたユーザーデータにアクセスできます。

さらに、インポートされたラベルとトレーニングされた変換は、顧客が提供した KMS キーを使用して暗号化できるようになりました。

### フィールド

- `MLUserDataEncryption` – [MLUserDataEncryption](#) オブジェクト。

暗号化モードとお客様が用意した KMS キー ID `MLUserDataEncryption` を含むオブジェクト。

- `TaskRunSecurityConfigurationName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

セキュリティ構成の名前。

## MLUserDataEncryption 構造

ユーザーデータへのアクセスに適用される変換の、保管時の暗号化設定。

### フィールド

- `MLUserDataEncryptionMode` – 必須: UTF-8 文字列 (有効な値: DISABLED | SSE-KMS="SSEKMS")。

ユーザーデータに適用される暗号化モード。有効な値は次のとおりです。

- `DISABLED`: 暗号化は無効です
- `SSEKMS`: Simple Storage Service (Amazon S3) に保存されたユーザーデータに、AWS Key Management Service (SSE-KMS) によるサーバー側暗号化を使用します。
- `KmsKeyId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ユーザーが用意した KMS キーの ID。

## ColumnImportance 構造

列の列名と列重要度スコアを含む構造。

列の重要度は、レコード内のどの列が他の列よりも重要かを識別することで、モデルへの列貢献度を理解するのに役立ちます。

## フィールド

- `ColumnName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

列の名前。

- `Importance` – 数値 (double)。1.0 以下。

列の列重要度スコア (10 進数)。

## 操作

- [CreateMLTransform アクション \(Python: create\\_ml\\_transform\)](#)
- [UpdateMLTransform アクション \(Python: update\\_ml\\_transform\)](#)
- [DeleteMLTransform アクション \(Python: delete\\_ml\\_transform\)](#)
- [GetMLTransform アクション \(Python: get\\_ml\\_transform\)](#)
- [GetMLTransforms アクション \(Python: get\\_ml\\_transforms\)](#)
- [ListMLTransforms アクション \(Python: list\\_ml\\_transforms\)](#)
- [StartMLEvaluationTaskRun アクション \(Python: start\\_ml\\_evaluation\\_task\\_run\)](#)
- [StartMLLabelingSetGenerationTaskRun アクション \(Python: start\\_ml\\_labeling\\_set\\_generation\\_task\\_run\)](#)
- [GetMLTaskRun アクション \(Python: get\\_ml\\_task\\_run\)](#)
- [GetMLTaskRuns アクション \(Python: get\\_ml\\_task\\_runs\)](#)
- [CancelMLTaskRun アクション \(Python: cancel\\_ml\\_task\\_run\)](#)
- [StartExportLabelsTaskRun アクション \(Python: start\\_export\\_labels\\_task\\_run\)](#)
- [StartImportLabelsTaskRun アクション \(Python: start\\_import\\_labels\\_task\\_run\)](#)

## CreateMLTransform アクション (Python: create\_ml\_transform)

AWS Glue 機械学習変換を作成します。このオペレーションでは、変換を作成します。また、この変換をトレーニングするために必要なすべてのパラメータも作成します。

機械学習変換 (FindMatches 変換など) を使用してデータの重複を除去するプロセスの最初のステップとして、このオペレーションを呼び出します。アルゴリズムのために使用するパラメータに加えて、オプションの Description を指定できます。

データから学習して高品質の機械学習変換を作成する一環として、ユーザーに代わって AWS Glue が実行するタスクに対して特定のパラメータを指定することもできます。これらのパラメータには Role と、必要に応じて AllocatedCapacity、Timeout、MaxRetries が含まれます。詳細については、「[ジョブ](#)」を参照してください。

## リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

変換の作成時に変換に付ける一意の名前。

- Description – 説明文字列、2048 バイト長以下、「[URI address multi-line string pattern](#)」に一致。

定義する機械学習変換の説明。デフォルトは空の文字列です。

- InputRecordTables – 必須: [GlueTable](#) オブジェクトの配列。構造 10 個以下。

変換で使用される AWS Glue テーブル定義のリスト。

- Parameters – 必須: [TransformParameters](#) オブジェクト。

使用する変換タイプに固有のアルゴリズムパラメータ。条件付きで変換タイプに依存します。

- Role – 必須: UTF-8 文字列。

必要なアクセス許可を持つ IAM ロールの Amazon リソースネーム (ARN) の名前。必要なアクセス許可には、AWS Glue リソースに対する AWS Glue サービスロールのアクセス許可と、変換に必要な Amazon S3 アクセス許可の両方が含まれます。

- このロールには、AWS Glue のリソースへのアクセスを許可する AWS Glue サービスロールのアクセス許可が必要です。「[Attach a Policy to IAM Users That Access AWS Glue](#)」を参照してください。
- このロールは、この変換のタスク実行で使用される Amazon Simple Storage Service (Amazon S3) のソース、ターゲット、一時ディレクトリ、スクリプト、およびライブラリへのアクセス許可を必要とします。
- GlueVersion – UTF-8 文字列、1~255 バイト長、「[Custom string pattern #20](#)」に一致。

この値により、この機械学習変換と互換性がある AWS Glue のバージョンが決定します。ほとんどのお客様に、Glue 1.0 が推奨されます。値が設定されていない場合、Glue の互換性はデフォルト

トで Glue 0.9 に設定されます。詳細については、デベロッパーガイドの「[AWS Glue Versions](#)」を参照してください。

- MaxCapacity – 数値 (double)。

この変換のタスク実行に割り当てられる AWS Glue データ処理ユニット (DPU) の数。DPU は、2 ~ 100 の範囲で割り当てることができます。デフォルトは 10 です。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

MaxCapacity は、NumberOfWorkers と WorkerType との相互排他的なオプションです。

- NumberOfWorkers または WorkerType のいずれかが設定されている場合、MaxCapacity は設定できません。
- MaxCapacity が設定されている場合、NumberOfWorkers または WorkerType は設定できません。
- WorkerType が設定されている場合、NumberOfWorkers は必須です ( 逆も同様です ) 。
- MaxCapacity と NumberOfWorkers は両方とも少なくとも 1 である必要があります。

WorkerType フィールドを Standard 以外の値に設定すると、MaxCapacity フィールドが自動的に設定され、読み取り専用になります。

WorkerType フィールドを Standard 以外の値に設定すると、MaxCapacity フィールドが自動的に設定され、読み取り専用になります。

- WorkerType – UTF-8 文字列 (有効な値: Standard="" | G.1X="" | G.2X="" | G.025X="" | G.4X="" | G.8X="" | Z.2X="")。

このタスクの実行時に割り当てられる事前定義済みワーカーのタイプ。使用できる値は、Standard、G.1X、または G.2X です。

- Standard ワーカータイプでは、各ワーカーは 4 vCPU、16 GB のメモリ、50 GB のディスク、ワーカーあたり 2 個のエグゼキュターを提供します。
- G.1X ワーカータイプでは、各ワーカーは 4 vCPU、16 GB のメモリ、64 GB のディスク、ワーカーあたり 1 個のエグゼキュターを提供します。
- G.2X ワーカータイプでは、各ワーカーは 8 vCPU、32 GB のメモリ、128 GB のディスク、ワーカーあたり 1 個のエグゼキュターを提供します。

MaxCapacity は、NumberOfWorkers と WorkerType との相互排他的なオプションです。

- NumberOfWorkers または WorkerType のいずれかが設定されている場合、MaxCapacity は設定できません。
- MaxCapacity が設定されている場合、NumberOfWorkers または WorkerType は設定できません。
- WorkerType が設定されている場合、NumberOfWorkers は必須です ( 逆も同様です ) 。
- MaxCapacity と NumberOfWorkers は両方とも少なくとも 1 である必要があります。
- NumberOfWorkers – 数値 (整数)。

このタスクの実行時に割り当てられる定義済み workerType のワーカー数。

WorkerType が設定されている場合、NumberOfWorkers は必須です ( 逆も同様です ) 。

- Timeout - 数値 (整数)。1 以上。

この変換のタスク実行のタイムアウト (分単位)。これは、この変換のタスク実行が終了して TIMEOUT ステータスに入るまでに、タスク実行がリソースを消費できる最大時間です。デフォルト値は 2,880 分 (48 時間) です。

- MaxRetries – 数値 (整数)。

タスク実行の失敗後に、この変換のタスクを再試行する最大回数。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

この機械学習変換で使用するタグ。タグを使用して、機械学習変換へのアクセスを制限できます。AWS Glue のタグの詳細については、デベロッパーガイドの「[AWS Tags in AWS Glue](#)」を参照してください。

- TransformEncryption – [TransformEncryption](#) オブジェクト。

ユーザーデータへのアクセスに適用される変換の、保管時の暗号化設定。機械学習の変換では、KMS を使用して Amazon S3 で暗号化されたユーザーデータにアクセスできます。

## レスポンス

- TransformId – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

変換用に生成された一意の識別子。

## エラー

- `AlreadyExistsException`
- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`
- `AccessDeniedException`
- `ResourceNumberLimitExceededException`
- `IdempotentParameterMismatchException`

## UpdateMLTransform アクション (Python: `update_ml_transform`)

既存の機械学習変換を更新します。アルゴリズムパラメータを調整して、より良い結果を達成するには、このオペレーションを呼び出します。

このオペレーションを呼び出した後で、`StartMLEvaluationTaskRun` オペレーションを呼び出して新しいパラメータが目標を達成した度合い (機械学習変換の品質やコスト効率の向上など) を評価できます。

### リクエスト

- `TransformId` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
変換の作成時に生成された一意の識別子。
- `Name` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
変換の作成時に付けた一意の変換名。
- `Description` – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
変換の説明。デフォルトは空の文字列です。
- `Parameters` – [TransformParameters](#) オブジェクト。  
使用する変換タイプ (アルゴリズム) に固有の設定パラメータ。条件付きで変換タイプに依存します。
- `Role` – UTF-8 文字列。  
必要なアクセス許可を持つ IAM ロールの Amazon リソースネーム (ARN) の名前。

- `GlueVersion` – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

この値により、この機械学習変換と互換性がある AWS Glue のバージョンが決定します。ほとんどのお客様に、Glue 1.0 が推奨されます。値が設定されていない場合、Glue の互換性はデフォルトで Glue 0.9 に設定されます。詳細については、デベロッパーガイドの「[AWS Glue Versions](#)」を参照してください。

- `MaxCapacity` – 数値 (double)。

この変換のタスク実行に割り当てられる AWS Glue データ処理ユニット (DPU) の数。DPU は、2~100 の範囲で割り当てることができます。デフォルトは 10 です。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

`WorkerType` フィールドを Standard 以外の値に設定すると、`MaxCapacity` フィールドが自動的に設定され、読み取り専用になります。

- `WorkerType` – UTF-8 文字列 (有効な値: Standard="" | G.1X="" | G.2X="" | G.025X="" | G.4X="" | G.8X="" | Z.2X="")。

このタスクの実行時に割り当てられる事前定義済みワーカーのタイプ。使用できる値は、Standard、G.1X、または G.2X です。

- Standard ワーカータイプでは、各ワーカーは 4 vCPU、16 GB のメモリ、50 GB のディスク、ワーカーあたり 2 個のエグゼキューターを提供します。
  - G.1X ワーカータイプでは、各ワーカーは 4 vCPU、16 GB のメモリ、64 GB のディスク、ワーカーあたり 1 個のエグゼキューターを提供します。
  - G.2X ワーカータイプでは、各ワーカーは 8 vCPU、32 GB のメモリ、128 GB のディスク、ワーカーあたり 1 個のエグゼキューターを提供します。
- `NumberOfWorkers` – 数値 (整数)。

このタスクの実行時に割り当てられる定義済み workerType のワーカー数。

- `Timeout` - 数値 (整数)。1 以上。

この変換のタスク実行のタイムアウト (分単位)。これは、この変換のタスク実行が終了して TIMEOUT ステータスに入るまでに、タスク実行がリソースを消費できる最大時間です。デフォルト値は 2,880 分 (48 時間) です。

- `MaxRetries` – 数値 (整数)。

タスク実行の失敗後に、この変換のタスクを再試行する最大回数。

## レスポンス

- TransformId – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
更新された変換の一意的識別子。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException
- AccessDeniedException

## DeleteMLTransform アクション (Python: delete\_ml\_transform)

AWS Glue 機械学習変換を削除します。機械学習変換は、人が提供した例を反映して、実行すべき変換の詳細を機械学習を通じて学習する特殊なタイプの変換です。これらの変換は AWS Glue が保存します。変換が不要になった場合は、DeleteMLTransforms を呼び出して削除できます。ただし、AWS Glue ジョブで削除後の変換を参照しようとする、正常に実行されません。

## リクエスト

- TransformId – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
削除する変換の一意的識別子。

## レスポンス

- TransformId – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
削除した変換の一意的識別子。

## エラー

- EntityNotFoundException
- InvalidInputException

- `OperationTimeoutException`
- `InternalServiceException`

## GetMLTransform アクション (Python: `get_ml_transform`)

AWS Glue 機械学習変換のアーティファクトおよびそのすべての対応するメタデータを取得します。機械学習変換は、人が提供した例を反映して、実行すべき変換の詳細を機械学習を通じて学習する特殊なタイプの変換です。これらの変換は AWS Glue が保存します。GetMLTransform を呼び出して変換のメタデータを取得できます。

### リクエスト

- `TransformId` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
変換の一意的識別子。変換の作成時に生成されます。

### レスポンス

- `TransformId` – UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
変換の一意的識別子。変換の作成時に生成されます。
- `Name` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
変換の作成時に付けた一意の変換名。
- `Description` – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
変換の説明。
- `Status` – UTF-8 文字列 (有効な値: `NOT_READY` | `READY` | `DELETING`)。  
変換の直近の既知のステータス (変換が使用可能かどうかを示します)。有効な値は「`NOT_READY`」、`READY`」、または「`DELETING`」です。
- `CreatedOn` – タイムスタンプ。  
変換の作成日時。
- `LastModifiedOn` – タイムスタンプ。  
変換の最終変更日時。
- `InputRecordTables` – [GlueTable](#) オブジェクトの配列。構造体 10 個以下。

変換で使用する AWS Glue テーブル定義のリスト。

- Parameters – [TransformParameters](#) オブジェクト。

使用するアルゴリズムに固有の設定パラメータ。

- EvaluationMetrics – [EvaluationMetrics](#) オブジェクト。

最新の評価メトリクス。

- LabelCount – 数値 (整数)。

この変換に使用できるラベルの数。

- Schema – [SchemaColumn](#) オブジェクトの配列。構造体 100 個以下。

この変換で使用できるスキーマを表す Map<Column, Type> オブジェクト。100 列が上限です。

- Role – UTF-8 文字列。

必要なアクセス許可を持つ IAM ロールの Amazon リソースネーム (ARN) の名前。

- GlueVersion – UTF-8 文字列、1~255 バイト長、[Custom string pattern #20](#) に一致。

この値により、この機械学習変換と互換性がある AWS Glue のバージョンが決定します。ほとんどのお客様に、Glue 1.0 が推奨されます。値が設定されていない場合、Glue の互換性はデフォルトで Glue 0.9 に設定されます。詳細については、デベロッパーガイドの「[AWS Glue Versions](#)」を参照してください。

- MaxCapacity – 数値 (double)。

この変換のタスク実行に割り当てられる AWS Glue データ処理ユニット (DPU) の数。DPU は、2 ~ 100 の範囲で割り当てることができます。デフォルトは 10 です。DPU は処理能力を相対的に測定するもので、4 個の vCPU のコンピューティング性能と 16 GB のメモリで構成されています。詳細については、「[AWS Glue 料金表ページ](#)」を参照してください。

WorkerType フィールドを Standard 以外の値に設定すると、MaxCapacity フィールドが自動的に設定され、読み取り専用になります。

- WorkerType – UTF-8 文字列 (有効な値: Standard="" | G.1X="" | G.2X="" | G.025X="" | G.4X="" | G.8X="" | Z.2X="")。

このタスクの実行時に割り当てられる事前定義済みワーカーのタイプ。使用できる値は、Standard、G.1X、または G.2X です。

- Standard ワーカータイプでは、各ワーカーは 4 vCPU、16 GB のメモリ、50 GB のディスク、ワーカーあたり 2 個のエグゼキュターを提供します。
- G.1X ワーカータイプでは、各ワーカーは 4 vCPU、16 GB のメモリ、64 GB のディスク、ワーカーあたり 1 個のエグゼキュターを提供します。
- G.2X ワーカータイプでは、各ワーカーは 8 vCPU、32 GB のメモリ、128 GB のディスク、ワーカーあたり 1 個のエグゼキュターを提供します。
- NumberOfWorkers - 数値 (整数)。

このタスクの実行時に割り当てられる定義済み workerType のワーカー数。

- Timeout - 数値 (整数)。1 以上。

この変換のタスク実行のタイムアウト (分単位)。これは、この変換のタスク実行が終了して TIMEOUT ステータスに入るまでに、タスク実行がリソースを消費できる最大時間です。デフォルト値は 2,880 分 (48 時間) です。

- MaxRetries - 数値 (整数)。

タスク実行の失敗後に、この変換のタスクを再試行する最大回数。

- TransformEncryption - [TransformEncryption](#) オブジェクト。

ユーザーデータへのアクセスに適用される変換の、保管時の暗号化設定。機械学習の変換では、KMS を使用して Amazon S3 で暗号化されたユーザーデータにアクセスできます。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## GetMLTransforms アクション (Python: get\_ml\_transforms)

既存の AWS Glue 機械学習変換のソートおよびフィルタリング可能なリストを取得します。機械学習変換は、人が提供した例を反映して、実行すべき変換の詳細を機械学習を通じて学習する特殊なタイプの変換です。これらの変換は AWS Glue が保存します。変換のメタデータを取得するには GetMLTransforms を呼び出します。

## リクエスト

- NextToken – UTF-8 文字列。

結果をオフセットするページ分割されたトークン。

- MaxResults – 1 ~ 1000 の数値 (整数)。

返される結果の最大数。

- Filter – [TransformFilterCriteria](#) オブジェクト。

フィルター変換基準。

- Sort – [TransformSortCriteria](#) オブジェクト。

並べ替え基準。

## レスポンス

- Transforms – 必須: [MLTransform](#) オブジェクトの配列。

機械学習変換のリスト。

- NextToken – UTF-8 文字列。

ページ分割トークン (さらに結果がある場合)。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## ListMLTransforms アクション (Python: list\_ml\_transforms)

この AWS アカウントの既存の AWS Glue 機械学習変換、または指定したタグの付いたリソースの、ソートおよびフィルタリング可能なリストを取得します。このオペレーションはオプションの Tags フィールドを受け取ります。このフィールドをレスポンスのフィルタとして使用すると、タグ

付きリソースをグループとして取得できます。タグフィルタリングの使用を選択した場合は、タグが付いたリソースのみが取得されます。

## リクエスト

- NextToken – UTF-8 文字列。

継続トークン (これが継続リクエストの場合)。

- MaxResults – 1 ~ 1000 の数値 (整数)。

返されるリストの最大サイズ。

- Filter – [TransformFilterCriteria](#) オブジェクト。

機械学習変換のフィルタ処理に使用される TransformFilterCriteria。

- Sort – [TransformSortCriteria](#) オブジェクト。

機械学習変換のソートに使用される TransformSortCriteria。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1 ~ 128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

これらのタグ付きリソースのみを返すように指定します。

## レスポンス

- TransformIds – 必須: UTF-8 文字列の配列。

アカウントのすべての機械学習変換、または指定したタグの付いた機械学習変換の識別子。

- NextToken – UTF-8 文字列。

継続トークン (戻されたリストに最後に使用可能なメトリクスが含まれていない場合)。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException

- `InternalServiceException`

## StartMLEvaluationTaskRun アクション (Python: `start_ml_evaluation_task_run`)

変換の品質を推定するためのタスクを開始します。

真の例としてラベルセットを指定すると、AWS Glue 機械学習は、これらの例のいくつかから学習します。残りのラベルは、品質を推定するためのテストとして使用されます。

実行用の一意の識別子を返します。GetMLTaskRun を呼び出して、EvaluationTaskRun の状態に関する詳細を取得できます。

### リクエスト

- `TransformId` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

機械学習変換の一意の識別子。

### レスポンス

- `TaskRunId` – UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

この実行に関連付けられた一意の識別子。

### エラー

- `EntityNotFoundException`
- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`
- `ConcurrentRunsExceededException`
- `MLTransformNotReadyException`

## StartMLLabelingSetGenerationTaskRun アクション (Python: start\_ml\_labeling\_set\_generation\_task\_run)

機械学習変換のアクティブな学習ワークフローを開始し、ラベルセットを生成してラベルを追加することで、変換の品質を向上させます。

StartMLLabelingSetGenerationTaskRun が完了すると、AWS Glue によって「ラベリングセット」または人が回答する一連の質問が生成されます。

FindMatches 変換の場合、これらの質問は「これらの行を、一致するレコードのみで構成されるグループに分ける適切な方法は何ですか?」という形式になります。

ラベリングプロセスが完了すると、StartImportLabelsTaskRun を呼び出してラベルをアップロードできます。StartImportLabelsTaskRun が完了すると、今後のすべての機械学習変換の実行で新規および改善されたラベルが使用され、高品質の変換が実行されます。

### リクエスト

- TransformId – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

機械学習変換の一意の識別子。

- OutputS3Path – 必須: UTF-8 文字列。

ラベリングセットを生成する Amazon Simple Storage Service (Amazon S3) のパス。

### レスポンス

- TaskRunId – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

このタスクに関連付けられている一意の実行識別子。

### エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException
- ConcurrentRunsExceededException

## GetMLTaskRun アクション (Python: get\_ml\_task\_run)

機械学習変換に対する特定のタスク実行の詳細を取得します。機械学習のタスク実行は、さまざまな機械学習ワークフローの一環として AWS Glue がユーザーに代わって実行する非同期タスクです。任意のタスク実行の統計情報を確認するには、TaskRunID およびその親変換の TransformID を使用して GetMLTaskRun を呼び出します。

### リクエスト

- TransformId – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

機械学習変換の一意の識別子。

- TaskRunId – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

タスク実行の一意の識別子。

### レスポンス

- TransformId – UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

タスク実行の一意の識別子。

- TaskRunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

この実行に関連付けられた一意の実行識別子。

- Status – UTF-8 文字列 (有効な値: STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT)。

このタスク実行のステータス。

- LogGroupName – UTF-8 文字列。

タスク実行に関連付けられているロググループの名前。

- Properties – [TaskRunProperties](#) オブジェクト。

タスク実行に関連付けられているプロパティのリスト。

- ErrorMessage – UTF-8 文字列。

タスク実行に関連付けられているエラー文字列。

- StartedOn – タイムスタンプ。

このタスク実行の開始日時。

- LastModifiedOn – タイムスタンプ。

このタスク実行の最終変更日時。

- CompletedOn – タイムスタンプ。

このタスク実行の完了日時。

- ExecutionTime – 数値 (整数)。

タスク実行でリソースを消費した時間 (秒)。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## GetMLTaskRuns アクション (Python: get\_ml\_task\_runs)

機械学習変換の実行のリストを取得します。機械学習のタスク実行は、さまざまな機械学習ワークフローの一環として AWS Glue がユーザーに代わって実行する非同期タスクです。機械学習のタスク実行の並べ替え可能およびフィルタリング可能なリストは、GetMLTaskRuns を呼び出して取得できます。この呼び出しで、親変換の TransformID とこのセクションに記載しているその他のオプションのパラメータを使用します。

このオペレーションは、実行履歴のリストを返します。また、ページ分割する必要があります。

## リクエスト

- TransformId – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

機械学習変換の一意の識別子。

- NextToken – UTF-8 文字列。

結果をページ分割するためのトークン。デフォルトは空です。

- `MaxResults` – 1 ~ 1000 の数値 (整数)。

返される結果の最大数。

- `Filter` – [TaskRunFilterCriteria](#) オブジェクト。

タスク実行のフィルタリング基準 (`TaskRunFilterCriteria` 構造)。

- `Sort` – [TaskRunSortCriteria](#) オブジェクト。

タスク実行の並べ替え基準件 (`TaskRunSortCriteria` 構造)。

## 応答

- `TaskRuns` – [TaskRun](#) オブジェクトの配列。

変換に関連付けられているタスク実行のリスト。

- `NextToken` – UTF-8 文字列。

ページ分割トークン (さらに結果がある場合)。

## エラー

- `EntityNotFoundException`
- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`

## CancelMLTaskRun アクション (Python: `cancel_ml_task_run`)

タスク実行をキャンセル (停止) します。機械学習のタスク実行は、さまざまな機械学習ワークフローの一環として AWS Glue がユーザーに代わって実行する非同期タスクです。機械学習のタスク実行は、`CancelMLTaskRun` を呼び出していつでもキャンセルできます。この呼び出しで、タスク実行の親変換の `TransformID` とタスク実行の `TaskRunId` を使用します。

## リクエスト

- `TransformId` – 必須: UTF-8 文字列、1 ~ 255 バイト長、[「Single-line string pattern」](#) に一致。

機械学習変換の一意の識別子。

- TaskRunId – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

タスク実行に対する一意の識別子。

## レスポンス

- TransformId – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

機械学習変換の一意の識別子。

- TaskRunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

タスク実行の一意の識別子。

- Status – UTF-8 文字列 (有効な値: STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT)。

この実行のステータス。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## StartExportLabelsTaskRun アクション (Python: start\_export\_labels\_task\_run)

特定の変換に関するすべてのラベル付きデータをエクスポートするための非同期タスクを開始します。このタスクは、通常のアクティブな学習ワークフローに属さない唯一のラベル関連の API コールです。すべての既存のラベルを同時に操作するときは、通常、StartExportLabelsTaskRun を使用します (以前に真として送信したラベルを削除または変更する場合など)。この API オペレーションでは、エクスポートするラベルの TransformId と、ラベルをエクスポートする先の Amazon Simple Storage Service (Amazon S3) のパスを使用できます。このオペレーションは、TaskRunId を返します。タスク実行のステータスは、GetMLTaskRun API を呼び出して確認できます。

## リクエスト

- TransformId – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

機械学習変換の一意の識別子。

- OutputS3Path – 必須: UTF-8 文字列。

ラベルをエクスポートする先となる Simple Storage Service (Amazon S3) へのパス。

## レスポンス

- TaskRunId – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

タスク実行の一意の識別子。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## StartImportLabelsTaskRun アクション (Python: start\_import\_labels\_task\_run)

機械学習変換をトレーニングしてその品質を向上させるために使用する追加のラベル (真の例) を提供できます。この API オペレーションは、通常、StartMLLabelingSetGenerationTaskRun の呼び出しで始まり、最終的に機械学習変換の品質向上につながるアクティブな学習ワークフローの一環として使用します。

StartMLLabelingSetGenerationTaskRun が完了すると、AWS Glue 機械学習によって人が回答すべき一連の質問が生成されます。(これらの質問に回答することは、機械学習ワークフローにおいて「ラベリング」とも呼ばれます)。FindMatches 変換の場合、これらの質問は「これらの行を、一致するレコードのみで構成されるグループに分ける適切な方法は何ですか?」という形式になります。ラベリングプロセスが完了すると、ユーザーは StartImportLabelsTaskRun を呼び出して回答/ラベルをアップロードできます。StartImportLabelsTaskRun が完了すると、今後のすべ

での機械学習変換の実行で新規および改善されたラベルが使用され、より高品質の変換が実行されます。

デフォルトでは、StartMLLabelingSetGenerationTaskRun は引き続き学習を継続し、ユーザーがアップロードしたすべてのラベルを結合します。ただし、Replace を true に設定した場合を除きます。Replace を true に設定すると、StartImportLabelsTaskRun は以前にアップロードされたすべてのラベルを削除して無視し、ユーザーがアップロードする正確なセットからのみ学習します。ラベルの置き換えが役立つのは、以前に間違ったラベルをアップロードしたことに気づき、これらが変換の品質に悪影響を及ぼしていると思われる場合です。

タスク実行のステータスは、GetMLTaskRun オペレーションを呼び出して確認できます。

### リクエスト

- TransformId – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

機械学習変換の一意の識別子。

- InputS3Path – 必須: UTF-8 文字列。

ラベルをインポートする元となる Amazon Simple Storage Service (Amazon S3) へのパス。

- ReplaceAllLabels – ブール。

既存のラベルを上書きするかどうかを示します。

### レスポンス

- TaskRunId – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

タスク実行の一意の識別子。

### エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- ResourceNumberLimitExceededException
- InternalServiceException

# Data Quality API

Data Quality API は、データ品質のデータ型を提供し、データ品質ルールセット、実行、および評価を作成、削除、または更新するための API を含んでいます。

## データ型

- [DataSource 構造](#)
- [DataQualityRulesetListDetails 構造](#)
- [DataQualityTargetTable 構造](#)
- [DataQualityRulesetEvaluationRunDescription 構造](#)
- [DataQualityRulesetEvaluationRunFilter 構造](#)
- [DataQualityEvaluationRunAdditionalRunOptions 構造](#)
- [DataQualityRuleRecommendationRunDescription 構造](#)
- [DataQualityRuleRecommendationRunFilter 構造](#)
- [DataQualityResult 構造](#)
- [DataQualityAnalyzerResult 構造](#)
- [DataQualityObservation 構造](#)
- [MetricBasedObservation 構造](#)
- [DataQualityMetricValues 構造](#)
- [DataQualityRuleResult 構造](#)
- [DataQualityResultDescription 構造](#)
- [DataQualityResultFilterCriteria 構造](#)
- [DataQualityRulesetFilterCriteria 構造](#)

## DataSource 構造

データ品質結果が必要なデータソース ( AWS Glue テーブル ) 。

### フィールド

- GlueTable – 必須: [GlueTable](#) オブジェクト。

AWS Glue テーブル。

## DataQualityRulesetListDetails 構造

GetDataQualityRuleset によって返されるデータ品質ルールセットを記述します。

フィールド

- Name – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
データ品質ルールセットの名前。
- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
データ品質ルールセットの説明。
- CreatedOn – タイムスタンプ。  
データ品質ルールセットの作成日時。
- LastModifiedOn – タイムスタンプ。  
データ品質ルールセットの最終変更日時。
- TargetTable – [DataQualityTargetTable](#) オブジェクト。  
AWS Glue テーブルを表すオブジェクト。
- RecommendationRunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
推奨実行からルールセットが作成されると、この 2 つをリンクさせるためにこの実行 ID が生成されます。
- RuleCount – 数値 (整数)。  
ルールセット内のルールの数。

## DataQualityTargetTable 構造

AWS Glue テーブルを表すオブジェクト。

フィールド

- TableName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

AWS Glue テーブルの名前。

- DatabaseName – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

AWS Glue テーブルが存在するデータベースの名前。

- CatalogId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

AWS Glue テーブルが存在するカタログ ID。

## DataQualityRulesetEvaluationRunDescription 構造

データ品質ルールセットの評価実行の結果を記述します。

フィールド

- RunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

この実行に関連付けられた一意の実行識別子。

- Status – UTF-8 文字列 (有効な値: STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT)。

この実行のステータス。

- StartedOn – タイムスタンプ。

実行が開始された日時。

- DataSource – [DataSource](#) オブジェクト。

実行に関連付けられたデータソース (AWS Glue テーブル)。

## DataQualityRulesetEvaluationRunFilter 構造

フィルタリング基準。

フィールド

- DataSource – 必須: [DataSource](#) オブジェクト。

実行に関連付けられたデータソース (AWS Glue テーブル) に基づいてフィルタリングします。

- StartedBefore – タイムスタンプ。

この時刻より前に開始された実行で結果をフィルタリングします。

- `StartedAfter` – タイムスタンプ。

この時刻以降に開始された実行で結果をフィルタリングします。

## DataQualityEvaluationRunAdditionalRunOptions 構造

評価実行で指定できるその他の実行オプション。

フィールド

- `CloudWatchMetricsEnabled` – ブール。

CloudWatch メトリクスを有効にするかどうか。

- `ResultsS3Prefix` – UTF-8 文字列。

結果を保存する Amazon S3 のプレフィックス。

- `CompositeRuleEvaluationMethod` – UTF-8 文字列 (有効な値: COLUMN | ROW)。

ルールセットの複合ルールの評価方法を ROW/COLUMN に設定する

## DataQualityRuleRecommendationRunDescription 構造

データ品質ルールの推奨実行の結果を記述します。

フィールド

- `RunId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

この実行に関連付けられた一意の実行識別子。

- `Status` – UTF-8 文字列 (有効な値: STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT)。

この実行のステータス。

- `StartedOn` – タイムスタンプ。

この実行の開始日時。

- `DataSource` – [DataSource](#) オブジェクト。

レコメンデーション実行に関連付けられたデータソース (AWS Glue テーブル)。

## DataQualityRuleRecommendationRunFilter 構造

データ品質の推奨実行を一覧表示するためのフィルター。

フィールド

- DataSource – 必須: [DataSource](#) オブジェクト。

指定されたデータソース (AWS Glue テーブル) に基づいてフィルタリングします。

- StartedBefore – タイムスタンプ。

指定した時刻より前に開始された結果の場合、時刻に基づいてフィルタリングします。

- StartedAfter – タイムスタンプ。

指定した時刻以降に開始された結果の場合、時刻に基づいてフィルタリングします。

## DataQualityResult 構造

データ品質結果を記述します。

フィールド

- ResultId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データ品質結果に対する一意の結果 ID。

- Score – 数値 (double)。1.0 以下。

集計されたデータ品質のスコア。ルールの合計数に対する合格ルールの比率を表します。

- DataSource – [DataSource](#) オブジェクト。

データ品質結果に関連付けられたテーブル (存在する場合)。

- RulesetName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データ品質結果に関連付けられたルールセットの名前。

- EvaluationContext – UTF-8 文字列。

AWS Glue Studio のジョブのコンテキストでは、キャンバス内の各ノードには通常、何らかの名前が割り当てられ、データ品質ノードには名前が割り当てられます。複数のノードの場合、`evaluationContext` でノードを区別できます。

- `StartedOn` – タイムスタンプ。

このデータ品質の実行の開始日時。

- `CompletedOn` – タイムスタンプ。

このデータ品質の実行の完了日時。

- `JobName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データ品質結果に関連付けられたジョブ名 (存在する場合)。

- `JobRunId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データ品質結果に関連付けられたジョブ実行 ID (存在する場合)。

- `RulesetEvaluationRunId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このデータ品質結果のルールセット評価に対する一意の実行 ID。

- `RuleResults` – [DataQualityRuleResult](#) オブジェクトの配列。構造 2000 個以下。

各ルールの結果を表す `DataQualityRuleResult` オブジェクトのリスト。

- `AnalyzerResults` – [DataQualityAnalyzerResult](#) オブジェクトの配列。構造 2000 個以下。

各アナライザーの結果を表す `DataQualityAnalyzerResult` オブジェクトのリスト。

- `Observations` – [DataQualityObservation](#) オブジェクトの配列。構造 50 個以下。

`DataQualityObservation` ルールとアナライザーを評価した後に生成された観測値を表すオブジェクトのリスト。

## DataQualityAnalyzerResult 構造

データ品質アナライザーの評価の結果を説明します。

フィールド

- `Name` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データ品質アナライザーの名前。

- Description - UTF-8 文字列。2,048 バイト長以下。 [URI address multi-line string pattern](#) に一致。

データ品質アナライザーの説明。

- EvaluationMessage - UTF-8 文字列。2,048 バイト長以下。 [URI address multi-line string pattern](#) に一致。

評価メッセージ。

- EvaluatedMetrics - キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1~255 バイト長で、 [Single-line string pattern](#) に一致します。

各値は数値 (double) です。

アナライザーの評価に関連付けられているメトリクスのマップ。

## DataQualityObservation 構造

ルールおよびアナライザーを評価した後に生成される観測値について説明します。

フィールド

- Description - UTF-8 文字列。2,048 バイト長以下。 [URI address multi-line string pattern](#) に一致。

データ品質観察の説明。

- MetricBasedObservation - [MetricBasedObservation](#) オブジェクト。

評価されたデータ品質メトリクスに基づく観測値を表す MetricBasedObservation タイプのオブジェクト。

## MetricBasedObservation 構造

評価されたデータ品質メトリクスに基づいて生成されたメトリクスベースの観測値について説明します。

## フィールド

- **MetricName** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
観測値を生成するために使用されるデータ品質メトリクスの名前。
- **MetricValues** – [DataQualityMetricValues](#) オブジェクト。  
データ品質メトリクス値の分析を表す DataQualityMetricValues タイプのオブジェクト。
- **NewRules** – UTF-8 文字列の配列。  
データ品質メトリクス値に基づく観測値の一部として生成された新しいデータ品質ルールの一覧。

## DataQualityMetricValues 構造

履歴データの分析に基づくデータ品質メトリクス値を説明します。

### フィールド

- **ActualValue** – 数値 (double)。  
データ品質メトリクスの実際の値。
- **ExpectedValue** – 数値 (double)。  
履歴データの分析に基づくデータ品質メトリクスの期待値。
- **LowerLimit** – 数値 (double)。  
履歴データの分析に基づくデータ品質メトリクスの下限値。
- **UpperLimit** – 数値 (double)。  
履歴データの分析に基づくデータ品質メトリクスの上限値。

## DataQualityRuleResult 構造

データ品質ルールの評価の結果を記述します。

### フィールド

- **Name** – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データ品質ルールの名前。

- **Description** - UTF-8 文字列。2,048 バイト長以下。 [URI address multi-line string pattern](#) に一致。

データ品質ルールの説明。

- **EvaluationMessage** - UTF-8 文字列。2,048 バイト長以下。 [URI address multi-line string pattern](#) に一致。

評価メッセージ。

- **Result** - UTF-8 文字列 (有効な値: PASS | FAIL | ERROR)。

ルールの合格または不合格のステータス。

- **EvaluatedMetrics** - キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1~255 バイト長で、 [Single-line string pattern](#) に一致します。

各値は数値 (double) です。

ルールの評価に関連するメトリクスのマップ。

## DataQualityResultDescription 構造

データ品質結果を記述します。

フィールド

- **ResultId** - UTF-8 文字列、1~255 バイト長、 [Single-line string pattern](#) に一致。

このデータ品質結果に対する一意の結果 ID。

- **DataSource** - [DataSource](#) オブジェクト。

データ品質結果に関連付けられたテーブル名。

- **JobName** - UTF-8 文字列、1~255 バイト長、 [Single-line string pattern](#) に一致。

データ品質結果に関連付けられたジョブ名。

- **JobRunId** - UTF-8 文字列、1~255 バイト長、 [Single-line string pattern](#) に一致。

データ品質結果に関連付けられたジョブ実行 ID。

- `StartedOn` – タイムスタンプ。

このデータ品質結果の実行が開始された時刻。

## DataQualityResultFilterCriteria 構造

データ品質結果を返すために使用される基準。

フィールド

- `DataSource` – [DataSource](#) オブジェクト。

指定したデータソースで結果をフィルタリングします。例えば、AWS Glue テーブルのすべての結果を取得します。

- `JobName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

指定したジョブ名で結果をフィルタリングします。

- `JobRunId` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

指定したジョブ実行 ID で結果をフィルタリングします。

- `StartedAfter` – タイムスタンプ。

この時刻以降に開始された実行で結果をフィルタリングします。

- `StartedBefore` – タイムスタンプ。

この時刻より前に開始された実行で結果をフィルタリングします。

## DataQualityRulesetFilterCriteria 構造

データ品質ルールセットのフィルタリングに使用される基準。

フィールド

- `Name` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ルールセットのフィルタリング基準の名前。

- `Description` – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

ルールセットのフィルタリング基準の説明。

- `CreatedBefore` – タイムスタンプ。

この日付より前に作成されたルールセットをフィルタリングします。

- `CreatedAfter` – タイムスタンプ。

この日付以降に作成されたルールセットをフィルタリングします。

- `LastModifiedBefore` – タイムスタンプ。

この日付より前に最終変更されたルールセットをフィルタリングします。

- `LastModifiedAfter` – タイムスタンプ。

この日付より後に最終変更されたルールセットをフィルタリングします。

- `TargetTable` – [DataQualityTargetTable](#) オブジェクト。

ターゲットテーブルの名前およびデータベース名。

## 操作

- [StartDataQualityRulesetEvaluationRun アクション \(Python: `start\_data\_quality\_ruleset\_evaluation\_run`\)](#)
- [CancelDataQualityRulesetEvaluationRun アクション \(Python: `cancel\_data\_quality\_ruleset\_evaluation\_run`\)](#)
- [GetDataQualityRulesetEvaluationRun アクション \(Python: `get\_data\_quality\_ruleset\_evaluation\_run`\)](#)
- [ListDataQualityRulesetEvaluationRuns アクション \(Python: `list\_data\_quality\_ruleset\_evaluation\_runs`\)](#)
- [StartDataQualityRuleRecommendationRun アクション \(Python: `start\_data\_quality\_rule\_recommendation\_run`\)](#)
- [CancelDataQualityRuleRecommendationRun アクション \(Python: `cancel\_data\_quality\_rule\_recommendation\_run`\)](#)
- [GetDataQualityRuleRecommendationRun アクション \(Python: `get\_data\_quality\_rule\_recommendation\_run`\)](#)
- [ListDataQualityRuleRecommendationRuns アクション \(Python: `list\_data\_quality\_rule\_recommendation\_runs`\)](#)
- [GetDataQualityResult アクション \(Python: `get\_data\_quality\_result`\)](#)

- [BatchGetDataQualityResult アクション \(Python: batch\\_get\\_data\\_quality\\_result\)](#)
- [ListDataQualityResults アクション \(Python: list\\_data\\_quality\\_results\)](#)
- [CreateDataQualityRuleset アクション \(Python: create\\_data\\_quality\\_ruleset\)](#)
- [DeleteDataQualityRuleset アクション \(Python: delete\\_data\\_quality\\_ruleset\)](#)
- [GetDataQualityRuleset アクション \(Python: get\\_data\\_quality\\_ruleset\)](#)
- [ListDataQualityRulesets アクション \(Python: list\\_data\\_quality\\_rulesets\)](#)
- [UpdateDataQualityRuleset アクション \(Python: update\\_data\\_quality\\_ruleset\)](#)

## StartDataQualityRulesetEvaluationRun アクション (Python: start\_data\_quality\_ruleset\_evaluation\_run )

ルールセット定義 (推奨または独自の) を取得したら、このオペレーションを呼び出して、データソース (AWS Glue テーブル) に対してルールセットを評価します。GetDataQualityResult API で取得できる結果は評価で計算されます。

### リクエスト

- DataSource – 必須: [DataSource](#) オブジェクト。

この実行に関連付けられたデータソース (AWS Glue テーブル)。

- Role – 必須: UTF-8 文字列。

実行の結果を暗号化するために提供される IAM ロール。

- NumberOfWorkers – 数値 (整数)。

実行に使用される G.1X ワーカーの数。デフォルトは 5 です。

- Timeout - 数値 (整数)。1 以上。

実行のタイムアウト (分)。これは、実行が終了済みになって TIMEOUT ステータスに入るまでに、その実行でリソースを消費できる最大時間です。デフォルト値は 2,880 分 (48 時間) です。

- ClientToken – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

冪等性を保つために使用されるもので、同じリソースの複数のインスタンスを作成または起動しないように、ランダムな ID (UUID など) に設定することをお勧めします。

- AdditionalRunOptions – [DataQualityEvaluationRunAdditionalRunOptions](#) オブジェクト。

評価実行で指定できるその他の実行オプション。

- RulesetNames – 必須: UTF-8 文字列の配列、1～10 個の文字列。

ルールセット名のリスト。

- AdditionalDataSources – キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1～255 バイト長で、[Single-line string pattern](#) に一致します。

各値は [DataSource](#) オブジェクトです。

評価実行で指定できる他のデータソースへの参照文字列のマップ。

## レスポンス

- RunId – UTF-8 文字列、1～255 バイト長、「[Single-line string pattern](#)」に一致。

この実行に関連付けられた一意の実行識別子。

## エラー

- InvalidInputException
- EntityNotFoundException
- OperationTimeoutException
- InternalServiceException
- ConflictException

## CancelDataQualityRulesetEvaluationRun アクション (Python: cancel\_data\_quality\_ruleset\_evaluation\_run )

ルールセットがデータソースに対して評価されている実行をキャンセルします。

## リクエスト

- RunId – 必須: UTF-8 文字列、1～255 バイト長、「[Single-line string pattern](#)」に一致。

この実行に関連付けられた一意の実行識別子。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## GetDataQualityRulesetEvaluationRun アクション (Python: `get_data_quality_ruleset_evaluation_run`)

ルールセットがデータソースに対して評価される特定の実行を取得します。

### リクエスト

- RunId – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

この実行に関連付けられた一意の実行識別子。

### レスポンス

- RunId – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

この実行に関連付けられた一意の実行識別子。

- DataSource – [DataSource](#) オブジェクト。

この評価実行に関連付けられたデータソース (AWS Glue テーブル)。

- Role – UTF-8 文字列。

実行の結果を暗号化するために提供される IAM ロール。

- NumberOfWorkers – 数値 (整数)。

実行に使用される G.1X ワーカーの数。デフォルトは 5 です。

- Timeout - 数値 (整数)。1 以上。

実行のタイムアウト (分)。これは、実行が終了済みになって TIMEOUT ステータスに入るまでに、その実行でリソースを消費できる最大時間です。デフォルト値は 2,880 分 (48 時間) です。

- `AdditionalRunOptions` – [DataQualityEvaluationRunAdditionalRunOptions](#) オブジェクト。

評価実行で指定できるその他の実行オプション。

- `Status` – UTF-8 文字列 (有効な値: STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT)。

この実行のステータス。

- `ErrorMessage` – UTF-8 文字列。

実行に関連付けられているエラー文字列。

- `StartedOn` – タイムスタンプ。

この実行の開始日時。

- `LastModifiedOn` – タイムスタンプ。

タイムスタンプ。このデータ品質ルールの推奨実行が変更された最後の時間ポイント。

- `CompletedOn` – タイムスタンプ。

この実行の完了日時。

- `ExecutionTime` – 数値 (整数)。

実行でリソースを消費した時間 (秒)。

- `RulesetNames` – UTF-8 文字列の配列、1 ~ 10 個の文字列。

実行用のルールセット名のリスト。現在、このパラメータは 1 つのルールセット名しかとりません。

- `ResultIds` – UTF-8 文字列の配列、1 ~ 10 個の文字列。

実行用のデータ品質結果の結果 ID のリスト。

- `AdditionalDataSources` – キーバリューペアのマップ配列。

各キーは UTF-8 文字列、1 ~ 255 バイト長で、[Single-line string pattern](#) に一致します。

各値は [DataSource](#) オブジェクトです。

評価実行で指定できる他のデータソースへの参照文字列のマップ。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## ListDataQualityRulesetEvaluationRuns アクション (Python: list\_data\_quality\_ruleset\_evaluation\_runs )

フィルタリング基準を満たすすべての実行が一覧表示され、ルールセットがデータソースに対して評価されます。

### リクエスト

- Filter – [DataQualityRulesetEvaluationRunFilter](#) オブジェクト。

フィルタリング基準。

- NextToken – UTF-8 文字列。

結果をオフセットするページ分割されたトークン。

- MaxResults – 1 ~ 1000 の数値 (整数)。

返される結果の最大数。

### 応答

- Runs – [DataQualityRulesetEvaluationRunDescription](#) オブジェクトの配列。

データ品質ルールセットの実行を表す DataQualityRulesetEvaluationRunDescription オブジェクトのリスト。

- NextToken – UTF-8 文字列。

ページ分割トークン (さらに結果がある場合)。

## エラー

- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`

## StartDataQualityRuleRecommendationRun アクション (Python: `start_data_quality_rule_recommendation_run`)

どのルールを書き込むべきかわからないときにルールを生成するために使用されるレコメンデーション実行を開始します。AWS Glue Data Quality はデータを分析し、潜在的なルールセットのレコメンデーションを作成します。その後、ルールセットをトリアーژیし、生成されたルールセットを好みに合わせて変更できます。

推奨事項の実行は、実行から 90 日後に自動的に削除されます。

### リクエスト

- `DataSource` – 必須: [DataSource](#) オブジェクト。

この実行に関連付けられたデータソース (AWS Glue テーブル)。

- `Role` – 必須: UTF-8 文字列。

実行の結果を暗号化するために提供される IAM ロール。

- `NumberOfWorkers` – 数値 (整数)。

実行に使用される G.1X ワーカーの数。デフォルトは 5 です。

- `Timeout` - 数値 (整数)。1 以上。

実行のタイムアウト (分)。これは、実行が終了済みになって `TIMEOUT` ステータスに入るまでに、その実行でリソースを消費できる最大時間です。デフォルト値は 2,880 分 (48 時間) です。

- `CreatedRulesetName` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ルールセットの名前。

- `ClientToken` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

冪等性を保つために使用されるもので、同じリソースの複数のインスタンスを作成または起動しないように、ランダムな ID (UUID など) に設定することをお勧めします。

## レスポンス

- RunId – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

この実行に関連付けられた一意の実行識別子。

## エラー

- InvalidInputException
- OperationTimeoutException
- InternalServiceException
- ConflictException

## CancelDataQualityRuleRecommendationRun アクション (Python: `cancel_data_quality_rule_recommendation_run`)

ルールの生成に使用されていた指定の推奨実行をキャンセルします。

## リクエスト

- RunId – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

この実行に関連付けられた一意の実行識別子。

## レスポンス

- 応答パラメータはありません。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## GetDataQualityRuleRecommendationRun アクション (Python: `get_data_quality_rule_recommendation_run`)

ルールの生成に使用された指定の推奨実行を取得します。

### リクエスト

- `RunId` – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
この実行に関連付けられた一意の実行識別子。

### レスポンス

- `RunId` – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。  
この実行に関連付けられた一意の実行識別子。
- `DataSource` – [DataSource](#) オブジェクト。  
この実行に関連付けられたデータソース (AWS Glue テーブル)。
- `Role` – UTF-8 文字列。  
実行の結果を暗号化するために提供される IAM ロール。
- `NumberOfWorkers` – 数値 (整数)。  
実行に使用される G.1X ワーカーの数。デフォルトは 5 です。
- `Timeout` - 数値 (整数)。1 以上。  
実行のタイムアウト (分)。これは、実行が終了済みになって TIMEOUT ステータスに入るまでに、その実行でリソースを消費できる最大時間です。デフォルト値は 2,880 分 (48 時間) です。
- `Status` – UTF-8 文字列 (有効な値: STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT)。  
この実行のステータス。
- `ErrorString` – UTF-8 文字列。  
実行に関連付けられているエラー文字列。
- `StartedOn` – タイムスタンプ。  
この実行の開始日時。

- LastModifiedOn – タイムスタンプ。

タイムスタンプ。このデータ品質ルールの推奨実行が変更された最後の時間ポイント。

- CompletedOn – タイムスタンプ。

この実行の完了日時。

- ExecutionTime – 数値 (整数)。

実行でリソースを消費した時間 (秒)。

- RecommendedRuleset – UTF-8 文字列、1~65536 バイト長。

開始ルールの推奨実行が完了すると、推奨ルールセット (ルールのセット) が作成されます。このメンバーには、これらのルールがデータ品質定義言語 (DQDL) 形式で格納されています。

- CreatedRulesetName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

実行によって作成されたルールセットの名前。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## ListDataQualityRuleRecommendationRuns アクション (Python: list\_data\_quality\_rule\_recommendation\_runs )

フィルタリング基準を満たす推奨実行を一覧表示します。

### リクエスト

- Filter – [DataQualityRuleRecommendationRunFilter](#) オブジェクト。

フィルタリング基準。

- NextToken – UTF-8 文字列。

結果をオフセットするページ分割されたトークン。

- `MaxResults` – 1 ~ 1000 の数値 (整数)。

返される結果の最大数。

## 応答

- `Runs` – [DataQualityRuleRecommendationRunDescription](#) オブジェクトの配列。  
DataQualityRuleRecommendationRunDescription オブジェクトのリスト。
- `NextToken` – UTF-8 文字列。  
ページ分割トークン (さらに結果がある場合)。

## エラー

- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`

## GetDataQualityResult アクション (Python: `get_data_quality_result`)

データ品質ルール評価の結果を取得します。

### リクエスト

- `ResultId` – 必須: UTF-8 文字列、1 ~ 255 バイト長、「[Single-line string pattern](#)」に一致。  
データ品質結果に対する一意の結果 ID。

### レスポンス

- `ResultId` – UTF-8 文字列、1 ~ 255 バイト長、「[Single-line string pattern](#)」に一致。  
データ品質結果に対する一意の結果 ID。
- `Score` – 数値 (double)。1.0 以下。  
集計されたデータ品質のスコア。ルールの合計数に対する合格ルールの比率を表します。
- `DataSource` – [DataSource](#) オブジェクト。

データ品質結果に関連付けられたテーブル (存在する場合)。

- RulesetName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データ品質結果に関連付けられたルールセットの名前。

- EvaluationContext – UTF-8 文字列。

AWS Glue Studio のジョブのコンテキストでは、キャンバス内の各ノードには通常、何らかの名前が割り当てられ、データ品質ノードには名前が割り当てられます。複数のノードの場合、evaluationContext でノードを区別できます。

- StartedOn – タイムスタンプ。

このデータ品質結果の実行が開始された日時。

- CompletedOn – タイムスタンプ。

このデータ品質結果の実行が完了した日時。

- JobName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データ品質結果に関連付けられたジョブ名 (存在する場合)。

- JobRunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

データ品質結果に関連付けられたジョブ実行 ID (存在する場合)。

- RulesetEvaluationRunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

ルールセット評価に関連付けられた一意の実行 ID。

- RuleResults – [DataQualityRuleResult](#) オブジェクトの配列。構造 2000 個以下。

各ルールの結果を表す DataQualityRuleResult オブジェクトのリスト。

- AnalyzerResults – [DataQualityAnalyzerResult](#) オブジェクトの配列。構造 2000 個以下。

各アナライザーの結果を表す DataQualityAnalyzerResult オブジェクトのリスト。

- Observations – [DataQualityObservation](#) オブジェクトの配列。構造 50 個以下。

DataQualityObservation ルールとアナライザーを評価した後に生成された観測値を表すオブジェクトのリスト。

## エラー

- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`
- `EntityNotFoundException`

## BatchGetDataQualityResult アクション (Python: `batch_get_data_quality_result`)

指定した結果 ID のデータ品質結果のリストを取得します。

### リクエスト

- `ResultIds` – 必須: UTF-8 文字列の配列、1~100 個の文字列。

データ品質結果に対する一意の結果 ID のリスト。

### レスポンス

- `Results` – 必須: [DataQualityResult](#) オブジェクトの配列。

データ品質結果を表す `DataQualityResult` オブジェクトのリスト。

- `ResultsNotFound` – UTF-8 文字列の配列、1~100 個の文字列。

結果が見つからなかった結果 ID のリスト。

## エラー

- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`

## ListDataQualityResults アクション (Python: `list_data_quality_results`)

アカウントのすべてのデータ品質実行結果を返します。

## リクエスト

- `Filter` – [DataQualityResultFilterCriteria](#) オブジェクト。

フィルタリング基準。

- `NextToken` – UTF-8 文字列。

結果をオフセットするページ分割されたトークン。

- `MaxResults` – 1 ~ 1000 の数値 (整数)。

返される結果の最大数。

## レスポンス

- `Results` – 必須: [DataQualityResultDescription](#) オブジェクトの配列。

`DataQualityResultDescription` オブジェクトのリスト。

- `NextToken` – UTF-8 文字列。

ページ分割トークン (さらに結果がある場合)。

## エラー

- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`

## CreateDataQualityRuleset アクション (Python: `create_data_quality_ruleset`)

指定された AWS Glue テーブルに適用される DQDL ルールを使用して、データ品質ルールセットを作成します。

データ品質定義言語 (DQDL) を使用して、ルールセットを作成します。詳細については、「[デ AWS Glue ベロッパーガイド](#)」を参照してください。

## リクエスト

- `Name` – 必須: UTF-8 文字列、1 ~ 255 バイト長、[「Single-line string pattern」](#) に一致。

データ品質ルールセットに対する一意の名前。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

データ品質ルールセットの説明。

- Ruleset – 必須: UTF-8 文字列、1~65536 バイト長。

データ品質定義言語 (DQDL) ルールセット。詳細については、「[デ AWS Glue ベロツパーガイド](#)」を参照してください。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

データ品質ルールセットに適用されるタグのリスト。

- TargetTable – [DataQualityTargetTable](#) オブジェクト。

データ品質ルールセットに関連付けられたターゲットテーブル。

- RecommendationRunId – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

推奨実行に対する一意の実行 ID。

- ClientToken – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

冪等性を保つために使用されるもので、同じリソースの複数のインスタンスを作成または起動しないように、ランダムな ID (UUID など) に設定することをお勧めします。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

データ品質ルールセットに対する一意の名前。

## エラー

- InvalidInputException
- AlreadyExistsException
- OperationTimeoutException

- `InternalServiceException`
- `ResourceNumberLimitExceededException`

## DeleteDataQualityRuleset アクション (Python: `delete_data_quality_ruleset`)

データ品質ルールセットを削除します。

### リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

データ品質ルールセットの名前。

### レスポンス

- 応答パラメータはありません。

### エラー

- `EntityNotFoundException`
- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`

## GetDataQualityRuleset アクション (Python: `get_data_quality_ruleset`)

既存のルールセットを識別子または名前で返します。

### リクエスト

- `Name` – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

ルールセットの名前。

### レスポンス

- `Name` – UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

ルールセットの名前。

- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

ルールセットの説明。

- Ruleset – UTF-8 文字列、1 ~ 65536 バイト長。

データ品質定義言語 (DQDL) ルールセット。詳細については、「[デ AWS Glue ベロツパーガイド](#)」を参照してください。

- TargetTable – [DataQualityTargetTable](#) オブジェクト。

ターゲットテーブルの名前およびデータベース名。

- CreatedOn – タイムスタンプ。

タイムスタンプ。このデータ品質ルールセットが作成された日時。

- LastModifiedOn – タイムスタンプ。

タイムスタンプ。このデータ品質ルールセットが変更された最後の時間ポイント。

- RecommendationRunId – UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。

推奨実行からルールセットが作成されると、この 2 つをリンクさせるためにこの実行 ID が生成されます。

## エラー

- EntityNotFoundException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException

## ListDataQualityRulesets アクション (Python: list\_data\_quality\_rulesets)

指定された AWS Glue テーブルのリストのルールセットのページ分割されたリストを返します。

## リクエスト

- NextToken – UTF-8 文字列。

結果をオフセットするページ分割されたトークン。

- `MaxResults` – 1 ~ 1000 の数値 (整数)。

返される結果の最大数。

- `Filter` – [DataQualityRulesetFilterCriteria](#) オブジェクト。

フィルタリング基準。

- `Tags` – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1 ~ 128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

キーと値のペアタグのリスト。

## 応答

- `Rulesets` – [DataQualityRulesetListDetails](#) オブジェクトの配列。

指定された AWS Glue テーブルのリストのルールセットのページ分割されたリスト。

- `NextToken` – UTF-8 文字列。

ページ分割トークン (さらに結果がある場合)。

## エラー

- `EntityNotFoundException`
- `InvalidInputException`
- `OperationTimeoutException`
- `InternalServiceException`

## UpdateDataQualityRuleset アクション (Python: `update_data_quality_ruleset`)

指定したデータ品質ルールセットを更新します。

## リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
データ品質ルールセットの名前。
- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
ルールセットの説明。
- Ruleset – UTF-8 文字列、1~65536 バイト長。  
データ品質定義言語 (DQDL) ルールセット。詳細については、「[デ AWS Glue ベロツパーガイド](#)」を参照してください。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
データ品質ルールセットの名前。
- Description – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。  
ルールセットの説明。
- Ruleset – UTF-8 文字列、1~65536 バイト長。  
データ品質定義言語 (DQDL) ルールセット。詳細については、「[デ AWS Glue ベロツパーガイド](#)」を参照してください。

## エラー

- EntityNotFoundException
- AlreadyExistsException
- IdempotentParameterMismatchException
- InvalidInputException
- OperationTimeoutException
- InternalServiceException
- ResourceNumberLimitExceededException

# 機密データの検出 API

機密データ検出 API では、構造化データの列と行全体にある機密データを検出するために使用される API について説明します。

## データ型

- [CustomentityType 構造](#)

## CustomentityType 構造

構造化データの列と行全体にある機密データを検出するためのカスタムパターンを表すオブジェクトです。

### フィールド

- Name – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

後で、取得したり削除したりすることができるカスタムパターンの名前です。この名前は AWS アカウントごとに一意である必要があります。

- RegexString – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

カスタムパターン内にある機密データの検出に使用される正規表現の文字列。

- ContextWords – UTF-8 文字列の配列、1~20 個の文字列。

コンテキスト単語のリスト。正規表現の範囲内でこれらのコンテキスト単語が見つからない場合、データは機密データとして検出されません。

コンテキスト単語が渡されない場合は、正規表現のみがチェックされます。

## 操作

- [CreateCustomType アクション \(Python: create\\_custom\\_entity\\_type\)](#)
- [DeleteCustomEntityType アクション \(Python: delete\\_custom\\_entity\\_type\)](#)
- [GetCustomEntityType アクション \(Python: get\\_custom\\_entity\\_type\)](#)
- [BatchGetCustomEntityTypes アクション \(Python: batch\\_get\\_custom\\_entity\\_types\)](#)
- [ListCustomTypes アクション \(Python: list\\_custom\\_entity\\_types\)](#)

## CreateCustomType アクション (Python: create\_custom\_entity\_type)

構造化データの列と行全体にある機密データを検出するために使用されるカスタムパターンを作成します。

作成するカスタムパターンごとに、正規表現とオプションのコンテキスト単語リストを作成し、指定します。コンテキスト単語が渡されない場合は、正規表現のみがチェックされます。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

後で、取得したり削除したりすることができるカスタムパターンの名前です。この名前は AWS アカウントごとに一意である必要があります。

- RegexString – 必須: UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

カスタムパターン内にある機密データの検出に使用される正規表現の文字列。

- ContextWords – UTF-8 文字列の配列、1~20 個の文字列。

コンテキスト単語のリスト。正規表現の範囲内でこれらのコンテキスト単語が見つからない場合、データは機密データとして検出されません。

コンテキスト単語が渡されない場合は、正規表現のみがチェックされます。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

カスタムエンティティタイプに適用されるタグのリスト。

### レスポンス

- Name – UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。

作成したカスタムパターンの名前。

### エラー

- AccessDeniedException

- AlreadyExistsException
- IdempotentParameterMismatchException
- InternalServiceException
- InvalidInputException
- OperationTimeoutException
- ResourceNumberLimitExceededException

## DeleteCustomEntityType アクション (Python: delete\_custom\_entity\_type)

カスタムパターン名を指定して、カスタムパターンを削除します。

### リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
削除するカスタムパターンの名前です。

### レスポンス

- Name – UTF-8 文字列、1~255 バイト長、[「Single-line string pattern」](#) に一致。  
削除したカスタムパターンの名前です。

### エラー

- EntityNotFoundException
- AccessDeniedException
- InternalServiceException
- InvalidInputException
- OperationTimeoutException

## GetCustomEntityType アクション (Python: get\_custom\_entity\_type)

カスタムパターン名を指定して、カスタムパターンの詳細を取得します。

## リクエスト

- Name – 必須: UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

取得するカスタムパターンの名前。

## レスポンス

- Name – UTF-8 文字列、1~255 バイト長、「[Single-line string pattern](#)」に一致。

取得したカスタムパターンの名前。

- RegexString – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

カスタムパターン内にある機密データの検出に使用される正規表現の文字列。

- ContextWords – UTF-8 文字列の配列、1~20 個の文字列。

カスタムパターンを作成した時に指定している場合は、コンテキスト単語のリストです。正規表現の範囲内でこれらのコンテキスト単語が見つからない場合、データは機密データとして検出されません。

## エラー

- EntityNotFoundException
- AccessDeniedException
- InternalServiceException
- InvalidInputException
- OperationTimeoutException

## BatchGetCustomEntityTypes アクション (Python: batch\_get\_custom\_entity\_types)

名前リストで指定されたカスタムパターンの詳細を取得します。

## リクエスト

- Names – 必須: UTF-8 文字列の配列、1~50 個の文字列。

取得するカスタムパターンの名前リストです。

## 応答

- CustomEntityTypes – [CustomEntityType](#) オブジェクトの配列。

作成したカスタムパターンを表す CustomEntityType オブジェクトのリストです。

- CustomEntityTypesNotFound – UTF-8 文字列の配列、1~50 個の文字列。

検出できなかったカスタムパターンの名前リストです。

## エラー

- InvalidInputException
- InternalServiceException
- OperationTimeoutException

## ListCustomEntityTypes アクション (Python: list\_custom\_entity\_types)

作成したすべてのカスタムパターンをリスト表示します。

### リクエスト

- NextToken – UTF-8 文字列。

結果をオフセットするページ分割されたトークン。

- MaxResults – 1~1000 の数値 (整数)。

返される結果の最大数。

- Tags – キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

キーと値のペアタグのリスト。

## 応答

- CustomEntityTypes – [CustomEntityType](#) オブジェクトの配列。  
カスタムパターンを表す CustomEntityType オブジェクトのリストです。
- NextToken – UTF-8 文字列。  
ページ分割トークン (さらに結果がある場合)。

## エラー

- InvalidInputException
- OperationTimeoutException
- InternalServiceException

# AWS Glue の API のタグ付け

## データ型

- [Tag 構造](#)

## Tag 構造

Tag オブジェクトは、AWS リソースに割り当てることができるラベルを表します。タグはそれぞれ、1つのキーとオプションの1つの値で設定されており、どちらもお客様側が定義します。

リソース AWS Glue へのアクセス制御や、[AWS タグ AWS Glue](#) および [指定 AWS Glue リソース ARN](#) の詳細については、開発者ガイドを参照してください。

## フィールド

- key – UTF-8 文字列、1~128 バイト長。  
タグキー。オブジェクトにタグを作成するときにキーが必要です。キーでは大文字と小文字が区別され、プレフィックス aws を含めることはできません。
- value – UTF-8 文字列、256 バイト長以下。

タグ値。オブジェクトにタグを作成するときの値はオプションです。値では大文字と小文字が区別され、プレフィックス `aws` を含めることはできません。

## 操作

- [TagResource アクション \(Python: tag\\_resource\)](#)
- [UntagResource アクション \(Python: untag\\_resource\)](#)
- [GetTags アクション \(Python: get\\_tags\)](#)

## TagResource アクション (Python: tag\_resource)

リソースにタグを追加します。タグとは、AWS のリソースに付けることができるラベルです。AWS Glue では、特定のリソースにのみタグを付けることができます。どのリソースにタグを付けることができるかについては、「[AWS Tags in AWS Glue](#)」を参照してください。

タグ関連の API を呼び出すためのタグ付けのアクセス許可に加えて、接続でタグ付け API を呼び出すための `glue:GetConnection` アクセス許可と、データベースでタグ付け API を呼び出すための `glue:GetDatabase` アクセス許可も必要です。

### リクエスト

- `ResourceArn` – 必須: UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。

タグを追加する AWS Glue リソースの ARN。AWS Glue リソースの ARN の詳細については、「[AWS Glue ARN string pattern](#)」を参照してください。

- `TagsToAdd` – 必須: キーと値のペアのマップ配列。50 ペア以下。

各キーは UTF-8 文字列で、1~128 バイト長です。

各値は UTF-8 文字列で、256 バイト長以下です。

このリソースに追加するタグ。

### レスポンス

- 応答パラメータはありません。

## エラー

- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`
- `EntityNotFoundException`

## UntagResource アクション (Python: `untag_resource`)

リソースからタグを削除します。

### リクエスト

- `ResourceArn` – 必須: UTF-8 文字列。1~10240 バイト長。 [Custom string pattern #22](#) に一致。

タグを削除するリソースの Amazon リソースネーム (ARN)。

- `TagsToRemove` – 必須: UTF-8 文字列の配列。文字列 50 個以下。

このリソースから削除するタグ。

### レスポンス

- 応答パラメータはありません。

## エラー

- `InvalidInputException`
- `InternalServiceException`
- `OperationTimeoutException`
- `EntityNotFoundException`

## GetTags アクション (Python: `get_tags`)

リソースに関連付けられているタグのリストを取得します。

## リクエスト

- ResourceArn – 必須: UTF-8 文字列。1~10240 バイト長。[Custom string pattern #22](#) に一致。  
タグを取得する対象のリソースの Amazon リソースネーム (ARN)。

## レスポンス

- Tags – キーと値のペアのマップ配列。50 ペア以下。  
各キーは UTF-8 文字列で、1~128 バイト長です。  
各値は UTF-8 文字列で、256 バイト長以下です。  
リクエストされたタグ。

## エラー

- InvalidInputException
- InternalServiceException
- OperationTimeoutException
- EntityNotFoundException

## 一般的なデータ型

一般的なデータ型では、AWS Glueでのその他の一般的なデータ型について説明します。

## Tag 構造

Tag オブジェクトは、AWS リソースに割り当てることができるラベルを表します。タグはそれぞれ、1つのキーとオプションの1つの値で設定されており、どちらもお客様側が定義します。

タグと のリソースへのアクセスの制御の詳細については、「[デベロッパーガイドAWS](#)」の「[のタグ AWS Glue](#)」および[AWS Glue 「リソース ARNs AWS Glue」](#)を参照してください。

## フィールド

- key – UTF-8 文字列、1~128 バイト長。

タグキー。オブジェクトにタグを作成するときにキーが必要です。キーでは大文字と小文字が区別され、プレフィックス `aws` を含めることはできません。

- `value` - UTF-8 文字列、256 バイト長以下。

タグ値。オブジェクトにタグを作成するときの値はオプションです。値では大文字と小文字が区別され、プレフィックス `aws` を含めることはできません。

## DecimalNumber 構造

10 進数形式の数値が含まれます。

フィールド

- `UnscaledValue` – 必須: Blob。

スケールされていない数値。

- `Scale` – 必須: 数値 (integer)。

スケールされていない値のどの位置に小数点を置くかを決定するスケール。

## ErrorDetail 構造

エラーに関する詳細が含まれています。

フィールド

- `ErrorCode` – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。

このエラーに関連付けられたコード。

- `ErrorMessage` – 説明文字列、2048 バイト長以下、[URI address multi-line string pattern](#) に一致。

エラーを説明するメッセージ。

## PropertyPredicate 構造

プロパティの述語を定義します。

## フィールド

- **Key – Value** 文字列、1024 バイト長以下。

プロパティのキー。

- **Value – Value** 文字列、1024 バイト長以下。

プロパティの値。

- **Comparator** – UTF-8 文字列 (有効な値: EQUALS | GREATER\_THAN | LESS\_THAN | GREATER\_THAN\_EQUALS | LESS\_THAN\_EQUALS)。

このプロパティを他のプロパティと比較するために使用されたコンパレータ。

## ResourceUri 構造

関数リソースの URI。

### フィールド

- **ResourceType** – UTF-8 文字列 (有効な値: JAR | FILE | ARCHIVE)。

リソースのタイプ。

- **Uri** – Uniform resource identifier (uri)、1 ~ 1024 バイト長、[URI address multi-line string pattern](#) に一致。

リソースにアクセスするための URI。

## ColumnStatistics 構造

テーブルまたはパーティションに対して生成された列レベルの統計を表します。

### フィールド

- **ColumnName** – 必須: UTF-8 文字列、1 ~ 255 バイト長、[Single-line string pattern](#) に一致。

統計が属する列の名前。

- **ColumnType** – 必須: 型名。20000 バイト長以下。[Single-line string pattern](#) に一致。

列のデータ型。

- AnalyzedTime – 必須: タイムスタンプ。  
列統計が生成された時点を示すタイムスタンプ。
- StatisticsData – 必須: [ColumnStatisticsData](#) オブジェクト。  
統計データ値を含む ColumnStatisticData オブジェクト。

## ColumnStatisticsError 構造

失敗した ColumnStatistics オブジェクトと失敗の理由をカプセル化します。

フィールド

- ColumnStatistics – [ColumnStatistics](#) オブジェクト。  
列の ColumnStatistics。
- Error – [ErrorDetail](#) オブジェクト。  
オペレーションが失敗した理由を示すエラーメッセージ。

## ColumnError 構造

エラーがある列の名前とエラーの理由について要約します。

フィールド

- ColumnName – UTF-8 文字列、1~255 バイト長、[Single-line string pattern](#) に一致。  
エラーがある列の名前。
- Error – [ErrorDetail](#) オブジェクト。  
オペレーションが失敗した理由を示すエラーメッセージ。

## ColumnStatisticsData 構造

個々のタイプの列統計データが含まれます。Type 属性で 1 つのデータオブジェクトのみを設定し、指示します。

## フィールド

- **Type** – 必須: UTF-8 文字列 (有効な値: BOOLEAN | DATE | DECIMAL | DOUBLE | LONG | STRING | BINARY)。

列の統計データの型。

- **BooleanColumnStatisticsData** – [BooleanColumnStatisticsData](#) オブジェクト。

Boolean の列統計データ。

- **DateColumnStatisticsData** – [DateColumnStatisticsData](#) オブジェクト。

Date の列統計データ。

- **DecimalColumnStatisticsData** – [DecimalColumnStatisticsData](#) オブジェクト。

10 進列統計 data.UnscaledValues within は Base64-encoded バイナリオブジェクトで、10 進数のスケールリングされていない値の 2 の補数表現であるビッグエンディアンが格納されます。

- **DoubleColumnStatisticsData** – [DoubleColumnStatisticsData](#) オブジェクト。

Double の列統計データ。

- **LongColumnStatisticsData** – [LongColumnStatisticsData](#) オブジェクト。

Long の列統計データ。

- **StringColumnStatisticsData** – [StringColumnStatisticsData](#) オブジェクト。

文字列の列統計データ。

- **BinaryColumnStatisticsData** – [BinaryColumnStatisticsData](#) オブジェクト。

Binary の列統計データ。

## BooleanColumnStatisticsData 構造

Boolean データ列でサポートされる列統計を定義します。

### フィールド

- **NumberOfTrues** – 必須: 数値 (long)。None 以下。

列内の true の値の数。

- **NumberOfFalses** – 必須: 数値 (long)。None 以下。

列内の false の値の数。

- `NumberOfNulls` – 必須: 数値 (long)。None 以下。

列内の null 値の数。

## DateColumnStatisticsData 構造

タイムスタンプデータ列でサポートされている列統計を定義します。

フィールド

- `MinimumValue` – タイムスタンプ。

列内の最小値。

- `MaximumValue` – タイムスタンプ。

列内の最大値。

- `NumberOfNulls` – 必須: 数値 (long)。None 以下。

列内の null 値の数。

- `NumberOfDistinctValues` – 必須: 数値 (long)。None 以下。

列内の異なる値の数です。

## DecimalColumnStatisticsData 構造

固定小数点数データ列でサポートされている列統計を定義します。

フィールド

- `MinimumValue` – [DecimalNumber](#) オブジェクト。

列内の最小値。

- `MaximumValue` – [DecimalNumber](#) オブジェクト。

列内の最大値。

- `NumberOfNulls` – 必須: 数値 (long)。None 以下。

列内の null 値の数。

- `NumberOfDistinctValues` – 必須: 数値 (long)。None 以下。

列内の異なる値の数です。

## DoubleColumnStatisticsData 構造

浮動小数点数データ列でサポートされている列統計を定義します。

フィールド

- `MinimumValue` – 数値 (double)。

列内の最小値。

- `MaximumValue` – 数値 (double)。

列内の最大値。

- `NumberOfNulls` – 必須: 数値 (long)。None 以下。

列内の null 値の数。

- `NumberOfDistinctValues` – 必須: 数値 (long)。None 以下。

列内の異なる値の数です。

## LongColumnStatisticsData 構造

整数データ列でサポートされる列統計を定義します。

フィールド

- `MinimumValue` – 数値 (long 型)。

列内の最小値。

- `MaximumValue` – 数値 (long 型)。

列内の最大値。

- `NumberOfNulls` – 必須: 数値 (long)。None 以下。

列内の null 値の数。

- `NumberOfDistinctValues` – 必須: 数値 (long)。None 以下。

列内の異なる値の数です。

## StringColumnStatisticsData 構造

文字シーケンスのデータ値に対してサポートされる列統計を定義します。

フィールド

- `MaxLength` – 必須: 数値 (long)。None 以下。

列内の最も長い文字列のサイズです。

- `AverageLength` – 必須: 数値 (double)。None 以下。

列内の文字列の平均長。

- `NumberOfNulls` – 必須: 数値 (long)。None 以下。

列内の null 値の数。

- `NumberOfDistinctValues` – 必須: 数値 (long)。None 以下。

列内の異なる値の数です。

## BinaryColumnStatisticsData 構造

ビットシーケンスデータ値に対してサポートされる列統計を定義します。

フィールド

- `MaxLength` – 必須: 数値 (long)。None 以下。

列内の最も長いビットシーケンスのサイズ。

- `AverageLength` – 必須: 数値 (double)。None 以下。

列内の平均ビットシーケンス長。

- `NumberOfNulls` – 必須: 数値 (long)。None 以下。

列内の null 値の数。

## 文字列パターン

API は、さまざまな文字列パラメータとメンバーの有効コンテンツを定義するために、以下の正規表現を使用します。

- 単一行文字列パターン - 「`[\u0020-\uD7FF\uE000-\uFFFF\uD800\uDC00-\uDBFF\uDFFF\t]*`」
- URI アドレスの複数行文字列パターン - 「`[\u0020-\uD7FF\uE000-\uFFFF\uD800\uDC00-\uDBFF\uDFFF\r\n\t]*`」
- Logstash Grok 文字列パターン - 「`[\u0020-\uD7FF\uE000-\uFFFF\uD800\uDC00-\uDBFF\uDFFF\r\t]*`」
- 識別子文字列パターン - 「`[A-Za-z_][A-Za-z0-9_]*`」
- AWS IAM ARN 文字列パターン - `"arn:aws:iam::\d{12}:role/.*"`
- バージョン文字列パターン - 「`^[a-zA-Z0-9-_]+$`」
- ロググループ文字列パターン - 「`[\.\-_\#A-Za-z0-9]+`」
- ログストリーム文字列パターン - 「`^[^:]*`」
- カスタム文字列パターン #10 - `"^[r\n]"`
- カスタム文字列パターン #11 - `"^arn:aws(-(cn|us-gov|iso(-[bef]))?)?:secretsmanager:.*$"`
- カスタム文字列パターン #12 - `"^(https?):\/\/[-a-zA-Z0-9+@#/%=?~_!|:,.;]*[-a-zA-Z0-9+@#/%=?~_]"`
- カスタム文字列パターン #13 - `"\S+"`
- カスタム文字列パターン #14 - `"^(https?):\//\[^\s/$.?\#].\[^\s]*$"`
- カスタム文字列パターン #15 - `"^subnet-[a-z0-9]+$"`
- カスタム文字列パターン #16 - `"[\p{L}\p{N}\p{P}]*"`
- カスタム文字列パターン #17 - `"[a-f0-9]{8}-[a-f0-9]{4}-[a-f0-9]{4}-[a-f0-9]{4}-[a-f0-9]{12}"`
- カスタム文字列パターン #18 - `"[a-zA-Z0-9-_$#.]+"`
- カスタム文字列パターン #19 - `"^\\w+\\.\\w+\\.\\w+$"`

- カスタム文字列パターン #20 – `^\w+\.\w+$`
- カスタム文字列パターン #21 – `^([2-3]|3[.]9)$`
- カスタム文字列パターン #22 – `arn:(aws|aws-us-gov|aws-cn):glue:.*`
- カスタム文字列パターン #23 – `(^arn:aws:iam::\w{12}:root)"`
- カスタム文字列パターン #24 – `^arn:aws(-(cn|us-gov|iso(-[bef]))?)?:iam::[0-9]{12}:role/.+)"`
- カスタム文字列パターン #25 – `arn:aws:kms:.*`
- カスタム文字列パターン #26 – `arn:aws[^:]*:iam::[0-9]*:role/.+)"`
- カスタム文字列パターン #27 – `[\.\-_\A-Za-z0-9]+"`
- カスタム文字列パターン #28 – `^s3://([^\s/]+)/([^\s/]+)*([^\s/]+)$"`
- カスタム文字列パターン #29 – `.*`
- カスタム文字列パターン #30 – `^(Sun|Mon|Tue|Wed|Thu|Fri|Sat):([01]?[0-9]|2[0-3])$"`
- カスタム文字列パターン #31 – `[a-zA-Z0-9_.-]+"`
- カスタム文字列パターン #32 – `.*\S.*"`
- カスタム文字列パターン #33 – `[a-zA-Z0-9-=_./@]+"`
- カスタム文字列パターン #34 – `[1-9][0-9]*|[1-9][0-9]*-[1-9][0-9]*"`
- カスタム文字列パターン #35 – `[\s\S]*"`
- カスタム文字列パターン #36 – `([\u0020-\uD7FF\uE000-\uFFFF\uD800\uDC00-\uDBFF\uDFFF]|[\^\S\r\n"'= ;])*"`
- カスタム文字列パターン #37 – `[*A-Za-z0-9_-]*"`
- カスタム文字列パターン #38 – `([\u0020-\u007E\r\s\n])*"`
- カスタム文字列パターン #39 – `[A-Za-z0-9_-]*"`
- カスタム文字列パターン #40 – `([\u0020-\uD7FF\uE000-\uFFFF\uD800\uDC00-\uDBFF\uDFFF]|[\^\S\r\n"' ])*"`
- カスタム文字列パターン #41 – `([\u0020-\uD7FF\uE000-\uFFFF\uD800\uDC00-\uDBFF\uDFFF]|[\^\S\r\n])*"`
- カスタム文字列パターン #42 – `([\u0020-\uD7FF\uE000-\uFFFF\uD800\uDC00-\uDBFF\uDFFF\s])*"`
- カスタム文字列パターン #43 – `([\u0020-\uD7FF\uE000-\uFFFF\uD800\uDC00-\uDBFF\uDFFF]|[\^\r\n])*"`

## 例外

このセクションでは、問題の原因を見つけて解決するために使用できる AWS Glue の例外について説明します。機械学習に関連する例外の HTTP エラーコードと文字列の詳細については、「[the section called “AWS Glue 機械学習の例外”](#)」を参照してください。

### AccessDeniedException 構造

リソースへのアクセスが拒否されました。

フィールド

- Message – UTF-8 文字列。  
問題を説明するメッセージ。

### AlreadyExistsException 構造

作成または追加するリソースはすでに存在します。

フィールド

- Message – UTF-8 文字列。  
問題を説明するメッセージ。

### ConcurrentModificationException 構造

2 つのプロセスが同時にリソースを変更しようとしています。

フィールド

- Message – UTF-8 文字列。  
問題を説明するメッセージ。

### ConcurrentRunsExceededException 構造

同時に実行されているジョブが多すぎます。

## フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## CrawlerNotRunningException 構造

指定されたクローラーが実行されていません。

### フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## CrawlerRunningException 構造

クローラーが既に実行されているため、操作を実行できません。

### フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## CrawlerStoppingException 構造

指定されたクローラーが停止しています。

### フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## EntityNotFoundException 構造

指定されたエンティティは存在しません

## フィールド

- Message – UTF-8 文字列。  
問題を説明するメッセージ。
- FromFederationSource – ブール。  
例外がフェデレーションされたソースに関連するものかどうかを示します。

## FederationSourceException 構造

ソースのフェデレーションに失敗しました。

### フィールド

- FederationSourceErrorCode – UTF-8 文字列 (有効な値:AccessDeniedException | EntityNotFoundException | InvalidCredentialsException | InvalidInputException | InvalidResponseException | OperationTimeoutException | OperationNotSupportedException | InternalServiceException | PartialFailureException | ThrottlingException)。  
問題のエラーコード。
- Message – UTF-8 文字列。  
問題を説明するメッセージ。

## FederationSourceRetryableException 構造

ソースのフェデレーションに失敗しましたが、操作を再試行できます。

### フィールド

- Message – UTF-8 文字列。  
問題を説明するメッセージ。

## GlueEncryptionException 構造

暗号化オペレーションが失敗しました。

## フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## IdempotentParameterMismatchException 構造

同じ一意の識別子が 2 つの異なるレコードに関連付けられていました。

## フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## IllegalWorkflowStateException 構造

リクエストされたオペレーションを実行するには、ワークフローの状態が無効です。

## フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## InternalServiceException 構造

内部サービスエラーが発生しました。

## フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## InvalidExecutionEngineException 構造

不明または無効な実行エンジンが指定されました。

## フィールド

- message – UTF-8 文字列。

問題を説明するメッセージ。

## InvalidInputException 構造

指定された入力は無効です。

### フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

- FromFederationSource – ブール。

例外がフェデレーションされたソースに関連するものかどうかを示します。

## InvalidStateException 構造

データが無効な状態であることを示すエラー。

### フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## InvalidTaskStatusTransitionException 構造

あるタスクから次のタスクへの適切な移行に失敗しました。

### フィールド

- message – UTF-8 文字列。

問題を説明するメッセージ。

## JobDefinitionErrorException 構造

ジョブ定義が無効です。

フィールド

- message – UTF-8 文字列。  
問題を説明するメッセージ。

## JobRunInTerminalStateException 構造

ジョブ実行の終了状態が障害を通知しています。

フィールド

- message – UTF-8 文字列。  
問題を説明するメッセージ。

## JobRunInvalidStateTransitionException 構造

ジョブ実行で、ソース状態からターゲット状態への無効な移行が発生しました。

フィールド

- jobRunId – UTF-8 文字列、1～255 バイト長、[Single-line string pattern](#) に一致。  
該当するジョブ実行の ID。
- message – UTF-8 文字列。  
問題を説明するメッセージ。
- sourceState – UTF-8 文字列 (有効な値:STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT | ERROR | WAITING | EXPIRED)。  
ソースの状態。
- targetState – UTF-8 文字列 (有効な値:STARTING | RUNNING | STOPPING | STOPPED | SUCCEEDED | FAILED | TIMEOUT | ERROR | WAITING | EXPIRED)。  
ターゲットの状態。

## JobRunNotInTerminalStateException 構造

ジョブ実行が終了状態ではありません。

フィールド

- message – UTF-8 文字列。

問題を説明するメッセージ。

## LateRunnerException 構造

ジョブランナーが遅延しています。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## NoScheduleException 構造

該当するスケジュールがありません。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## OperationTimeoutException 構造

オペレーションがタイムアウトしました。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## ResourceNotReadyException 構造

リソースのトランザクション準備ができていませんでした。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## ResourceNumberLimitExceededException 構造

リソースの数値制限を超えました。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## SchedulerNotRunningException 構造

指定されたスケジューラが実行されていません。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## SchedulerRunningException 構造

指定されたスケジューラは既に実行されています。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## SchedulerTransitioningException 構造

指定されたスケジューラが移行中です。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## UnrecognizedRunnerException 構造

ジョブランナーが認識されませんでした。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## ValidationException 構造

値を検証できませんでした。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

## VersionMismatchException 構造

バージョンの競合がありました。

フィールド

- Message – UTF-8 文字列。

問題を説明するメッセージ。

# AWS Glue AWS SDKs を使用した API コード例

次のコード例は、Software Development Kit (SDK) AWS Glue で を使用する方法を示しています。

アクションはより大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。アクションは個々のサービス機能呼び出す方法を示していますが、関連するシナリオやサービス間の例ではアクションのコンテキストが確認できます。

「シナリオ」は、同じサービス内で複数の関数を呼び出して、特定のタスクを実行する方法を示すコード例です。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## 開始方法

### こんにちは AWS Glue は

次のコード例は、AWS Glue の使用を開始する方法を示しています。

## .NET

### AWS SDK for .NET

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
namespace GlueActions;

public class HelloGlue
{
 private static ILogger logger = null!;

 static async Task Main(string[] args)
 {
 // Set up dependency injection for AWS Glue.
```

```
using var host = Host.CreateDefaultBuilder(args)
 .ConfigureLogging(logging =>
 logging.AddFilter("System", LogLevel.Debug)
 .AddFilter<DebugLoggerProvider>("Microsoft",
 LogLevel.Information)
 .AddFilter<ConsoleLoggerProvider>("Microsoft",
 LogLevel.Trace))
 .ConfigureServices((_, services) =>
 services.AddAWSService<IAmazonGlue>()
 .AddTransient<GlueWrapper>()
)
 .Build();

logger = LoggerFactory.Create(builder => { builder.AddConsole(); })
 .CreateLogger<HelloGlue>();
var glueClient = host.Services.GetRequiredService<IAmazonGlue>();

var request = new ListJobsRequest();

var jobNames = new List<string>();

do
{
 var response = await glueClient.ListJobsAsync(request);
 jobNames.AddRange(response.JobNames);
 request.NextToken = response.NextToken;
}
while (request.NextToken is not null);

Console.Clear();
Console.WriteLine("Hello, Glue. Let's list your existing Glue Jobs:");
if (jobNames.Count == 0)
{
 Console.WriteLine("You don't have any AWS Glue jobs.");
}
else
{
 jobNames.ForEach(Console.WriteLine);
}
}
```

- APIの詳細については、「API リファレンス [ListJobs](#)」の「」を参照してください。AWS SDK for .NET

## C++

### SDK for C++

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

C MakeLists.txt CMake ファイルのコード。

```
Set the minimum required version of CMake for this project.
cmake_minimum_required(VERSION 3.13)

Set the AWS service components used by this project.
set(SERVICE_COMPONENTS glue)

Set this project's name.
project("hello_glue")

Set the C++ standard to use to build this target.
At least C++ 11 is required for the AWS SDK for C++.
set(CMAKE_CXX_STANDARD 11)

Use the MSVC variable to determine if this is a Windows build.
set(WINDOWS_BUILD ${MSVC})

if (WINDOWS_BUILD) # Set the location where CMake can find the installed
 libraries for the AWS SDK.
 string(REPLACE ";" "/aws-cpp-sdk-all;" SYSTEM_MODULE_PATH
 "${CMAKE_SYSTEM_PREFIX_PATH}/aws-cpp-sdk-all")
 list(APPEND CMAKE_PREFIX_PATH ${SYSTEM_MODULE_PATH})
endif ()

Find the AWS SDK for C++ package.
find_package(AWSSDK REQUIRED COMPONENTS ${SERVICE_COMPONENTS})

if (WINDOWS_BUILD AND AWSSDK_INSTALL_AS_SHARED_LIBS)
```

```
Copy relevant AWS SDK for C++ libraries into the current binary directory
for running and debugging.

set(BIN_SUB_DIR "/Debug") # if you are building from the command line you
may need to uncomment this
 # and set the proper subdirectory to the
executables' location.

 AWSSDK_CPY_DYN_LIBS(SERVICE_COMPONENTS ""
 ${CMAKE_CURRENT_BINARY_DIR}${BIN_SUB_DIR})
endif ()

add_executable(${PROJECT_NAME}
 hello_glue.cpp)

target_link_libraries(${PROJECT_NAME}
 ${AWSSDK_LINK_LIBRARIES})
```

hello\_glue.cpp ソースファイルのコード。

```
#include <aws/core/Aws.h>
#include <aws/glue/GlueClient.h>
#include <aws/glue/model/ListJobsRequest.h>
#include <iostream>

/*
 * A "Hello Glue" starter application which initializes an AWS Glue client and
 lists the
 * AWS Glue job definitions.
 *
 * main function
 *
 * Usage: 'hello_glue'
 *
 */

int main(int argc, char **argv) {
 Aws::SDKOptions options;
 // Optionally change the log level for debugging.
 // options.loggingOptions.logLevel = Utils::Logging::LogLevel::Debug;
 Aws::InitAPI(options); // Should only be called once.
 int result = 0;
```

```
{
 Aws::Client::ClientConfiguration clientConfig;
 // Optional: Set to the AWS Region (overrides config file).
 // clientConfig.region = "us-east-1";

 Aws::Glue::GlueClient glueClient(clientConfig);

 std::vector<Aws::String> jobs;

 Aws::String nextToken; // Used for pagination.
 do {
 Aws::Glue::Model::ListJobsRequest listJobsRequest;
 if (!nextToken.empty()) {
 listJobsRequest.SetNextToken(nextToken);
 }

 Aws::Glue::Model::ListJobsOutcome listRunsOutcome =
glueClient.ListJobs(
 listJobsRequest);

 if (listRunsOutcome.IsSuccess()) {
 const std::vector<Aws::String> &jobNames =
listRunsOutcome.GetResult().GetJobNames();
 jobs.insert(jobs.end(), jobNames.begin(), jobNames.end());

 nextToken = listRunsOutcome.GetResult().GetNextToken();
 } else {
 std::cerr << "Error listing jobs. "
 << listRunsOutcome.GetError().GetMessage()
 << std::endl;

 result = 1;
 break;
 }
 } while (!nextToken.empty());

 std::cout << "Your account has " << jobs.size() << " jobs."
 << std::endl;
 for (size_t i = 0; i < jobs.size(); ++i) {
 std::cout << " " << i + 1 << ". " << jobs[i] << std::endl;
 }
}
Aws::ShutdownAPI(options); // Should only be called once.
return result;
}
```

- APIの詳細については、「API リファレンス [ListJobs](#)」の「」を参照してください。AWS SDK for C++

## Java

### SDK for Java 2.x

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
package com.example.glue;

import software.amazon.awssdk.regions.Region;
import software.amazon.awssdk.services.glue.GlueClient;
import software.amazon.awssdk.services.glue.model.ListJobsRequest;
import software.amazon.awssdk.services.glue.model.ListJobsResponse;
import java.util.List;

public class HelloGlue {
 public static void main(String[] args) {
 GlueClient glueClient = GlueClient.builder()
 .region(Region.US_EAST_1)
 .build();

 listJobs(glueClient);
 }

 public static void listJobs(GlueClient glueClient) {
 ListJobsRequest request = ListJobsRequest.builder()
 .maxResults(10)
 .build();
 ListJobsResponse response = glueClient.listJobs(request);
 List<String> jobList = response.jobNames();
 jobList.forEach(job -> {
 System.out.println("Job Name: " + job);
 });
 }
}
```

```
}
}
```

- APIの詳細については、「API リファレンス[ListJobs](#)」の「」を参照してください。AWS SDK for Java 2.x

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
import { ListJobsCommand, GlueClient } from "@aws-sdk/client-glue";

const client = new GlueClient({});

export const main = async () => {
 const command = new ListJobsCommand({});

 const { JobNames } = await client.send(command);
 const formattedJobNames = JobNames.join("\n");
 console.log("Job names: ");
 console.log(formattedJobNames);
 return JobNames;
};
```

- APIの詳細については、「API リファレンス[ListJobs](#)」の「」を参照してください。AWS SDK for JavaScript

## Rust

### SDK for Rust

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
let mut list_jobs = glue.list_jobs().into_paginator().send();
while let Some(list_jobs_output) = list_jobs.next().await {
 match list_jobs_output {
 Ok(list_jobs) => {
 let names = list_jobs.job_names();
 info!(?names, "Found these jobs")
 }
 Err(err) => return Err(GlueMvpError::from_glue_sdk(err)),
 }
}
```

- API の詳細については、[ListJobsAWS](#) 「SDK for Rust API リファレンス」の「」を参照してください。

### コードの例

- [AWS SDKs AWS Glue を使用するためのアクション](#)
  - [AWS SDK または CLI CreateCrawlerで を使用する](#)
  - [AWS SDK または CLI CreateJobで を使用する](#)
  - [AWS SDK または CLI DeleteCrawlerで を使用する](#)
  - [AWS SDK または CLI DeleteDatabaseで を使用する](#)
  - [AWS SDK または CLI DeleteJobで を使用する](#)
  - [AWS SDK または CLI DeleteTableで を使用する](#)
  - [AWS SDK または CLI GetCrawlerで を使用する](#)
  - [AWS SDK または CLI GetDatabaseで を使用する](#)
  - [AWS SDK または CLI GetDatabasesで を使用する](#)

- [AWS SDK または CLI GetJobで を使用する](#)
- [AWS SDK または CLI GetJobRunで を使用する](#)
- [AWS SDK または CLI GetJobRunsで を使用する](#)
- [AWS SDK または CLI GetTablesで を使用する](#)
- [AWS SDK または CLI ListJobsで を使用する](#)
- [AWS SDK または CLI StartCrawlerで を使用する](#)
- [AWS SDK または CLI StartJobRunで を使用する](#)
- [AWS SDKs AWS Glue を使用するシナリオ](#)
  - [AWS SDK を使用して AWS Glue クローラーとジョブの実行を開始する](#)

## AWS SDKs AWS Glue を使用するためのアクション

次のコード例は、AWS SDKsで個々のAWS Glue アクションを実行する方法を示しています。これらの抜粋はAWS Glue API を呼び出し、コンテキスト内で実行する必要がある大規模なプログラムからのコードの抜粋です。各例にはGitHub、コードの設定と実行の手順を示すへのリンクが含まれています。

以下の例には、最も一般的に使用されるアクションのみ含まれています。詳細な一覧については、「[AWS Glue API リファレンス](#)」を参照してください。

### 例

- [AWS SDK または CLI CreateCrawlerで を使用する](#)
- [AWS SDK または CLI CreateJobで を使用する](#)
- [AWS SDK または CLI DeleteCrawlerで を使用する](#)
- [AWS SDK または CLI DeleteDatabaseで を使用する](#)
- [AWS SDK または CLI DeleteJobで を使用する](#)
- [AWS SDK または CLI DeleteTableで を使用する](#)
- [AWS SDK または CLI GetCrawlerで を使用する](#)
- [AWS SDK または CLI GetDatabaseで を使用する](#)
- [AWS SDK または CLI GetDatabasesで を使用する](#)
- [AWS SDK または CLI GetJobで を使用する](#)
- [AWS SDK または CLI GetJobRunで を使用する](#)
- [AWS SDK または CLI GetJobRunsで を使用する](#)

- [AWS SDK または CLI GetTables で を使用する](#)
- [AWS SDK または CLI ListJobs で を使用する](#)
- [AWS SDK または CLI StartCrawler で を使用する](#)
- [AWS SDK または CLI StartJobRun で を使用する](#)

## AWS SDK または CLI **CreateCrawler** で を使用する

以下のコード例は、CreateCrawler の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Create an AWS Glue crawler.
/// </summary>
/// <param name="crawlerName">The name for the crawler.</param>
/// <param name="crawlerDescription">A description of the crawler.</param>
/// <param name="role">The AWS Identity and Access Management (IAM) role to
/// be assumed by the crawler.</param>
/// <param name="schedule">The schedule on which the crawler will be
executed.</param>
/// <param name="s3Path">The path to the Amazon Simple Storage Service
(Amazon S3)
/// bucket where the Python script has been stored.</param>
/// <param name="dbName">The name to use for the database that will be
/// created by the crawler.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
```

```
public async Task<bool> CreateCrawlerAsync(
 string crawlerName,
 string crawlerDescription,
 string role,
 string schedule,
 string s3Path,
 string dbName)
{
 var s3Target = new S3Target
 {
 Path = s3Path,
 };

 var targetList = new List<S3Target>
 {
 s3Target,
 };

 var targets = new CrawlerTargets
 {
 S3Targets = targetList,
 };

 var crawlerRequest = new CreateCrawlerRequest
 {
 DatabaseName = dbName,
 Name = crawlerName,
 Description = crawlerDescription,
 Targets = targets,
 Role = role,
 Schedule = schedule,
 };

 var response = await _amazonGlue.CreateCrawlerAsync(crawlerRequest);
 return response.HttpStatusCode == System.Net.HttpStatusCode.OK;
}
```

- APIの詳細については、「APIリファレンス[CreateCrawler](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

## SDK for C++

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::S3Target s3Target;
s3Target.SetPath("s3://crawler-public-us-east-1/flight/2016/csv");
Aws::Glue::Model::CrawlerTargets crawlerTargets;
crawlerTargets.AddS3Targets(s3Target);

Aws::Glue::Model::CreateCrawlerRequest request;
request.SetTargets(crawlerTargets);
request.SetName(CRAWLER_NAME);
request.SetDatabaseName(CRAWLER_DATABASE_NAME);
request.SetTablePrefix(CRAWLER_DATABASE_PREFIX);
request.SetRole(roleArn);

Aws::Glue::Model::CreateCrawlerOutcome outcome =
client.CreateCrawler(request);

if (outcome.IsSuccess()) {
 std::cout << "Successfully created the crawler." << std::endl;
}
else {
 std::cerr << "Error creating a crawler. " <<
outcome.GetError().GetMessage()
 << std::endl;
 deleteAssets("", CRAWLER_DATABASE_NAME, "", bucketName,
clientConfig);
 return false;
}
```

```
}
```

- APIの詳細については、「APIリファレンス[CreateCrawler](#)」の「」を参照してください。  
AWS SDK for C++

## Java

### SDK for Java 2.x

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
import software.amazon.awssdk.regions.Region;
import software.amazon.awssdk.services.glue.GlueClient;
import software.amazon.awssdk.services.glue.model.CreateCrawlerRequest;
import software.amazon.awssdk.services.glue.model.CrawlerTargets;
import software.amazon.awssdk.services.glue.model.GlueException;
import software.amazon.awssdk.services.glue.model.S3Target;
import java.util.ArrayList;
import java.util.List;

/**
 * Before running this Java V2 code example, set up your development
 * environment, including your credentials.
 *
 * For more information, see the following documentation topic:
 *
 * https://docs.aws.amazon.com/sdk-for-java/latest/developer-guide/get-
 * started.html
 */
public class CreateCrawler {
 public static void main(String[] args) {
 final String usage = ""

 Usage:
 <IAM> <s3Path> <cron> <dbName> <crawlerName>
```

```

 Where:
 IAM - The ARN of the IAM role that has AWS Glue and S3
permissions.\s
 s3Path - The Amazon Simple Storage Service (Amazon S3) target
that contains data (for example, CSV data).
 cron - A cron expression used to specify the schedule (i.e.,
cron(15 12 * * ? *).
 dbName - The database name.\s
 crawlerName - The name of the crawler.\s
 """;

 if (args.length != 5) {
 System.out.println(usage);
 System.exit(1);
 }

 String iam = args[0];
 String s3Path = args[1];
 String cron = args[2];
 String dbName = args[3];
 String crawlerName = args[4];
 Region region = Region.US_EAST_1;
 GlueClient glueClient = GlueClient.builder()
 .region(region)
 .build();

 createGlueCrawler(glueClient, iam, s3Path, cron, dbName, crawlerName);
 glueClient.close();
}

public static void createGlueCrawler(GlueClient glueClient,
 String iam,
 String s3Path,
 String cron,
 String dbName,
 String crawlerName) {

 try {
 S3Target s3Target = S3Target.builder()
 .path(s3Path)
 .build();

 // Add the S3Target to a list.
 List<S3Target> targetList = new ArrayList<>();

```

```
targetList.add(s3Target);

CrawlerTargets targets = CrawlerTargets.builder()
 .s3Targets(targetList)
 .build();

CreateCrawlerRequest crawlerRequest = CreateCrawlerRequest.builder()
 .databaseName(dbName)
 .name(crawlerName)
 .description("Created by the AWS Glue Java API")
 .targets(targets)
 .role(iam)
 .schedule(cron)
 .build();

glueClient.createCrawler(crawlerRequest);
System.out.println(crawlerName + " was successfully created");

} catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
}
}
```

- APIの詳細については、「APIリファレンス[CreateCrawler](#)」の「」を参照してください。  
AWS SDK for Java 2.x

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const createCrawler = (name, role, dbName, tablePrefix, s3TargetPath) => {
 const client = new GlueClient({});
```

```
const command = new CreateCrawlerCommand({
 Name: name,
 Role: role,
 DatabaseName: dbName,
 TablePrefix: tablePrefix,
 Targets: {
 S3Targets: [{ Path: s3TargetPath }],
 },
});

return client.send(command);
};
```

- APIの詳細については、「APIリファレンス[CreateCrawler](#)」の「」を参照してください。  
AWS SDK for JavaScript

## Kotlin

### SDK for Kotlin

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
suspend fun createGlueCrawler(
 iam: String?,
 s3Path: String?,
 cron: String?,
 dbName: String?,
 crawlerName: String,
) {
 val s3Target =
 S3Target {
 path = s3Path
 }

 // Add the S3Target to a list.
```

```
val targetList = mutableListOf<S3Target>()
targetList.add(s3Target)

val target0b =
 CrawlerTargets {
 s3Targets = targetList
 }

val request =
 CreateCrawlerRequest {
 databaseName = dbName
 name = crawlerName
 description = "Created by the AWS Glue Kotlin API"
 targets = target0b
 role = iam
 schedule = cron
 }

GlueClient { region = "us-west-2" }.use { glueClient ->
 glueClient.createCrawler(request)
 println("$crawlerName was successfully created")
}
}
```

- APIの詳細については、AWS SDK for Kotlin API リファレンス [CreateCrawler](#) の「」を参照してください。

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
$crawlerName = "example-crawler-test-" . $uniqid;

$role = $iamService->getRole("AWSGlueServiceRole-DocExample");
```

```

 $path = 's3://crawler-public-us-east-1/flight/2016/csv';
 $glueService->createCrawler($crawlerName, $role['Role']['Arn'],
 $databaseName, $path);

 public function createCrawler($crawlerName, $role, $databaseName, $path):
 Result
 {
 return $this->customWaiter(function () use ($crawlerName, $role,
 $databaseName, $path) {
 return $this->glueClient->createCrawler([
 'Name' => $crawlerName,
 'Role' => $role,
 'DatabaseName' => $databaseName,
 'Targets' => [
 'S3Targets' =>
 [[
 'Path' => $path,
]]
],
]);
 });
 }

```

- APIの詳細については、「APIリファレンス[CreateCrawler](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください GitHub。 [AWSコード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```

class GlueWrapper:
 """Encapsulates AWS Glue actions."""

```

```
def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

def create_crawler(self, name, role_arn, db_name, db_prefix, s3_target):
 """
 Creates a crawler that can crawl the specified target and populate a
 database in your AWS Glue Data Catalog with metadata that describes the
 data
 in the target.

 :param name: The name of the crawler.
 :param role_arn: The Amazon Resource Name (ARN) of an AWS Identity and
 Access
 Management (IAM) role that grants permission to let AWS
 Glue
 access the resources it needs.
 :param db_name: The name to give the database that is created by the
 crawler.
 :param db_prefix: The prefix to give any database tables that are created
 by
 the crawler.
 :param s3_target: The URL to an S3 bucket that contains data that is
 the target of the crawler.
 """
 try:
 self.glue_client.create_crawler(
 Name=name,
 Role=role_arn,
 DatabaseName=db_name,
 TablePrefix=db_prefix,
 Targets={"S3Targets": [{"Path": s3_target}]},
)
 except ClientError as err:
 logger.error(
 "Couldn't create crawler. Here's why: %s: %s",
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
```

- APIの詳細については、[CreateCrawler](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Creates a new crawler with the specified configuration.
 #
 # @param name [String] The name of the crawler.
 # @param role_arn [String] The ARN of the IAM role to be used by the crawler.
 # @param db_name [String] The name of the database where the crawler stores its
 metadata.
 # @param db_prefix [String] The prefix to be added to the names of tables that
 the crawler creates.
 # @param s3_target [String] The S3 path that the crawler will crawl.
 # @return [void]
 def create_crawler(name, role_arn, db_name, db_prefix, s3_target)
 @glue_client.create_crawler(
 name: name,
```

```

 role: role_arn,
 database_name: db_name,
 targets: {
 s3_targets: [
 {
 path: s3_target
 }
]
 }
)
rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not create crawler: \n#{e.message}")
 raise
end

```

- APIの詳細については、「APIリファレンス [CreateCrawler](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```

let create_crawler = glue
 .create_crawler()
 .name(self.crawler())
 .database_name(self.database())
 .role(self.iam_role.expose_secret())
 .targets(
 CrawlerTargets::builder()
 .s3_targets(S3Target::builder().path(CRAWLER_TARGET).build())
 .build(),
)
 .send()
 .await;

```

```
match create_crawler {
 Err(err) => {
 let glue_err: aws_sdk_glue::Error = err.into();
 match glue_err {
 aws_sdk_glue::Error::AlreadyExistsException(_) => {
 info!("Using existing crawler");
 Ok(())
 }
 _ => Err(GlueMvpError::GlueSdk(glue_err)),
 }
 }
 Ok(_) => Ok(()),
}??;
```

- APIの詳細については、[CreateCrawler](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI **CreateJob**で を使用する

以下のコード例は、CreateJob の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Create an AWS Glue job.
/// </summary>
/// <param name="jobName">The name of the job.</param>
/// <param name="roleName">The name of the IAM role to be assumed by
/// the job.</param>
/// <param name="description">A description of the job.</param>
/// <param name="scriptUrl">The URL to the script.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> CreateJobAsync(string dbName, string tableName,
string bucketUrl, string jobName, string roleName, string description, string
scriptUrl)
{
 var command = new JobCommand
 {
 PythonVersion = "3",
 Name = "glueetl",
 ScriptLocation = scriptUrl,
 };

 var arguments = new Dictionary<string, string>
 {
 { "--input_database", dbName },
 { "--input_table", tableName },
 { "--output_bucket_url", bucketUrl }
 };

 var request = new CreateJobRequest
 {
 Command = command,
 DefaultArguments = arguments,
 Description = description,
 GlueVersion = "3.0",
 Name = jobName,
 NumberOfWorkers = 10,
 Role = roleName,
 WorkerType = "G.1X"
 };

 var response = await _amazonGlue.CreateJobAsync(request);
 return response.HttpStatusCode == HttpStatusCode.OK;
}
```

- APIの詳細については、「API リファレンス [CreateJob](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

### SDK for C++

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::CreateJobRequest request;
request.SetName(JOB_NAME);
request.SetRole(roleArn);
request.SetGlueVersion(GLUE_VERSION);

Aws::Glue::Model::JobCommand command;
command.SetName(JOB_COMMAND_NAME);
command.SetPythonVersion(JOB_PYTHON_VERSION);
command.SetScriptLocation(
 Aws::String("s3://") + bucketName + "/" + PYTHON_SCRIPT);
request.SetCommand(command);

Aws::Glue::Model::CreateJobOutcome outcome = client.CreateJob(request);

if (outcome.IsSuccess()) {
 std::cout << "Successfully created the job." << std::endl;
}
else {
```

```
 std::cerr << "Error creating the job. " <<
outcome.GetError().GetMessage()
 << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, "", bucketName,
 clientConfig);
 return false;
 }
```

- APIの詳細については、「API リファレンス [CreateJob](#)」の「」を参照してください。  
AWS SDK for C++

## CLI

### AWS CLI

データを変換するジョブを作成するには

次の `create-job` 例では、S3 に保存されているスクリプトを実行するストリーミングジョブを作成します。

```
aws glue create-job \
 --name my-testing-job \
 --role AWSGlueServiceRoleDefault \
 --command '{ \
 "Name": "gluestreaming", \
 "ScriptLocation": "s3://DOC-EXAMPLE-BUCKET/folder/" \
 }' \
 --region us-east-1 \
 --output json \
 --default-arguments '{ \
 "--job-language":"scala", \
 "--class":"GlueApp" \
 }' \
 --profile my-profile \
 --endpoint https://glue.us-east-1.amazonaws.com
```

`test_script.scala` の内容:

```
import com.amazonaws.services.glue.ChoiceOption
import com.amazonaws.services.glue.GlueContext
import com.amazonaws.services.glue.MappingSpec
```

```
import com.amazonaws.services.glue.ResolveSpec
import com.amazonaws.services.glue.errors.CallSite
import com.amazonaws.services.glue.util.GlueArgParser
import com.amazonaws.services.glue.util.Job
import com.amazonaws.services.glue.util.JsonOptions
import org.apache.spark.SparkContext
import scala.collection.JavaConverters._

object GlueApp {
 def main(sysArgs: Array[String]) {
 val spark: SparkContext = new SparkContext()
 val glueContext: GlueContext = new GlueContext(spark)
 // @params: [JOB_NAME]
 val args = GlueArgParser.getResolvedOptions(sysArgs,
Seq("JOB_NAME").toArray)
 Job.init(args("JOB_NAME"), glueContext, args.asJava)
 // @type: DataSource
 // @args: [database = "tempdb", table_name = "s3-source",
transformation_ctx = "datasource0"]
 // @return: datasource0
 // @inputs: []
 val datasource0 = glueContext.getCatalogSource(database = "tempdb",
tableName = "s3-source", redshiftTmpDir = "", transformationContext =
"datasource0").getDynamicFrame()
 // @type: ApplyMapping
 // @args: [mapping = [("sensorid", "int", "sensorid", "int"),
("currenttemperature", "int", "currenttemperature", "int"), ("status", "string",
"status", "string")], transformation_ctx = "applymapping1"]
 // @return: applymapping1
 // @inputs: [frame = datasource0]
 val applymapping1 = datasource0.applyMapping(mappings = Seq(("sensorid",
"int", "sensorid", "int"), ("currenttemperature", "int", "currenttemperature",
"int"), ("status", "string", "status", "string")), caseSensitive = false,
transformationContext = "applymapping1")
 // @type: SelectFields
 // @args: [paths = ["sensorid", "currenttemperature", "status"],
transformation_ctx = "selectfields2"]
 // @return: selectfields2
 // @inputs: [frame = applymapping1]
 val selectfields2 = applymapping1.selectFields(paths = Seq("sensorid",
"currenttemperature", "status"), transformationContext = "selectfields2")
 // @type: ResolveChoice
 // @args: [choice = "MATCH_CATALOG", database = "tempdb", table_name =
"my-s3-sink", transformation_ctx = "resolvechoice3"]
```

```
 // @return: resolvechoice3
 // @inputs: [frame = selectfields2]
 val resolvechoice3 = selectfields2.resolveChoice(choiceOption =
Some(ChoiceOption("MATCH_CATALOG")), database = Some("tempdb"), tableName =
Some("my-s3-sink"), transformationContext = "resolvechoice3")
 // @type: DataSink
 // @args: [database = "tempdb", table_name = "my-s3-sink",
transformation_ctx = "datasink4"]
 // @return: datasink4
 // @inputs: [frame = resolvechoice3]
 val datasink4 = glueContext.getCatalogSink(database = "tempdb",
tableName = "my-s3-sink", redshiftTmpDir = "", transformationContext =
"datasink4").writeDynamicFrame(resolvechoice3)
 Job.commit()
 }
}
```

出力:

```
{
 "Name": "my-testing-job"
}
```

詳細については、「[AWS Glue デベロッパーガイド](#)」の「[Glue でのジョブの作成](#)」を参照してください。AWS

- API の詳細については、「[コマンドリファレンス `CreateJob`](#)」の「」を参照してください。  
AWS CLI

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const createJob = (name, role, scriptBucketName, scriptKey) => {
 const client = new GlueClient({});
```

```
const command = new CreateJobCommand({
 Name: name,
 Role: role,
 Command: {
 Name: "glueetl",
 PythonVersion: "3",
 ScriptLocation: `s3://${scriptBucketName}/${scriptKey}`,
 },
 GlueVersion: "3.0",
});

return client.send(command);
};
```

- APIの詳細については、「APIリファレンス[CreateJob](#)」の「」を参照してください。  
AWS SDK for JavaScript

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
$role = $iamService->getRole("AWSGlueServiceRole-DocExample");

$jobName = 'test-job-' . $uniqid;

$scriptLocation = "s3://$bucketName/run_job.py";
$job = $glueService->createJob($jobName, $role['Role']['Arn'],
$scriptLocation);

public function createJob($jobName, $role, $scriptLocation, $pythonVersion =
'3', $glueVersion = '3.0'): Result
{
 return $this->glueClient->createJob([
```

```

 'Name' => $jobName,
 'Role' => $role,
 'Command' => [
 'Name' => 'glueetl',
 'ScriptLocation' => $scriptLocation,
 'PythonVersion' => $pythonVersion,
],
 'GlueVersion' => $glueVersion,
]);
}

```

- APIの詳細については、「APIリファレンス[CreateJob](#)」の「」を参照してください。  
AWS SDK for PHP

## PowerShell

### のツール PowerShell

例 1: この例では、AWS Glue で新しいジョブを作成します。コマンド名の値は常に **glueetl**。AWS Glue は Python または Scala で記述されたジョブスクリプトの実行をサポートしています。この例では、ジョブスクリプト (MyTestGlueJob.py) は Python で記述されています。Python パラメータは **\$DefArgs** 変数で指定され、**DefaultArguments** パラメータの PowerShell コマンドに渡されます。このパラメータはハッシュ可能なを受け入れます。**\$JobParams** 変数のパラメータは、AWS Glue CreateJob API リファレンスの「Jobs (https://docs.aws.amazon.com/glue/latest/dg/aws-glue-api-jobs-job.html)」トピックに記載されている API から取得されます。

```

$Command = New-Object Amazon.Glue.Model.JobCommand
$Command.Name = 'glueetl'
$Command.ScriptLocation = 's3://aws-glue-scripts-000000000000-us-west-2/admin/MyTestGlueJob.py'
$Command

$Source = "source_test_table"
$Target = "target_test_table"
$Connections = $Source, $Target

$DefArgs = @{
 '--TempDir' = 's3://aws-glue-temporary-000000000000-us-west-2/admin'
 '--job-bookmark-option' = 'job-bookmark-disable'
}

```

```

 '--job-language' = 'python'
 }
$DefArgs

$ExecutionProp = New-Object Amazon.Glue.Model.ExecutionProperty
$ExecutionProp.MaxConcurrentRuns = 1
$ExecutionProp

$JobParams = @{
 "AllocatedCapacity" = "5"
 "Command" = $Command
 "Connections_Connection" = $Connections
 "DefaultArguments" = $DefArgs
 "Description" = "This is a test"
 "ExecutionProperty" = $ExecutionProp
 "MaxRetries" = "1"
 "Name" = "MyOregonTestGlueJob"
 "Role" = "Amazon-GlueServiceRoleForSSM"
 "Timeout" = "20"
}

New-GlueJob @JobParams

```

- APIの詳細については、「コマンドレットリファレンス[CreateJob](#)」の「」を参照してください。AWS Tools for PowerShell

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```

class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """

```

```
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

 def create_job(self, name, description, role_arn, script_location):
 """
 Creates a job definition for an extract, transform, and load (ETL) job
 that can
 be run by AWS Glue.

 :param name: The name of the job definition.
 :param description: The description of the job definition.
 :param role_arn: The ARN of an IAM role that grants AWS Glue the
 permissions
 it requires to run the job.
 :param script_location: The Amazon S3 URL of a Python ETL script that is
 run as
 part of the job. The script defines how the data
 is
 transformed.
 """
 try:
 self.glue_client.create_job(
 Name=name,
 Description=description,
 Role=role_arn,
 Command={
 "Name": "glueetl",
 "ScriptLocation": script_location,
 "PythonVersion": "3",
 },
 GlueVersion="3.0",
)
 except ClientError as err:
 logger.error(
 "Couldn't create job %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
```

- APIの詳細については、[CreateJob](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Creates a new job with the specified configuration.
 #
 # @param name [String] The name of the job.
 # @param description [String] The description of the job.
 # @param role_arn [String] The ARN of the IAM role to be used by the job.
 # @param script_location [String] The location of the ETL script for the job.
 # @return [void]
 def create_job(name, description, role_arn, script_location)
 @glue_client.create_job(
 name: name,
 description: description,
 role: role_arn,
 command: {
```

```
 name: "glueetl",
 script_location: script_location,
 python_version: "3"
 },
 glue_version: "3.0"
)
rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not create job #{name}: \n#{e.message}")
 raise
end
```

- APIの詳細については、「APIリファレンス[CreateJob](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください [GitHub](#)。 [AWSコード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
let create_job = glue
 .create_job()
 .name(self.job())
 .role(self.iam_role.expose_secret())
 .command(
 JobCommand::builder()
 .name("glueetl")
 .python_version("3")
 .script_location(format!("s3://{}/job.py", self.bucket()))
 .build(),
)
 .glue_version("3.0")
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;
```

```
let job_name = create_job.name().ok_or_else(|| {
 GlueMvpError::Unknown("Did not get job name after creating
job".into())
})?;
```

- APIの詳細については、[CreateJob](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください。[AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI **DeleteCrawler**で を使用する

以下のコード例は、DeleteCrawler の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Delete an AWS Glue crawler.
/// </summary>
/// <param name="crawlerName">The name of the crawler.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> DeleteCrawlerAsync(string crawlerName)
{
```

```
var response = await _amazonGlue.DeleteCrawlerAsync(new
DeleteCrawlerRequest { Name = crawlerName });
return response.HttpStatusCode == HttpStatusCode.OK;
}
```

- APIの詳細については、「APIリファレンス[DeleteCrawler](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

### SDK for C++

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::DeleteCrawlerRequest request;
request.SetName(crawler);

Aws::Glue::Model::DeleteCrawlerOutcome outcome =
client.DeleteCrawler(request);

if (outcome.IsSuccess()) {
 std::cout << "Successfully deleted the crawler." << std::endl;
}
else {
 std::cerr << "Error deleting the crawler. "
 << outcome.GetError().GetMessage() << std::endl;
 result = false;
}
```

- APIの詳細については、「API リファレンス [DeleteCrawler](#)」の「」を参照してください。  
AWS SDK for C++

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const deleteCrawler = (crawlerName) => {
 const client = new GlueClient({});

 const command = new DeleteCrawlerCommand({
 Name: crawlerName,
 });

 return client.send(command);
};
```

- APIの詳細については、「API リファレンス [DeleteCrawler](#)」の「」を参照してください。  
AWS SDK for JavaScript

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
echo "Delete the crawler.\n";
$glueClient->deleteCrawler([
 'Name' => $crawlerName,
]);

public function deleteCrawler($crawlerName)
{
 return $this->glueClient->deleteCrawler([
 'Name' => $crawlerName,
]);
}
```

- APIの詳細については、「APIリファレンス[DeleteCrawler](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

 def delete_crawler(self, name):
 """
 Deletes a crawler.

 :param name: The name of the crawler to delete.
```

```
"""
try:
 self.glue_client.delete_crawler(Name=name)
except ClientError as err:
 logger.error(
 "Couldn't delete crawler %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
```

- APIの詳細については、[DeleteCrawler](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Deletes a crawler with the specified name.
```

```
#
@param name [String] The name of the crawler to delete.
@return [void]
def delete_crawler(name)
 @glue_client.delete_crawler(name: name)
rescue Aws::Glue::Errors::ServiceError => e
 @logger.error("Glue could not delete crawler #{name}: \n#{e.message}")
 raise
end
```

- APIの詳細については、「APIリファレンス[DeleteCrawler](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください [GitHub](#)。 [AWSコード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
glue.delete_crawler()
 .name(self.crawler())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;
```

- APIの詳細については、[DeleteCrawler](#) AWS 「SDK for Rust APIリファレンス」の「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDKでこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前のSDKバージョンの詳細も含まれています。

## AWS SDK または CLI `DeleteDatabase` で使用する

以下のコード例は、`DeleteDatabase` の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Delete the AWS Glue database.
/// </summary>
/// <param name="dbName">The name of the database.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> DeleteDatabaseAsync(string dbName)
{
 var response = await _amazonGlue.DeleteDatabaseAsync(new
DeleteDatabaseRequest { Name = dbName });
 return response.HttpStatusCode == HttpStatusCode.OK;
}
```

- API の詳細については、「API リファレンス [DeleteDatabase](#)」の「」を参照してください。 AWS SDK for .NET

## C++

## SDK for C++

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::DeleteDatabaseRequest request;
request.SetName(database);

Aws::Glue::Model::DeleteDatabaseOutcome outcome = client.DeleteDatabase(
 request);

if (outcome.IsSuccess()) {
 std::cout << "Successfully deleted the database." << std::endl;
}
else {
 std::cerr << "Error deleting database. " <<
outcome.GetError().GetMessage()
 << std::endl;
 result = false;
}
```

- APIの詳細については、「API リファレンス [DeleteDatabase](#)」の「」を参照してください。AWS SDK for C++

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const deleteDatabase = (databaseName) => {
 const client = new GlueClient({});

 const command = new DeleteDatabaseCommand({
 Name: databaseName,
 });

 return client.send(command);
};
```

- API の詳細については、「API リファレンス [DeleteDatabase](#)」の「」を参照してください。 AWS SDK for JavaScript

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
echo "Delete the databases.\n";
$glueClient->deleteDatabase([
 'Name' => $databaseName,
]);

public function deleteDatabase($databaseName)
```

```
{
 return $this->glueClient->deleteDatabase([
 'Name' => $databaseName,
]);
}
```

- APIの詳細については、「APIリファレンス[DeleteDatabase](#)」の「」を参照してください。AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

 def delete_database(self, name):
 """
 Deletes a metadata database from your Data Catalog.

 :param name: The name of the database to delete.
 """
 try:
 self.glue_client.delete_database(Name=name)
 except ClientError as err:
 logger.error(
 "Couldn't delete database %s. Here's why: %s: %s",
```

```
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
```

- APIの詳細については、[DeleteDatabase](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Removes a specified database from a Data Catalog.
 #
 # @param database_name [String] The name of the database to delete.
 # @return [void]
 def delete_database(database_name)
 @glue_client.delete_database(name: database_name)
 rescue Aws::Glue::Errors::ServiceError => e
```

```
@logger.error("Glue could not delete database: \n#{e.message}")
end
```

- APIの詳細については、「API リファレンス [DeleteDatabase](#)」の「」を参照してください。AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
glue.delete_database()
 .name(self.database())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;
```

- APIの詳細については、 [DeleteDatabase](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI **DeleteJob**で を使用する

以下のコード例は、DeleteJob の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

## .NET

### AWS SDK for .NET

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Delete an AWS Glue job.
/// </summary>
/// <param name="jobName">The name of the job.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> DeleteJobAsync(string jobName)
{
 var response = await _amazonGlue.DeleteJobAsync(new DeleteJobRequest
{ JobName = jobName });
 return response.HttpStatusCode == HttpStatusCode.OK;
}
```

- API の詳細については、「API リファレンス [DeleteJob](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

### SDK for C++

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
```

```
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::DeleteJobRequest request;
request.SetJobName(job);

Aws::Glue::Model::DeleteJobOutcome outcome = client.DeleteJob(request);

if (outcome.IsSuccess()) {
 std::cout << "Successfully deleted the job." << std::endl;
}
else {
 std::cerr << "Error deleting the job. " <<
outcome.GetError().GetMessage()
 << std::endl;
 result = false;
}
```

- APIの詳細については、「APIリファレンス[DeleteJob](#)」の「」を参照してください。  
AWS SDK for C++

## CLI

### AWS CLI

ジョブを削除するには

次の `delete-job` 例では、不要になったジョブを削除します。

```
aws glue delete-job \
 --job-name my-testing-job
```

出力:

```
{
 "JobName": "my-testing-job"
}
```

詳細については、「[AWS Glue デベロッパーガイド](#)」の「[Glue コンソールでのジョブの使用](#)」を参照してください。AWS

- API の詳細については、「[コマンドリファレンスDeleteJob](#)」の「」を参照してください。  
AWS CLI

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const deleteJob = (jobName) => {
 const client = new GlueClient({});

 const command = new DeleteJobCommand({
 JobName: jobName,
 });

 return client.send(command);
};
```

- API の詳細については、「[API リファレンスDeleteJob](#)」の「」を参照してください。  
AWS SDK for JavaScript

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
echo "Delete the job.\n";
$glueClient->deleteJob([
 'JobName' => $job['Name'],
]);

public function deleteJob($jobName)
{
 return $this->glueClient->deleteJob([
 'JobName' => $jobName,
]);
}
```

- APIの詳細については、「APIリファレンス[DeleteJob](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

 def delete_job(self, job_name):
 """
 Deletes a job definition. This also deletes data about all runs that are
 associated with this job definition.
 """
```

```
:param job_name: The name of the job definition to delete.
"""
try:
 self.glue_client.delete_job(JobName=job_name)
except ClientError as err:
 logger.error(
 "Couldn't delete job %s. Here's why: %s: %s",
 job_name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
```

- APIの詳細については、[DeleteJob](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end
```

```
Deletes a job with the specified name.
#
@param job_name [String] The name of the job to delete.
@return [void]
def delete_job(job_name)
 @glue_client.delete_job(job_name: job_name)
rescue Aws::Glue::Errors::ServiceError => e
 @logger.error("Glue could not delete job: \n#{e.message}")
end
```

- APIの詳細については、「APIリファレンス[DeleteJob](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
glue.delete_job()
 .job_name(self.job())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;
```

- APIの詳細については、 [DeleteJob](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI `DeleteTable`で を使用する

以下のコード例は、`DeleteTable` の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Delete a table from an AWS Glue database.
/// </summary>
/// <param name="tableName">The table to delete.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> DeleteTableAsync(string dbName, string tableName)
{
 var response = await _amazonGlue.DeleteTableAsync(new DeleteTableRequest
{ Name = tableName, DatabaseName = dbName });
 return response.HttpStatusCode == HttpStatusCode.OK;
}
```

- API の詳細については、「API リファレンス [DeleteTable](#)」の「」を参照してください。  
AWS SDK for .NET

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const deleteTable = (databaseName, tableName) => {
 const client = new GlueClient({});

 const command = new DeleteTableCommand({
 DatabaseName: databaseName,
 Name: tableName,
 });

 return client.send(command);
};
```

- API の詳細については、「API リファレンス [DeleteTable](#)」の「」を参照してください。  
AWS SDK for JavaScript

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
echo "Delete the tables.\n";
foreach ($tables['TableList'] as $table) {
 $glueService->deleteTable($table['Name'], $databaseName);
}
```

```
public function deleteTable($tableName, $databaseName)
{
 return $this->glueClient->deleteTable([
 'DatabaseName' => $databaseName,
 'Name' => $tableName,
]);
}
```

- APIの詳細については、「APIリファレンス[DeleteTable](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください GitHub。 [AWSコード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

 def delete_table(self, db_name, table_name):
 """
 Deletes a table from a metadata database.

 :param db_name: The name of the database that contains the table.
 :param table_name: The name of the table to delete.
 """
 try:
 self.glue_client.delete_table(DatabaseName=db_name, Name=table_name)
```

```
except ClientError as err:
 logger.error(
 "Couldn't delete table %s. Here's why: %s: %s",
 table_name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
```

- API の詳細については、[DeleteTable](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Deletes a table with the specified name.
 #
 # @param database_name [String] The name of the catalog database in which the
table resides.
```

```
@param table_name [String] The name of the table to be deleted.
@return [void]
def delete_table(database_name, table_name)
 @glue_client.delete_table(database_name: database_name, name: table_name)
rescue Aws::Glue::Errors::ServiceError => e
 @logger.error("Glue could not delete job: \n#{e.message}")
end
```

- APIの詳細については、「APIリファレンス[DeleteTable](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください GitHub。 [AWSコード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
for t in &self.tables {
 glue.delete_table()
 .name(t.name())
 .database_name(self.database())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;
}
```

- APIの詳細については、[DeleteTable](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI **GetCrawler**で を使用する

以下のコード例は、GetCrawler の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Get information about an AWS Glue crawler.
/// </summary>
/// <param name="crawlerName">The name of the crawler.</param>
/// <returns>A Crawler object describing the crawler.</returns>
public async Task<Crawler?> GetCrawlerAsync(string crawlerName)
{
 var crawlerRequest = new GetCrawlerRequest
 {
 Name = crawlerName,
 };

 var response = await _amazonGlue.GetCrawlerAsync(crawlerRequest);
 if (response.HttpStatusCode == System.Net.HttpStatusCode.OK)
 {
 var databaseName = response.Crawler.DatabaseName;
 Console.WriteLine($"{crawlerName} has the database {databaseName}");
 return response.Crawler;
 }

 Console.WriteLine($"No information regarding {crawlerName} could be
found.");
}
```

```
 return null;
}
```

- APIの詳細については、「APIリファレンス[GetCrawler](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

### SDK for C++

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::GetCrawlerRequest request;
request.SetName(CRAWLER_NAME);

Aws::Glue::Model::GetCrawlerOutcome outcome = client.GetCrawler(request);

if (outcome.IsSuccess()) {
 Aws::Glue::Model::CrawlerState crawlerState =
outcome.GetResult().GetCrawler().GetState();
 std::cout << "Retrieved crawler with state " <<

Aws::Glue::Model::CrawlerStateMapper::GetNameForCrawlerState(
 crawlerState)
 << "." << std::endl;
}
else {
 std::cerr << "Error retrieving a crawler. "
```

```
 << outcome.GetError().GetMessage() << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, "", bucketName,
 clientConfig);
 return false;
 }
```

- APIの詳細については、「APIリファレンス[GetCrawler](#)」の「」を参照してください。  
AWS SDK for C++

## Java

### SDK for Java 2.x

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
import software.amazon.awssdk.regions.Region;
import software.amazon.awssdk.services.glue.GlueClient;
import software.amazon.awssdk.services.glue.model.GetCrawlerRequest;
import software.amazon.awssdk.services.glue.model.GetCrawlerResponse;
import software.amazon.awssdk.services.glue.model.GlueException;
import java.time.Instant;
import java.time.ZoneId;
import java.time.format.DateTimeFormatter;
import java.time.format.FormatStyle;
import java.util.Locale;

/**
 * Before running this Java V2 code example, set up your development
 * environment, including your credentials.
 *
 * For more information, see the following documentation topic:
 *
 * https://docs.aws.amazon.com/sdk-for-java/latest/developer-guide/get-
 * started.html
 */
public class GetCrawler {
```

```
public static void main(String[] args) {
 final String usage = ""

 Usage:
 <crawlerName>

 Where:
 crawlerName - The name of the crawler.\s
 """;

 if (args.length != 1) {
 System.out.println(usage);
 System.exit(1);
 }

 String crawlerName = args[0];
 Region region = Region.US_EAST_1;
 GlueClient glueClient = GlueClient.builder()
 .region(region)
 .build();

 getSpecificCrawler(glueClient, crawlerName);
 glueClient.close();
}

public static void getSpecificCrawler(GlueClient glueClient, String
crawlerName) {
 try {
 GetCrawlerRequest crawlerRequest = GetCrawlerRequest.builder()
 .name(crawlerName)
 .build();

 GetCrawlerResponse response = glueClient.getCrawler(crawlerRequest);
 Instant createDate = response.crawler().creationTime();

 // Convert the Instant to readable date
 DateTimeFormatter formatter =
 DateTimeFormatter.ofLocalizedDateTime(FormatStyle.SHORT)
 .withLocale(Locale.US)
 .withZone(ZoneId.systemDefault());

 formatter.format(createDate);
 System.out.println("The create date of the Crawler is " +
createDate);
 }
}
```

```
 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
}
}
```

- APIの詳細については、「APIリファレンス[GetCrawler](#)」の「」を参照してください。  
AWS SDK for Java 2.x

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const getCrawler = (name) => {
 const client = new GlueClient({});

 const command = new GetCrawlerCommand({
 Name: name,
 });

 return client.send(command);
};
```

- APIの詳細については、「APIリファレンス[GetCrawler](#)」の「」を参照してください。  
AWS SDK for JavaScript

## Kotlin

### SDK for Kotlin

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
suspend fun getSpecificCrawler(crawlerName: String?) {
 val request =
 GetCrawlerRequest {
 name = crawlerName
 }
 GlueClient { region = "us-east-1" }.use { glueClient ->
 val response = glueClient.getCrawler(request)
 val role = response.crawler?.role
 println("The role associated with this crawler is $role")
 }
}
```

- API の詳細については、AWS SDK for Kotlin API リファレンス [GetCrawler](#) の「」を参照してください。

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
echo "Waiting for crawler";
do {
 $crawler = $glueService->getCrawler($crawlerName);
```

```
 echo ".";
 sleep(10);
 } while ($crawler['Crawler']['State'] != "READY");
 echo "\n";

 public function getCrawler($crawlerName)
 {
 return $this->customWaiter(function () use ($crawlerName) {
 return $this->glueClient->getCrawler([
 'Name' => $crawlerName,
]);
 });
 }
}
```

- APIの詳細については、「APIリファレンス[GetCrawler](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

 def get_crawler(self, name):
 """
 Gets information about a crawler.
 """
```

```
:param name: The name of the crawler to look up.
:return: Data about the crawler.
"""
crawler = None
try:
 response = self.glue_client.get_crawler(Name=name)
 crawler = response["Crawler"]
except ClientError as err:
 if err.response["Error"]["Code"] == "EntityNotFoundException":
 logger.info("Crawler %s doesn't exist.", name)
 else:
 logger.error(
 "Couldn't get crawler %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
return crawler
```

- APIの詳細については、[GetCrawler](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
```

```
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Retrieves information about a specific crawler.
 #
 # @param name [String] The name of the crawler to retrieve information about.
 # @return [Aws::Glue::Types::Crawler, nil] The crawler object if found, or nil
 if not found.
 def get_crawler(name)
 @glue_client.get_crawler(name: name)
 rescue Aws::Glue::Errors::EntityNotFoundException
 @logger.info("Crawler #{name} doesn't exist.")
 false
 rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not get crawler #{name}: \n#{e.message}")
 raise
 end
end
```

- APIの詳細については、「APIリファレンス[GetCrawler](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
let tmp_crawler = glue
 .get_crawler()
 .name(self.crawler())
 .send()
```

```
.await
.map_err(GlueMvpError::from_glue_sdk)?;
```

- APIの詳細については、[GetCrawler](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI **GetDatabase**で を使用する

以下のコード例は、GetDatabase の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Get information about an AWS Glue database.
/// </summary>
/// <param name="dbName">The name of the database.</param>
/// <returns>A Database object containing information about the database.</
returns>
public async Task<Database> GetDatabaseAsync(string dbName)
{
 var databasesRequest = new GetDatabaseRequest
```

```
{
 Name = dbName,
};

var response = await _amazonGlue.GetDatabaseAsync(databasesRequest);
return response.Database;
}
```

- APIの詳細については、「APIリファレンス[GetDatabase](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

### SDK for C++

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::GetDatabaseRequest request;
request.SetName(CRAWLER_DATABASE_NAME);

Aws::Glue::Model::GetDatabaseOutcome outcome =
client.GetDatabase(request);

if (outcome.IsSuccess()) {
 const Aws::Glue::Model::Database &database =
outcome.GetResult().GetDatabase();

 std::cout << "Successfully retrieve the database\n" <<
```

```
 database.Jsonize().View().WriteReadable() << "." <<
std::endl;
 }
 else {
 std::cerr << "Error getting the database. "
 << outcome.GetError().GetMessage() << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, "", bucketName,
 clientConfig);
 return false;
 }
}
```

- APIの詳細については、「APIリファレンス[GetDatabase](#)」の「」を参照してください。  
AWS SDK for C++

## Java

### SDK for Java 2.x

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
import software.amazon.awssdk.regions.Region;
import software.amazon.awssdk.services.glue.GlueClient;
import software.amazon.awssdk.services.glue.model.GetDatabaseRequest;
import software.amazon.awssdk.services.glue.model.GetDatabaseResponse;
import software.amazon.awssdk.services.glue.model.GlueException;
import java.time.Instant;
import java.time.ZoneId;
import java.time.format.DateTimeFormatter;
import java.time.format.FormatStyle;
import java.util.Locale;

/**
 * Before running this Java V2 code example, set up your development
 * environment, including your credentials.
 *
 * For more information, see the following documentation topic:
```

```
*
* https://docs.aws.amazon.com/sdk-for-java/latest/developer-guide/get-
started.html
*/
public class GetDatabase {
 public static void main(String[] args) {
 final String usage = ""

 Usage:
 <databaseName>

 Where:
 databaseName - The name of the database.\s
 """;

 if (args.length != 1) {
 System.out.println(usage);
 System.exit(1);
 }

 String databaseName = args[0];
 Region region = Region.US_EAST_1;
 GlueClient glueClient = GlueClient.builder()
 .region(region)
 .build();

 getSpecificDatabase(glueClient, databaseName);
 glueClient.close();
 }

 public static void getSpecificDatabase(GlueClient glueClient, String
databaseName) {
 try {
 GetDatabaseRequest databasesRequest = GetDatabaseRequest.builder()
 .name(databaseName)
 .build();

 GetDatabaseResponse response =
glueClient.getDatabase(databasesRequest);
 Instant createDate = response.database().createTime();

 // Convert the Instant to readable date.
 DateTimeFormatter formatter =
DateTimeFormatter.ofLocalizedDateTime(FormatStyle.SHORT)

```

```
 .withLocale(Locale.US)
 .withZone(ZoneId.systemDefault());

 formatter.format(createDate);
 System.out.println("The create date of the database is " +
createDate);

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
}
}
```

- APIの詳細については、「APIリファレンス[GetDatabase](#)」の「」を参照してください。  
AWS SDK for Java 2.x

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWSコード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const getDatabase = (name) => {
 const client = new GlueClient({});

 const command = new GetDatabaseCommand({
 Name: name,
 });

 return client.send(command);
};
```

- APIの詳細については、「APIリファレンス[GetDatabase](#)」の「」を参照してください。  
AWS SDK for JavaScript

## Kotlin

### SDK for Kotlin

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
suspend fun getSpecificDatabase(databaseName: String?) {
 val request =
 GetDatabaseRequest {
 name = databaseName
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 val response = glueClient.getDatabase(request)
 val dbDesc = response.database?.description
 println("The database description is $dbDesc")
 }
}
```

- API の詳細については、AWS SDK for Kotlin API リファレンス [GetDatabase](#) の「」を参照してください。

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
$databaseName = "doc-example-database-$uniqid";
```

```
$database = $glueService->getDatabase($databaseName);
echo "Found a database named " . $database['Database']['Name'] . "\n";

public function getDatabase(string $databaseName): Result
{
 return $this->customWaiter(function () use ($databaseName) {
 return $this->glueClient->getDatabase([
 'Name' => $databaseName,
]);
 });
}
```

- APIの詳細については、「APIリファレンス[GetDatabase](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

 def get_database(self, name):
 """
 Gets information about a database in your Data Catalog.

 :param name: The name of the database to look up.
```

```
:return: Information about the database.
"""
try:
 response = self.glue_client.get_database(Name=name)
except ClientError as err:
 logger.error(
 "Couldn't get database %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
else:
 return response["Database"]
```

- APIの詳細については、[GetDatabase](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
```

```
end

Retrieves information about a specific database.
#
@param name [String] The name of the database to retrieve information about.
@return [Aws::Glue::Types::Database, nil] The database object if found, or
nil if not found.
def get_database(name)
 response = @glue_client.get_database(name: name)
 response.database
rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not get database #{name}: \n#{e.message}")
 raise
end
```

- APIの詳細については、「APIリファレンス[GetDatabase](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
let database = glue
 .get_database()
 .name(self.database())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?
 .to_owned();
let database = database
 .database()
 .ok_or_else(|| GlueMvpError::Unknown("Could not find
database".into()))?;
```

- API の詳細については、[GetDatabase](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください。[AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI **GetDatabases**で を使用する

以下のコード例は、GetDatabases の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

### CLI

#### AWS CLI

AWS Glue データカタログの一部またはすべてのデータベースの定義を一覧表示するには、次の `get-databases` の例では、データカタログのデータベースに関する情報を返します。

```
aws glue get-databases
```

出力:

```
{
 "DatabaseList": [
 {
 "Name": "default",
 "Description": "Default Hive database",
 "LocationUri": "file:/spark-warehouse",
 "CreateTime": 1602084052.0,
 "CreateTableDefaultPermissions": [
 {
 "Principal": {
 "DataLakePrincipalIdentifier": "IAM_ALLOWED_PRINCIPALS"
 },
 "Permissions": [
```

```
 "ALL"
]
 }
],
 "CatalogId": "111122223333"
},
{
 "Name": "flights-db",
 "CreateTime": 1587072847.0,
 "CreateTableDefaultPermissions": [
 {
 "Principal": {
 "DataLakePrincipalIdentifier": "IAM_ALLOWED_PRINCIPALS"
 },
 "Permissions": [
 "ALL"
]
 }
],
 "CatalogId": "111122223333"
},
{
 "Name": "legislators",
 "CreateTime": 1601415625.0,
 "CreateTableDefaultPermissions": [
 {
 "Principal": {
 "DataLakePrincipalIdentifier": "IAM_ALLOWED_PRINCIPALS"
 },
 "Permissions": [
 "ALL"
]
 }
],
 "CatalogId": "111122223333"
},
{
 "Name": "tempdb",
 "CreateTime": 1601498566.0,
 "CreateTableDefaultPermissions": [
 {
 "Principal": {
 "DataLakePrincipalIdentifier": "IAM_ALLOWED_PRINCIPALS"
 },

```

```
 "Permissions": [
 "ALL"
]
 },
],
 "CatalogId": "111122223333"
}
]
```

詳細については、「AWS Glue デベロッパーガイド」の「[データカタログにデータベースを定義する](#)」を参照してください。

- API の詳細については、「[コマンドリファレンスGetDatabases](#)」の「」を参照してください。AWS CLI

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const getDatabases = () => {
 const client = new GlueClient({});

 const command = new GetDatabasesCommand({});

 return client.send(command);
};
```

- API の詳細については、「[API リファレンスGetDatabases](#)」の「」を参照してください。AWS SDK for JavaScript

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI **GetJob**で を使用する

以下のコード例は、GetJob の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

### CLI

#### AWS CLI

ジョブに関する情報を取得するには

次の `get-job` の例では、ジョブに関する情報を取得します。

```
aws glue get-job \
 --job-name my-testing-job
```

出力:

```
{
 "Job": {
 "Name": "my-testing-job",
 "Role": "Glue_DefaultRole",
 "CreatedOn": 1602805698.167,
 "LastModifiedOn": 1602805698.167,
 "ExecutionProperty": {
 "MaxConcurrentRuns": 1
 },
 "Command": {
 "Name": "gluestreaming",
 "ScriptLocation": "s3://janetst-bucket-01/Scripts/test_script.scala",
 "PythonVersion": "2"
 },
 "DefaultArguments": {
```

```
 "--class": "GlueApp",
 "--job-language": "scala"
 },
 "MaxRetries": 0,
 "AllocatedCapacity": 10,
 "MaxCapacity": 10.0,
 "GlueVersion": "1.0"
}
}
```

詳細については、「AWS Glue デベロッパーガイド」の「[ジョブ](#)」を参照してください。

- API の詳細については、「[コマンドリファレンスGetJob](#)」の「」を参照してください。  
AWS CLI

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const getJob = (jobName) => {
 const client = new GlueClient({});

 const command = new GetJobCommand({
 JobName: jobName,
 });

 return client.send(command);
};
```

- API の詳細については、「[API リファレンスGetJob](#)」の「」を参照してください。 AWS SDK for JavaScript

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI `GetJobRun` で使用する

以下のコード例は、`GetJobRun` の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Get information about a specific AWS Glue job run.
/// </summary>
/// <param name="jobName">The name of the job.</param>
/// <param name="jobRunId">The Id of the job run.</param>
/// <returns>A JobRun object with information about the job run.</returns>
public async Task<JobRun> GetJobRunAsync(string jobName, string jobRunId)
{
 var response = await _amazonGlue.GetJobRunAsync(new GetJobRunRequest
 { JobName = jobName, RunId = jobRunId });
 return response.JobRun;
}
```

- API の詳細については、「API リファレンス [GetJobRun](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

## SDK for C++

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::GetJobRunRequest jobRunRequest;
jobRunRequest.SetJobName(jobName);
jobRunRequest.SetRunId(jobRunID);

Aws::Glue::Model::GetJobRunOutcome jobRunOutcome = client.GetJobRun(
 jobRunRequest);

if (jobRunOutcome.IsSuccess()) {
 std::cout << "Displaying the job run JSON description." << std::endl;
 std::cout
 <<
jobRunOutcome.GetResult().GetJobRun().Jsonize().View().WriteReadable()
 << std::endl;
}
else {
 std::cerr << "Error get a job run. "
 << jobRunOutcome.GetError().GetMessage()
 << std::endl;
}
```

- API の詳細については、「API リファレンス [GetJobRun](#)」の「」を参照してください。  
AWS SDK for C++

## CLI

## AWS CLI

ジョブの実行に関する情報を取得するには

次の `get-job-run` の例では、ジョブ実行に関する情報を取得します。

```
aws glue get-job-run \
 --job-name "Combine legislators data" \
 --run-id jr_012e176506505074d94d761755e5c62538ee1aad6f17d39f527e9140cf0c9a5e
```

出力:

```
{
 "JobRun": {
 "Id":
 "jr_012e176506505074d94d761755e5c62538ee1aad6f17d39f527e9140cf0c9a5e",
 "Attempt": 0,
 "JobName": "Combine legislators data",
 "StartedOn": 1602873931.255,
 "LastModifiedOn": 1602874075.985,
 "CompletedOn": 1602874075.985,
 "JobRunState": "SUCCEEDED",
 "Arguments": {
 "--enable-continuous-cloudwatch-log": "true",
 "--enable-metrics": "",
 "--enable-spark-ui": "true",
 "--job-bookmark-option": "job-bookmark-enable",
 "--spark-event-logs-path": "s3://aws-glue-assets-111122223333-us-east-1/sparkHistoryLogs/"
 },
 "PredecessorRuns": [],
 "AllocatedCapacity": 10,
 "ExecutionTime": 117,
 "Timeout": 2880,
 "MaxCapacity": 10.0,
 "WorkerType": "G.1X",
 "NumberOfWorkers": 10,
 "LogGroupName": "/aws-glue/jobs",
 "GlueVersion": "2.0"
 }
}
```

詳細については、「AWS Glue デベロッパーガイド」の「[ジョブの実行](#)」を参照してください。

- API の詳細については、「[コマンドリファレンスGetJobRun](#)」の「」を参照してください。AWS CLI

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const getJobRun = (jobName, jobRunId) => {
 const client = new GlueClient({});
 const command = new GetJobRunCommand({
 JobName: jobName,
 RunId: jobRunId,
 });

 return client.send(command);
};
```

- API の詳細については、「[API リファレンスGetJobRun](#)」の「」を参照してください。AWS SDK for JavaScript

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
$jobName = 'test-job-' . $uniqid;

$outputBucketUrl = "s3://$bucketName";
$runId = $glueService->startJobRun($jobName, $databaseName, $tables,
$outputBucketUrl)['JobRunId'];

echo "waiting for job";
do {
 $jobRun = $glueService->getJobRun($jobName, $runId);
 echo ".";
 sleep(10);
} while (!array_intersect([$jobRun['JobRun']['JobRunState']],
['SUCCEEDED', 'STOPPED', 'FAILED', 'TIMEOUT']));
echo "\n";

public function getJobRun($jobName, $runId, $predecessorsIncluded = false):
Result
{
 return $this->glueClient->getJobRun([
 'JobName' => $jobName,
 'RunId' => $runId,
 'PredecessorsIncluded' => $predecessorsIncluded,
]);
}
```

- APIの詳細については、「APIリファレンス[GetJobRun](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""
```

```
def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

def get_job_run(self, name, run_id):
 """
 Gets information about a single job run.

 :param name: The name of the job definition for the run.
 :param run_id: The ID of the run.
 :return: Information about the run.
 """
 try:
 response = self.glue_client.get_job_run(JobName=name, RunId=run_id)
 except ClientError as err:
 logger.error(
 "Couldn't get job run %s/%s. Here's why: %s: %s",
 name,
 run_id,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
 else:
 return response["JobRun"]
```

- APIの詳細については、[GetJobRun](#) AWS 「SDK for Python (Boto3) API リファレンス」の「」を参照してください。

## Ruby

## SDK for Ruby

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
 a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
 for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
 calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Retrieves data for a specific job run.
 #
 # @param job_name [String] The name of the job run to retrieve data for.
 # @return [Glue::Types::GetJobRunResponse]
 def get_job_run(job_name, run_id)
 @glue_client.get_job_run(job_name: job_name, run_id: run_id)
 rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not get job runs: \n#{e.message}")
 end
end
```

- APIの詳細については、「API リファレンス [GetJobRun](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

## SDK for Rust

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
let get_job_run = || async {
 Ok:::<JobRun, GlueMvpError>(
 glue.get_job_run()
 .job_name(self.job())
 .run_id(job_run_id.to_string())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?
 .job_run()
 .ok_or_else(|| GlueMvpError::Unknown("Failed to get
job_run".into()))?
 .to_owned(),
)
};

let mut job_run = get_job_run().await?;
let mut state =
job_run.job_run_state().unwrap_or(&unknown_state).to_owned();

while matches!(
 state,
 JobRunState::Starting | JobRunState::Stopping | JobRunState::Running
) {
 info!(?state, "Waiting for job to finish");
 tokio::time::sleep(self.wait_delay).await;

 job_run = get_job_run().await?;
 state = job_run.job_run_state().unwrap_or(&unknown_state).to_owned();
}
```

- APIの詳細については、[GetJobRun](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください。[AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI **GetJobRuns**で を使用する

以下のコード例は、GetJobRuns の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Get information about all AWS Glue runs of a specific job.
/// </summary>
/// <param name="jobName">The name of the job.</param>
/// <returns>A list of JobRun objects.</returns>
public async Task<List<JobRun>> GetJobRunsAsync(string jobName)
{
 var jobRuns = new List<JobRun>();

 var request = new GetJobRunsRequest
 {
 JobName = jobName,
 };
};
```

```
// No need to loop to get all the log groups--the SDK does it for us
behind the scenes
var paginatorForJobRuns =
 _amazonGlue.Paginators.GetJobRuns(request);

await foreach (var response in paginatorForJobRuns.Responses)
{
 response.JobRuns.ForEach(jobRun =>
 {
 jobRuns.Add(jobRun);
 });
}

return jobRuns;
}
```

- APIの詳細については、「APIリファレンス[GetJobRuns](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

### SDK for C++

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::GetJobRunsRequest getJobRunsRequest;
getJobRunsRequest.SetJobName(jobName);
```

```
Aws::String nextToken; // Used for pagination.
std::vector<Aws::Glue::Model::JobRun> allJobRuns;
do {
 if (!nextToken.empty()) {
 getJobRunsRequest.SetNextToken(nextToken);
 }
 Aws::Glue::Model::GetJobRunsOutcome jobRunsOutcome =
client.GetJobRuns(
 getJobRunsRequest);

 if (jobRunsOutcome.IsSuccess()) {
 const std::vector<Aws::Glue::Model::JobRun> &jobRuns =
jobRunsOutcome.GetResult().GetJobRuns();
 allJobRuns.insert(allJobRuns.end(), jobRuns.begin(),
jobRuns.end());

 nextToken = jobRunsOutcome.GetResult().GetNextToken();
 }
 else {
 std::cerr << "Error getting job runs. "
 << jobRunsOutcome.GetError().GetMessage()
 << std::endl;

 break;
 }
} while (!nextToken.empty());
```

- APIの詳細については、「APIリファレンス[GetJobRuns](#)」の「」を参照してください。  
AWS SDK for C++

## CLI

### AWS CLI

ジョブのすべてのジョブ実行に関する情報を取得するには

次の `get-job-runs` の例では、ジョブのジョブ実行に関する情報を取得します。

```
aws glue get-job-runs \
 --job-name "my-testing-job"
```

出力:

```
{
 "JobRuns": [
 {
 "Id":
"jr_012e176506505074d94d761755e5c62538ee1aad6f17d39f527e9140cf0c9a5e",
 "Attempt": 0,
 "JobName": "my-testing-job",
 "StartedOn": 1602873931.255,
 "LastModifiedOn": 1602874075.985,
 "CompletedOn": 1602874075.985,
 "JobRunState": "SUCCEEDED",
 "Arguments": {
 "--enable-continuous-cloudwatch-log": "true",
 "--enable-metrics": "",
 "--enable-spark-ui": "true",
 "--job-bookmark-option": "job-bookmark-enable",
 "--spark-event-logs-path": "s3://aws-glue-assets-111122223333-us-
east-1/sparkHistoryLogs/"
 },
 "PredecessorRuns": [],
 "AllocatedCapacity": 10,
 "ExecutionTime": 117,
 "Timeout": 2880,
 "MaxCapacity": 10.0,
 "WorkerType": "G.1X",
 "NumberOfWorkers": 10,
 "LogGroupName": "/aws-glue/jobs",
 "GlueVersion": "2.0"
 },
 {
 "Id":
"jr_03cc19ddab11c4e244d3f735567de74ff93b0b3ef468a713ffe73e53d1aec08f_attempt_2",
 "Attempt": 2,
 "PreviousRunId":
"jr_03cc19ddab11c4e244d3f735567de74ff93b0b3ef468a713ffe73e53d1aec08f_attempt_1",
 "JobName": "my-testing-job",
 "StartedOn": 1602811168.496,
 "LastModifiedOn": 1602811282.39,
 "CompletedOn": 1602811282.39,
 "JobRunState": "FAILED",
 "ErrorMessage": "An error occurred while calling
o122.pyWriteDynamicFrame."
 }
]
}
```

```

 Access Denied (Service: Amazon S3; Status Code: 403; Error Code:
AccessDenied;
 Request ID: 021AAB703DB20A2D;
 S3 Extended Request ID: teZk24Y09TkXzBvMPG502L5VJBhe9DJuWA9/
TXtuG0qfByajkfL/Tlqt5JBGdEGpigAqzdMDM/U=)",
 "PredecessorRuns": [],
 "AllocatedCapacity": 10,
 "ExecutionTime": 110,
 "Timeout": 2880,
 "MaxCapacity": 10.0,
 "WorkerType": "G.1X",
 "NumberOfWorkers": 10,
 "LogGroupName": "/aws-glue/jobs",
 "GlueVersion": "2.0"
 },
 {
 "Id":
"jr_03cc19ddab11c4e244d3f735567de74ff93b0b3ef468a713ffe73e53d1aec08f_attempt_1",
 "Attempt": 1,
 "PreviousRunId":
"jr_03cc19ddab11c4e244d3f735567de74ff93b0b3ef468a713ffe73e53d1aec08f",
 "JobName": "my-testing-job",
 "StartedOn": 1602811020.518,
 "LastModifiedOn": 1602811138.364,
 "CompletedOn": 1602811138.364,
 "JobRunState": "FAILED",
 "ErrorMessage": "An error occurred while calling
o122.pyWriteDynamicFrame.
 Access Denied (Service: Amazon S3; Status Code: 403; Error Code:
AccessDenied;
 Request ID: 2671D37856AE7ABB;
 S3 Extended Request ID: RLJCJw20brV
+PpC6Gp0RahyF2fp9f1B5SSb2bTGPnUSPVizLXR11PN3QZ1db+v1o9qRVktNYbW8=)",
 "PredecessorRuns": [],
 "AllocatedCapacity": 10,
 "ExecutionTime": 113,
 "Timeout": 2880,
 "MaxCapacity": 10.0,
 "WorkerType": "G.1X",
 "NumberOfWorkers": 10,
 "LogGroupName": "/aws-glue/jobs",
 "GlueVersion": "2.0"
 }
]

```

```
}
```

詳細については、「AWS Glue デベロッパーガイド」の「[ジョブの実行](#)」を参照してください。

- API の詳細については、「[コマンドリファレンスGetJobRuns](#)」の「」を参照してください。AWS CLI

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const getJobRuns = (jobName) => {
 const client = new GlueClient({});
 const command = new GetJobRunsCommand({
 JobName: jobName,
 });

 return client.send(command);
};
```

- API の詳細については、「[API リファレンスGetJobRuns](#)」の「」を参照してください。AWS SDK for JavaScript

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
$jobName = 'test-job-' . $uniqid;

$jobRuns = $glueService->getJobRuns($jobName);

public function getJobRuns($jobName, $maxResults = 0, $nextToken = ''):
Result
{
 $arguments = ['JobName' => $jobName];
 if ($maxResults) {
 $arguments['MaxResults'] = $maxResults;
 }
 if ($nextToken) {
 $arguments['NextToken'] = $nextToken;
 }
 return $this->glueClient->getJobRuns($arguments);
}
```

- APIの詳細については、「APIリファレンス[GetJobRuns](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWSコード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client
```

```
def get_job_runs(self, job_name):
 """
 Gets information about runs that have been performed for a specific job
 definition.

 :param job_name: The name of the job definition to look up.
 :return: The list of job runs.
 """
 try:
 response = self.glue_client.get_job_runs(JobName=job_name)
 except ClientError as err:
 logger.error(
 "Couldn't get job runs for %s. Here's why: %s: %s",
 job_name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
 else:
 return response["JobRuns"]
```

- APIの詳細については、[GetJobRuns](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
```

```
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Retrieves a list of job runs for the specified job.
 #
 # @param job_name [String] The name of the job to retrieve job runs for.
 # @return [Array<Aws::Glue::Types::JobRun>]
 def get_job_runs(job_name)
 response = @glue_client.get_job_runs(job_name: job_name)
 response.job_runs
 rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not get job runs: \n#{e.message}")
 end
end
```

- APIの詳細については、「APIリファレンス[GetJobRuns](#)」の「」を参照してください。  
AWS SDK for Ruby

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください[AWS SDKでこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前のSDKバージョンの詳細も含まれています。

## AWS SDK または CLI **GetTables**で使用する

以下のコード例は、GetTables の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

## .NET

### AWS SDK for .NET

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Get a list of tables for an AWS Glue database.
/// </summary>
/// <param name="dbName">The name of the database.</param>
/// <returns>A list of Table objects.</returns>
public async Task<List<Table>> GetTablesAsync(string dbName)
{
 var request = new GetTablesRequest { DatabaseName = dbName };
 var tables = new List<Table>();

 // Get a paginator for listing the tables.
 var tablePaginator = _amazonGlue.Paginators.GetTables(request);

 await foreach (var response in tablePaginator.Responses)
 {
 tables.AddRange(response.TableList);
 }

 return tables;
}
```

- API の詳細については、「API リファレンス [GetTables](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

## SDK for C++

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::GetTablesRequest request;
request.SetDatabaseName(CRAWLER_DATABASE_NAME);
std::vector<Aws::Glue::Model::Table> all_tables;
Aws::String nextToken; // Used for pagination.
do {
 Aws::Glue::Model::GetTablesOutcome outcome =
client.GetTables(request);

 if (outcome.IsSuccess()) {
 const std::vector<Aws::Glue::Model::Table> &tables =
outcome.GetResult().GetTableList();
 all_tables.insert(all_tables.end(), tables.begin(),
tables.end());
 nextToken = outcome.GetResult().GetNextToken();
 }
 else {
 std::cerr << "Error getting the tables. "
 << outcome.GetError().GetMessage()
 << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, "", bucketName,
clientConfig);
 return false;
 }
} while (!nextToken.empty());
```

```
std::cout << "The database contains " << all_tables.size()
 << (all_tables.size() == 1 ?
 " table." : "tables.") << std::endl;
std::cout << "Here is a list of the tables in the database.";
for (size_t index = 0; index < all_tables.size(); ++index) {
 std::cout << " " << index + 1 << ": " <<
all_tables[index].GetName()
 << std::endl;
}

if (!all_tables.empty()) {
 int tableIndex = askQuestionForIntRange(
 "Enter an index to display the database detail ",
 1, static_cast<int>(all_tables.size()));
 std::cout << all_tables[tableIndex -
1].Jsonize().View().WriteReadable()
 << std::endl;

 tableName = all_tables[tableIndex - 1].GetName();
}
```

- API の詳細については、「API リファレンス[GetTables](#)」の「」を参照してください。  
AWS SDK for C++

## CLI

### AWS CLI

指定されたデータベース内で、一部またはすべてのテーブルの定義を一覧表示するには

次の `get-tables` の例では、指定されたデータベース内のテーブルに関する情報を返します。

```
aws glue get-tables --database-name 'tempdb'
```

出力:

```
{
 "TableList": [
 {
 "Name": "my-s3-sink",
```

```
"DatabaseName": "tempdb",
"CreateTime": 1602730539.0,
"UpdateTime": 1602730539.0,
"Retention": 0,
"StorageDescriptor": {
 "Columns": [
 {
 "Name": "sensorid",
 "Type": "int"
 },
 {
 "Name": "currenttemperature",
 "Type": "int"
 },
 {
 "Name": "status",
 "Type": "string"
 }
],
 "Location": "s3://janetst-bucket-01/test-s3-output/",
 "Compressed": false,
 "NumberOfBuckets": 0,
 "SerdeInfo": {
 "SerializationLibrary": "org.openx.data.jsonserde.JsonSerDe"
 },
 "SortColumns": [],
 "StoredAsSubDirectories": false
},
"Parameters": {
 "classification": "json"
},
"CreatedBy": "arn:aws:iam::007436865787:user/JRSTERN",
"IsRegisteredWithLakeFormation": false,
"CatalogId": "007436865787"
},
{
 "Name": "s3-source",
 "DatabaseName": "tempdb",
 "CreateTime": 1602730658.0,
 "UpdateTime": 1602730658.0,
 "Retention": 0,
 "StorageDescriptor": {
 "Columns": [
 {
```

```
 "Name": "sensorid",
 "Type": "int"
 },
 {
 "Name": "currenttemperature",
 "Type": "int"
 },
 {
 "Name": "status",
 "Type": "string"
 }
],
"Location": "s3://janetst-bucket-01/",
"Compressed": false,
"NumberOfBuckets": 0,
"SortColumns": [],
"StoredAsSubDirectories": false
},
"Parameters": {
 "classification": "json"
},
"CreatedBy": "arn:aws:iam::007436865787:user/JRSTERN",
"IsRegisteredWithLakeFormation": false,
"CatalogId": "007436865787"
},
{
 "Name": "test-kinesis-input",
 "DatabaseName": "tempdb",
 "CreateTime": 1601507001.0,
 "UpdateTime": 1601507001.0,
 "Retention": 0,
 "StorageDescriptor": {
 "Columns": [
 {
 "Name": "sensorid",
 "Type": "int"
 },
 {
 "Name": "currenttemperature",
 "Type": "int"
 },
 {
 "Name": "status",
 "Type": "string"
 }
]
 }
}
```

```
 }
],
 "Location": "my-testing-stream",
 "Compressed": false,
 "NumberOfBuckets": 0,
 "SerdeInfo": {
 "SerializationLibrary": "org.openx.data.jsonserde.JsonSerDe"
 },
 "SortColumns": [],
 "Parameters": {
 "kinesisUrl": "https://kinesis.us-east-1.amazonaws.com",
 "streamName": "my-testing-stream",
 "typeOfData": "kinesis"
 },
 "StoredAsSubDirectories": false
},
"Parameters": {
 "classification": "json"
},
"CreatedBy": "arn:aws:iam::007436865787:user/JRSTERN",
"IsRegisteredWithLakeFormation": false,
"CatalogId": "007436865787"
}
]
}
```

詳細については、[AWS Glue デベロッパーガイドの「Glue データカタログでのテーブルの定義」](#)を参照してください。AWS

- API の詳細については、「[コマンドリファレンスGetTables](#)」の「」を参照してください。  
AWS CLI

## Java

### SDK for Java 2.x

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
import software.amazon.awssdk.regions.Region;
import software.amazon.awssdk.services.glue.GlueClient;
import software.amazon.awssdk.services.glue.model.GetTableRequest;
import software.amazon.awssdk.services.glue.model.GetTableResponse;
import software.amazon.awssdk.services.glue.model.GlueException;
import java.time.Instant;
import java.time.ZoneId;
import java.time.format.DateTimeFormatter;
import java.time.format.FormatStyle;
import java.util.Locale;

/**
 * Before running this Java V2 code example, set up your development
 * environment, including your credentials.
 *
 * For more information, see the following documentation topic:
 *
 * https://docs.aws.amazon.com/sdk-for-java/latest/developer-guide/get-started.html
 */
public class GetTable {
 public static void main(String[] args) {
 final String usage = ""

 Usage:
 <dbName> <tableName>

 Where:
 dbName - The database name.\s
 tableName - The name of the table.\s
 """;

 if (args.length != 2) {
 System.out.println(usage);
 System.exit(1);
 }

 String dbName = args[0];
 String tableName = args[1];
 Region region = Region.US_EAST_1;
 GlueClient glueClient = GlueClient.builder()
 .region(region)
 .build();
 }
}
```

```
 getGlueTable(glueClient, dbName, tableName);
 glueClient.close();
 }

 public static void getGlueTable(GlueClient glueClient, String dbName, String
tableName) {
 try {
 GetTableRequest tableRequest = GetTableRequest.builder()
 .databaseName(dbName)
 .name(tableName)
 .build();

 GetTableResponse tableResponse = glueClient.getTable(tableRequest);
 Instant createDate = tableResponse.table().createTime();

 // Convert the Instant to readable date.
 DateTimeFormatter formatter =
 DateTimeFormatter.ofLocalizedDateTime(FormatStyle.SHORT)
 .withLocale(Locale.US)
 .withZone(ZoneId.systemDefault());

 formatter.format(createDate);
 System.out.println("The create date of the table is " + createDate);

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
 }
}
```

- APIの詳細については、「APIリファレンス[GetTables](#)」の「」を参照してください。  
AWS SDK for Java 2.x

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const getTables = (databaseName) => {
 const client = new GlueClient({});

 const command = new GetTablesCommand({
 DatabaseName: databaseName,
 });

 return client.send(command);
};
```

- API の詳細については、「API リファレンス [GetTables](#)」の「」を参照してください。  
AWS SDK for JavaScript

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
$databaseName = "doc-example-database-$uniqid";

$tables = $glueService->getTables($databaseName);

public function getTables($databaseName): Result
{
```

```
return $this->glueClient->getTables([
 'DatabaseName' => $databaseName,
]);
}
```

- APIの詳細については、「APIリファレンス[GetTables](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

 def get_tables(self, db_name):
 """
 Gets a list of tables in a Data Catalog database.

 :param db_name: The name of the database to query.
 :return: The list of tables in the database.
 """
 try:
 response = self.glue_client.get_tables(DatabaseName=db_name)
 except ClientError as err:
 logger.error(
 "Couldn't get tables %s. Here's why: %s: %s",
```

```
 db_name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
else:
 return response["TableList"]
```

- APIの詳細については、[GetTables](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Retrieves a list of tables in the specified database.
 #
 # @param db_name [String] The name of the database to retrieve tables from.
 # @return [Array<Aws::Glue::Types::Table>]
 def get_tables(db_name)
```

```
response = @glue_client.get_tables(database_name: db_name)
response.table_list
rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not get tables #{db_name}: \n#{e.message}")
 raise
end
```

- APIの詳細については、「APIリファレンス[GetTables](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
let tables = glue
 .get_tables()
 .database_name(self.database())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;

let tables = tables.table_list();
```

- APIの詳細については、 [GetTables](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI `ListJobs` を使用する

以下のコード例は、`ListJobs` の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// List AWS Glue jobs using a paginator.
/// </summary>
/// <returns>A list of AWS Glue job names.</returns>
public async Task<List<string>> ListJobsAsync()
{
 var jobNames = new List<string>();

 var listJobsPaginator = _amazonGlue.Paginators.ListJobs(new
ListJobsRequest { MaxResults = 10 });
 await foreach (var response in listJobsPaginator.Responses)
 {
 jobNames.AddRange(response.JobNames);
 }

 return jobNames;
}
```

- API の詳細については、「API リファレンス [ListJobs](#)」の「」を参照してください。 AWS SDK for .NET

## C++

## SDK for C++

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::ListJobsRequest listJobsRequest;

Aws::String nextToken;
std::vector<Aws::String> allJobNames;

do {
 if (!nextToken.empty()) {
 listJobsRequest.SetNextToken(nextToken);
 }
 Aws::Glue::Model::ListJobsOutcome listRunsOutcome = client.ListJobs(
 listJobsRequest);

 if (listRunsOutcome.IsSuccess()) {
 const std::vector<Aws::String> &jobNames =
listRunsOutcome.GetResult().GetJobNames();
 allJobNames.insert(allJobNames.end(), jobNames.begin(),
jobNames.end());
 nextToken = listRunsOutcome.GetResult().GetNextToken();
 }
 else {
 std::cerr << "Error listing jobs. "
 << listRunsOutcome.GetError().GetMessage()
 << std::endl;
 }
} while (!nextToken.empty());
```

- API の詳細については、「API リファレンス[ListJobs](#)」の「」を参照してください。AWS SDK for C++

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const listJobs = () => {
 const client = new GlueClient({});

 const command = new ListJobsCommand({});

 return client.send(command);
};
```

- API の詳細については、「API リファレンス[ListJobs](#)」の「」を参照してください。AWS SDK for JavaScript

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
$jobs = $glueService->listJobs();
echo "Current jobs:\n";
```

```
 foreach ($jobs['JobNames'] as $jobsName) {
 echo "{$jobsName}\n";
 }

 public function listJobs($maxResults = null, $nextToken = null, $tags = []):
 Result
 {
 $arguments = [];
 if ($maxResults) {
 $arguments['MaxResults'] = $maxResults;
 }
 if ($nextToken) {
 $arguments['NextToken'] = $nextToken;
 }
 if (!empty($tags)) {
 $arguments['Tags'] = $tags;
 }
 return $this->glueClient->listJobs($arguments);
 }
 }
```

- APIの詳細については、「APIリファレンス[ListJobs](#)」の「」を参照してください。AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください。GitHub。 [AWSコード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
```

```
self.glue_client = glue_client

def list_jobs(self):
 """
 Lists the names of job definitions in your account.

 :return: The list of job definition names.
 """
 try:
 response = self.glue_client.list_jobs()
 except ClientError as err:
 logger.error(
 "Couldn't list jobs. Here's why: %s: %s",
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
 else:
 return response["JobNames"]
```

- APIの詳細については、[ListJobs](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
```

```
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Retrieves a list of jobs in AWS Glue.
 #
 # @return [Aws::Glue::Types::ListJobsResponse]
 def list_jobs
 @glue_client.list_jobs
 rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not list jobs: \n#{e.message}")
 raise
 end
end
```

- APIの詳細については、「APIリファレンス[ListJobs](#)」の「」を参照してください。AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
let mut list_jobs = glue.list_jobs().into_paginator().send();
while let Some(list_jobs_output) = list_jobs.next().await {
 match list_jobs_output {
 Ok(list_jobs) => {
 let names = list_jobs.job_names();
 info!(?names, "Found these jobs")
 }
 Err(err) => return Err(GlueMvpError::from_glue_sdk(err)),
 }
}
```

```
}
```

- APIの詳細については、[ListJobs](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください。[AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI `StartCrawler`で を使用する

以下のコード例は、`StartCrawler` の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Start an AWS Glue crawler.
/// </summary>
/// <param name="crawlerName">The name of the crawler.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> StartCrawlerAsync(string crawlerName)
{
 var crawlerRequest = new StartCrawlerRequest
 {
 Name = crawlerName,
 };
}
```

```
var response = await _amazonGlue.StartCrawlerAsync(crawlerRequest);

return response.HttpStatusCode == System.Net.HttpStatusCode.OK;
}
```

- APIの詳細については、「API リファレンス [StartCrawler](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

### SDK for C++

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::StartCrawlerRequest request;
request.SetName(CRAWLER_NAME);

Aws::Glue::Model::StartCrawlerOutcome outcome =
client.StartCrawler(request);

if (outcome.IsSuccess() || (Aws::Glue::GlueErrors::CRAWLER_RUNNING ==
 outcome.GetError().GetErrorType())) {
 if (!outcome.IsSuccess()) {
 std::cout << "Crawler was already started." << std::endl;
 }
}
else {
```

```

 std::cout << "Successfully started crawler." << std::endl;
 }

 std::cout << "This may take a while to run." << std::endl;

 Aws::Glue::Model::CrawlerState crawlerState =
 Aws::Glue::Model::CrawlerState::NOT_SET;
 int iterations = 0;
 while (Aws::Glue::Model::CrawlerState::READY != crawlerState) {
 std::this_thread::sleep_for(std::chrono::seconds(1));
 ++iterations;
 if ((iterations % 10) == 0) { // Log status every 10 seconds.
 std::cout << "Crawler status " <<

 Aws::Glue::Model::CrawlerStateMapper::GetNameForCrawlerState(
 crawlerState)
 << ". After " << iterations
 << " seconds elapsed."
 << std::endl;
 }
 Aws::Glue::Model::GetCrawlerRequest getCrawlerRequest;
 getCrawlerRequest.SetName(CRAWLER_NAME);

 Aws::Glue::Model::GetCrawlerOutcome getCrawlerOutcome =
 client.GetCrawler(
 getCrawlerRequest);

 if (getCrawlerOutcome.IsSuccess()) {
 crawlerState =
 getCrawlerOutcome.GetResult().GetCrawler().GetState();
 }
 else {
 std::cerr << "Error getting crawler. "
 << getCrawlerOutcome.GetError().GetMessage() <<
 std::endl;
 break;
 }
 }

 if (Aws::Glue::Model::CrawlerState::READY == crawlerState) {
 std::cout << "Crawler finished running after " << iterations
 << " seconds."
 << std::endl;
 }
}

```

```
 }
 else {
 std::cerr << "Error starting a crawler. "
 << outcome.GetError().GetMessage()
 << std::endl;

 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, "", bucketName,
 clientConfig);
 return false;
 }
}
```

- APIの詳細については、「APIリファレンス[StartCrawler](#)」の「」を参照してください。  
AWS SDK for C++

## CLI

### AWS CLI

クローラーを開始するには

次の `start-crawler` の例では、クローラーを開始します。

```
aws glue start-crawler --name my-crawler
```

出力:

```
None
```

詳細については、「AWS Glue デベロッパーガイド」の「[クローラーの定義](#)」を参照してください。

- APIの詳細については、「コマンドリファレンス[StartCrawler](#)」の「」を参照してください。  
AWS CLI

## Java

## SDK for Java 2.x

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
import software.amazon.awssdk.regions.Region;
import software.amazon.awssdk.services.glue.GlueClient;
import software.amazon.awssdk.services.glue.model.GlueException;
import software.amazon.awssdk.services.glue.model.StartCrawlerRequest;

/**
 * Before running this Java V2 code example, set up your development
 * environment, including your credentials.
 *
 * For more information, see the following documentation topic:
 *
 * https://docs.aws.amazon.com/sdk-for-java/latest/developer-guide/get-
 * started.html
 */
public class StartCrawler {
 public static void main(String[] args) {
 final String usage = ""

 Usage:
 <crawlerName>

 Where:
 crawlerName - The name of the crawler.\s
 """;

 if (args.length != 1) {
 System.out.println(usage);
 System.exit(1);
 }

 String crawlerName = args[0];
 Region region = Region.US_EAST_1;
```

```
 GlueClient glueClient = GlueClient.builder()
 .region(region)
 .build();

 startSpecificCrawler(glueClient, crawlerName);
 glueClient.close();
 }

 public static void startSpecificCrawler(GlueClient glueClient, String
crawlerName) {
 try {
 StartCrawlerRequest crawlerRequest = StartCrawlerRequest.builder()
 .name(crawlerName)
 .build();

 glueClient.startCrawler(crawlerRequest);

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
 }
}
```

- APIの詳細については、「APIリファレンス[StartCrawler](#)」の「」を参照してください。  
AWS SDK for Java 2.x

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const startCrawler = (name) => {
 const client = new GlueClient({});
```

```
const command = new StartCrawlerCommand({
 Name: name,
});

return client.send(command);
};
```

- APIの詳細については、「API リファレンス[StartCrawler](#)」の「」を参照してください。  
AWS SDK for JavaScript

## Kotlin

### SDK for Kotlin

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
suspend fun startSpecificCrawler(crawlerName: String?) {
 val request =
 StartCrawlerRequest {
 name = crawlerName
 }

 GlueClient { region = "us-west-2" }.use { glueClient ->
 glueClient.startCrawler(request)
 println("$crawlerName was successfully started.")
 }
}
```

- APIの詳細については、AWS SDK for Kotlin API リファレンス[StartCrawler](#)の「」を参照してください。

## PHP

## SDK for PHP

**Note**

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
$crawlerName = "example-crawler-test-" . $uniqid;

$databaseName = "doc-example-database-$uniqid";

$glueService->startCrawler($crawlerName);

public function startCrawler($crawlerName): Result
{
 return $this->glueClient->startCrawler([
 'Name' => $crawlerName,
]);
}
```

- API の詳細については、「API リファレンス [StartCrawler](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

## SDK for Python (Boto3)

**Note**

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""
```

```
def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

def start_crawler(self, name):
 """
 Starts a crawler. The crawler crawls its configured target and creates
 metadata that describes the data it finds in the target data source.

 :param name: The name of the crawler to start.
 """
 try:
 self.glue_client.start_crawler(Name=name)
 except ClientError as err:
 logger.error(
 "Couldn't start crawler %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
```

- APIの詳細については、[StartCrawler](#) AWS SDK for Python (Boto3) API リファレンスの「」を参照してください。

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Starts a crawler with the specified name.
 #
 # @param name [String] The name of the crawler to start.
 # @return [void]
 def start_crawler(name)
 @glue_client.start_crawler(name: name)
 rescue Aws::Glue::Errors::ServiceError => e
 @logger.error("Glue could not start crawler #{name}: \n#{e.message}")
 raise
 end
end
```

- APIの詳細については、「APIリファレンス[StartCrawler](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
let start_crawler =
glue.start_crawler().name(self.crawler()).send().await;

match start_crawler {
```

```
Ok(_) => Ok(()),
Err(err) => {
 let glue_err: aws_sdk_glue::Error = err.into();
 match glue_err {
 aws_sdk_glue::Error::CrawlerRunningException(_) => Ok(()),
 _ => Err(GlueMvpError::GlueSdk(glue_err)),
 }
}
}?:;
```

- APIの詳細については、[StartCrawler](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください。[AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDK または CLI `StartJobRun` を使用する

以下のコード例は、`StartJobRun` の使用方法を示しています。

アクション例は、より大きなプログラムからのコードの抜粋であり、コンテキスト内で実行する必要があります。次のコード例で、このアクションのコンテキストを確認できます。

- [クローラーとジョブを開始する](#)

.NET

AWS SDK for .NET

### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/// <summary>
/// Start an AWS Glue job run.
/// </summary>
```

```
/// <param name="jobName">The name of the job.</param>
/// <returns>A string representing the job run Id.</returns>
public async Task<string> StartJobRunAsync(
 string jobName,
 string inputDatabase,
 string inputTable,
 string bucketName)
{
 var request = new StartJobRunRequest
 {
 JobName = jobName,
 Arguments = new Dictionary<string, string>
 {
 {"--input_database", inputDatabase},
 {"--input_table", inputTable},
 {"--output_bucket_url", $"s3://{bucketName}/"}
 }
 };

 var response = await _amazonGlue.StartJobRunAsync(request);
 return response.JobRunId;
}
```

- APIの詳細については、「APIリファレンス[StartJobRun](#)」の「」を参照してください。  
AWS SDK for .NET

## C++

### SDK for C++

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
Aws::Client::ClientConfiguration clientConfig;
// Optional: Set to the AWS Region in which the bucket was created
(overrides config file).
```

```
// clientConfig.region = "us-east-1";

Aws::Glue::GlueClient client(clientConfig);

Aws::Glue::Model::StartJobRunRequest request;
request.SetJobName(JOB_NAME);

Aws::Map<Aws::String, Aws::String> arguments;
arguments["--input_database"] = CRAWLER_DATABASE_NAME;
arguments["--input_table"] = tableName;
arguments["--output_bucket_url"] = Aws::String("s3://") + bucketName +
"/";
request.SetArguments(arguments);

Aws::Glue::Model::StartJobRunOutcome outcome =
client.StartJobRun(request);

if (outcome.IsSuccess()) {
 std::cout << "Successfully started the job." << std::endl;

 Aws::String jobRunId = outcome.GetResult().GetJobRunId();

 int iterator = 0;
 bool done = false;
 while (!done) {
 ++iterator;
 std::this_thread::sleep_for(std::chrono::seconds(1));
 Aws::Glue::Model::GetJobRunRequest jobRunRequest;
 jobRunRequest.SetJobName(JOB_NAME);
 jobRunRequest.SetRunId(jobRunId);

 Aws::Glue::Model::GetJobRunOutcome jobRunOutcome =
client.GetJobRun(
 jobRunRequest);

 if (jobRunOutcome.IsSuccess()) {
 const Aws::Glue::Model::JobRun &jobRun =
jobRunOutcome.GetResult().GetJobRun();
 Aws::Glue::Model::JobRunState jobRunState =
jobRun.GetJobRunState();

 if ((jobRunState == Aws::Glue::Model::JobRunState::STOPPED)
||
 (jobRunState == Aws::Glue::Model::JobRunState::FAILED) ||
```

```

 (jobRunState == Aws::Glue::Model::JobRunState::TIMEOUT))
 {
 std::cerr << "Error running job. "
 << jobRun.GetErrorMessage()
 << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME,
JOB_NAME,
 bucketName,
 clientConfig);
 return false;
 }
 else if (jobRunState ==
 Aws::Glue::Model::JobRunState::SUCCEEDED) {
 std::cout << "Job run succeeded after " << iterator <<
 " seconds elapsed." << std::endl;
 done = true;
 }
 else if ((iterator % 10) == 0) { // Log status every 10
seconds.
 std::cout << "Job run status " <<
 Aws::Glue::Model::JobRunStateMapper::GetNameForJobRunState(
 jobRunState) <<
 ". " << iterator <<
 " seconds elapsed." << std::endl;
 }
 }
 else {
 std::cerr << "Error retrieving job run state. "
 << jobRunOutcome.GetError().GetMessage()
 << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, JOB_NAME,
 bucketName, clientConfig);
 return false;
 }
 }
 }
 else {
 std::cerr << "Error starting a job. " <<
outcome.GetError().GetMessage()
 << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, JOB_NAME,
bucketName,
 clientConfig);
 }
}

```

```
 return false;
 }
```

- APIの詳細については、「APIリファレンス[StartJobRun](#)」の「」を参照してください。  
AWS SDK for C++

## CLI

### AWS CLI

ジョブの実行を開始するには

次の `start-job-run` の例ではジョブを開始します。

```
aws glue start-job-run \
 --job-name my-job
```

出力:

```
{
 "JobRunId":
 "jr_22208b1f44eb5376a60569d4b21dd20fcb8621e1a366b4e7b2494af764b82ded"
}
```

詳細については、「AWS Glue デベロッパーガイド」の「[ジョブの作成](#)」を参照してください。

- APIの詳細については、「コマンドリファレンス[StartJobRun](#)」の「」を参照してください。  
AWS CLI

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
const startJobRun = (jobName, dbName, tableName, bucketName) => {
 const client = new GlueClient({});

 const command = new StartJobRunCommand({
 JobName: jobName,
 Arguments: {
 "--input_database": dbName,
 "--input_table": tableName,
 "--output_bucket_url": `s3://${bucketName}/`,
 },
 });

 return client.send(command);
};
```

- APIの詳細については、「APIリファレンス[StartJobRun](#)」の「」を参照してください。  
AWS SDK for JavaScript

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
$jobName = 'test-job-' . $uniqid;

$databaseName = "doc-example-database-$uniqid";

$tables = $glueService->getTables($databaseName);

$outputBucketUrl = "s3://$bucketName";
$runId = $glueService->startJobRun($jobName, $databaseName, $tables,
$outputBucketUrl)['JobRunId'];

public function startJobRun($jobName, $databaseName, $tables,
$outputBucketUrl): Result
```

```
{
 return $this->glueClient->startJobRun([
 'JobName' => $jobName,
 'Arguments' => [
 'input_database' => $databaseName,
 'input_table' => $tables['TableList'][0]['Name'],
 'output_bucket_url' => $outputBucketUrl,
 '--input_database' => $databaseName,
 '--input_table' => $tables['TableList'][0]['Name'],
 '--output_bucket_url' => $outputBucketUrl,
],
]);
}
```

- APIの詳細については、「APIリファレンス[StartJobRun](#)」の「」を参照してください。  
AWS SDK for PHP

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

 def start_job_run(self, name, input_database, input_table,
output_bucket_name):
 """
```

Starts a job run. A job run extracts data from the source, transforms it, and loads it to the output bucket.

:param name: The name of the job definition.

:param input\_database: The name of the metadata database that contains tables that describe the source data. This is typically created by a crawler.

:param input\_table: The name of the table in the metadata database that describes the source data.

:param output\_bucket\_name: The S3 bucket where the output is written.

:return: The ID of the job run.

"""

try:

# The custom Arguments that are passed to this function are used by the Python ETL script to determine the location of input and output data.

```
response = self.glue_client.start_job_run(
 JobName=name,
 Arguments={
 "--input_database": input_database,
 "--input_table": input_table,
 "--output_bucket_url": f"s3://{output_bucket_name}/",
 },
)
```

except ClientError as err:

```
logger.error(
 "Couldn't start job run %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
```

raise

else:

```
return response["JobRunId"]
```

- APIの詳細については、[StartJobRun](#) AWS 「SDK for Python (Boto3) API リファレンス」の「」を参照してください。

## Ruby

## SDK for Ruby

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
 a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
 for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
 calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Starts a job run for the specified job.
 #
 # @param name [String] The name of the job to start the run for.
 # @param input_database [String] The name of the input database for the job.
 # @param input_table [String] The name of the input table for the job.
 # @param output_bucket_name [String] The name of the output S3 bucket for the
 job.
 # @return [String] The ID of the started job run.
 def start_job_run(name, input_database, input_table, output_bucket_name)
 response = @glue_client.start_job_run(
 job_name: name,
 arguments: {
 '--input_database': input_database,
 '--input_table': input_table,
 '--output_bucket_url': "s3://#{output_bucket_name}/"
 }
)
 response.job_run_id
 rescue Aws::Glue::Errors::GlueException => e
```

```
@logger.error("Glue could not start job run #{name}: \n#{e.message}")
raise
end
```

- APIの詳細については、「APIリファレンス[StartJobRun](#)」の「」を参照してください。  
AWS SDK for Ruby

## Rust

### SDK for Rust

#### Note

については、「」を参照してください [GitHub](#)。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
let job_run_output = glue
 .start_job_run()
 .job_name(self.job())
 .arguments("--input_database", self.database())
 .arguments(
 "--input_table",
 self.tables
 .first()
 .ok_or_else(|| GlueMvpError::Unknown("Missing crawler
table".into()))?
 .name(),
)
 .arguments("--output_bucket_url", self.bucket())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;

let job = job_run_output
 .job_run_id()
 .ok_or_else(|| GlueMvpError::Unknown("Missing run id from just
started job".into()))?
 .to_string();
```

- API の詳細については、[StartJobRun](#) AWS SDK for Rust API リファレンスの「」を参照してください。

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください。[AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

## AWS SDKs AWS Glue を使用するシナリオ

次のコード例は、AWS SDKs を使用して AWS Glue 一般的なシナリオを実装する方法を示しています。これらのシナリオは、内で複数の関数を呼び出して特定のタスクを実行する方法を示しています AWS Glue。各シナリオには GitHub、コードのセットアップと実行の手順を示すへのリンクが含まれています。

### 例

- [AWS SDK を使用して AWS Glue クローラーとジョブの実行を開始する](#)

## AWS SDK を使用して AWS Glue クローラーとジョブの実行を開始する

次のコード例は、以下を実行する方法を示しています。

- パブリック Amazon S3 バケットをクローラし、CSV 形式のメタデータのデータベースを生成するクローラーを作成する。
- のデータベースとテーブルに関する情報を一覧表示します AWS Glue Data Catalog。
- S3 バケットから CSV 形式のデータを抽出するジョブを作成し、そのデータを変換して JSON 形式の出力を別の S3 バケットにロードする。
- ジョブ実行に関する情報を一覧表示し、変換されたデータを表示してリソースをクリーンアップする。

詳細については、「[チュートリアル: AWS Glue Studio の開始方法](#)」を参照してください。

## .NET

### AWS SDK for .NET

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

シナリオで使用される AWS Glue 関数をラップするクラスを作成します。

```
using System.Net;

namespace GlueActions;

public class GlueWrapper
{
 private readonly IAmazonGlue _amazonGlue;

 /// <summary>
 /// Constructor for the AWS Glue actions wrapper.
 /// </summary>
 /// <param name="amazonGlue"></param>
 public GlueWrapper(IAmazonGlue amazonGlue)
 {
 _amazonGlue = amazonGlue;
 }

 /// <summary>
 /// Create an AWS Glue crawler.
 /// </summary>
 /// <param name="crawlerName">The name for the crawler.</param>
 /// <param name="crawlerDescription">A description of the crawler.</param>
 /// <param name="role">The AWS Identity and Access Management (IAM) role to
 /// be assumed by the crawler.</param>
 /// <param name="schedule">The schedule on which the crawler will be
 executed.</param>
 /// <param name="s3Path">The path to the Amazon Simple Storage Service
 (Amazon S3)
 /// bucket where the Python script has been stored.</param>
 /// <param name="dbName">The name to use for the database that will be
```

```
/// created by the crawler.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> CreateCrawlerAsync(
 string crawlerName,
 string crawlerDescription,
 string role,
 string schedule,
 string s3Path,
 string dbName)
{
 var s3Target = new S3Target
 {
 Path = s3Path,
 };

 var targetList = new List<S3Target>
 {
 s3Target,
 };

 var targets = new CrawlerTargets
 {
 S3Targets = targetList,
 };

 var crawlerRequest = new CreateCrawlerRequest
 {
 DatabaseName = dbName,
 Name = crawlerName,
 Description = crawlerDescription,
 Targets = targets,
 Role = role,
 Schedule = schedule,
 };

 var response = await _amazonGlue.CreateCrawlerAsync(crawlerRequest);
 return response.HttpStatusCode == System.Net.HttpStatusCode.OK;
}

/// <summary>
/// Create an AWS Glue job.
/// </summary>
/// <param name="jobName">The name of the job.</param>
```

```
/// <param name="roleName">The name of the IAM role to be assumed by
/// the job.</param>
/// <param name="description">A description of the job.</param>
/// <param name="scriptUrl">The URL to the script.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> CreateJobAsync(string dbName, string tableName,
string bucketUrl, string jobName, string roleName, string description, string
scriptUrl)
{
 var command = new JobCommand
 {
 PythonVersion = "3",
 Name = "glueetl",
 ScriptLocation = scriptUrl,
 };

 var arguments = new Dictionary<string, string>
 {
 { "--input_database", dbName },
 { "--input_table", tableName },
 { "--output_bucket_url", bucketUrl }
 };

 var request = new CreateJobRequest
 {
 Command = command,
 DefaultArguments = arguments,
 Description = description,
 GlueVersion = "3.0",
 Name = jobName,
 NumberOfWorkers = 10,
 Role = roleName,
 WorkerType = "G.1X"
 };

 var response = await _amazonGlue.CreateJobAsync(request);
 return response.HttpStatusCode == HttpStatusCode.OK;
}

/// <summary>
/// Delete an AWS Glue crawler.
/// </summary>
/// <param name="crawlerName">The name of the crawler.</param>
```

```
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> DeleteCrawlerAsync(string crawlerName)
{
 var response = await _amazonGlue.DeleteCrawlerAsync(new
DeleteCrawlerRequest { Name = crawlerName });
 return response.HttpStatusCode == HttpStatusCode.OK;
}

/// <summary>
/// Delete the AWS Glue database.
/// </summary>
/// <param name="dbName">The name of the database.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> DeleteDatabaseAsync(string dbName)
{
 var response = await _amazonGlue.DeleteDatabaseAsync(new
DeleteDatabaseRequest { Name = dbName });
 return response.HttpStatusCode == HttpStatusCode.OK;
}

/// <summary>
/// Delete an AWS Glue job.
/// </summary>
/// <param name="jobName">The name of the job.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> DeleteJobAsync(string jobName)
{
 var response = await _amazonGlue.DeleteJobAsync(new DeleteJobRequest
{ JobName = jobName });
 return response.HttpStatusCode == HttpStatusCode.OK;
}

/// <summary>
/// Delete a table from an AWS Glue database.
/// </summary>
/// <param name="tableName">The table to delete.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> DeleteTableAsync(string dbName, string tableName)
{
 var response = await _amazonGlue.DeleteTableAsync(new DeleteTableRequest
{ Name = tableName, DatabaseName = dbName });
}
```

```
 return response.HttpStatusCode == HttpStatusCode.OK;
 }

 /// <summary>
 /// Get information about an AWS Glue crawler.
 /// </summary>
 /// <param name="crawlerName">The name of the crawler.</param>
 /// <returns>A Crawler object describing the crawler.</returns>
 public async Task<Crawler?> GetCrawlerAsync(string crawlerName)
 {
 var crawlerRequest = new GetCrawlerRequest
 {
 Name = crawlerName,
 };

 var response = await _amazonGlue.GetCrawlerAsync(crawlerRequest);
 if (response.HttpStatusCode == System.Net.HttpStatusCode.OK)
 {
 var databaseName = response.Crawler.DatabaseName;
 Console.WriteLine($"{crawlerName} has the database {databaseName}");
 return response.Crawler;
 }

 Console.WriteLine($"No information regarding {crawlerName} could be
found.");
 return null;
 }

 /// <summary>
 /// Get information about the state of an AWS Glue crawler.
 /// </summary>
 /// <param name="crawlerName">The name of the crawler.</param>
 /// <returns>A value describing the state of the crawler.</returns>
 public async Task<CrawlerState> GetCrawlerStateAsync(string crawlerName)
 {
 var response = await _amazonGlue.GetCrawlerAsync(
 new GetCrawlerRequest { Name = crawlerName });
 return response.Crawler.State;
 }

 /// <summary>
```

```
/// Get information about an AWS Glue database.
/// </summary>
/// <param name="dbName">The name of the database.</param>
/// <returns>A Database object containing information about the database.</
returns>
public async Task<Database> GetDatabaseAsync(string dbName)
{
 var databasesRequest = new GetDatabaseRequest
 {
 Name = dbName,
 };

 var response = await _amazonGlue.GetDatabaseAsync(databasesRequest);
 return response.Database;
}

/// <summary>
/// Get information about a specific AWS Glue job run.
/// </summary>
/// <param name="jobName">The name of the job.</param>
/// <param name="jobRunId">The Id of the job run.</param>
/// <returns>A JobRun object with information about the job run.</returns>
public async Task<JobRun> GetJobRunAsync(string jobName, string jobRunId)
{
 var response = await _amazonGlue.GetJobRunAsync(new GetJobRunRequest
{ JobName = jobName, RunId = jobRunId });
 return response.JobRun;
}

/// <summary>
/// Get information about all AWS Glue runs of a specific job.
/// </summary>
/// <param name="jobName">The name of the job.</param>
/// <returns>A list of JobRun objects.</returns>
public async Task<List<JobRun>> GetJobRunsAsync(string jobName)
{
 var jobRuns = new List<JobRun>();

 var request = new GetJobRunsRequest
 {
 JobName = jobName,
 };
};
```

```
// No need to loop to get all the log groups--the SDK does it for us
behind the scenes
var paginatorForJobRuns =
 _amazonGlue.Paginators.GetJobRuns(request);

await foreach (var response in paginatorForJobRuns.Responses)
{
 response.JobRuns.ForEach(jobRun =>
 {
 jobRuns.Add(jobRun);
 });
}

return jobRuns;
}

/// <summary>
/// Get a list of tables for an AWS Glue database.
/// </summary>
/// <param name="dbName">The name of the database.</param>
/// <returns>A list of Table objects.</returns>
public async Task<List<Table>> GetTablesAsync(string dbName)
{
 var request = new GetTablesRequest { DatabaseName = dbName };
 var tables = new List<Table>();

 // Get a paginator for listing the tables.
 var tablePaginator = _amazonGlue.Paginators.GetTables(request);

 await foreach (var response in tablePaginator.Responses)
 {
 tables.AddRange(response.TableList);
 }

 return tables;
}

/// <summary>
/// List AWS Glue jobs using a paginator.
/// </summary>
/// <returns>A list of AWS Glue job names.</returns>
```

```
public async Task<List<string>> ListJobsAsync()
{
 var jobNames = new List<string>();

 var listJobsPaginator = _amazonGlue.Paginators.ListJobs(new
ListJobsRequest { MaxResults = 10 });
 await foreach (var response in listJobsPaginator.Responses)
 {
 jobNames.AddRange(response.JobNames);
 }

 return jobNames;
}

/// <summary>
/// Start an AWS Glue crawler.
/// </summary>
/// <param name="crawlerName">The name of the crawler.</param>
/// <returns>A Boolean value indicating the success of the action.</returns>
public async Task<bool> StartCrawlerAsync(string crawlerName)
{
 var crawlerRequest = new StartCrawlerRequest
 {
 Name = crawlerName,
 };

 var response = await _amazonGlue.StartCrawlerAsync(crawlerRequest);

 return response.HttpStatusCode == System.Net.HttpStatusCode.OK;
}

/// <summary>
/// Start an AWS Glue job run.
/// </summary>
/// <param name="jobName">The name of the job.</param>
/// <returns>A string representing the job run Id.</returns>
public async Task<string> StartJobRunAsync(
 string jobName,
 string inputDatabase,
 string inputTable,
 string bucketName)
{
```

```
var request = new StartJobRunRequest
{
 JobName = jobName,
 Arguments = new Dictionary<string, string>
 {
 {"--input_database", inputDatabase},
 {"--input_table", inputTable},
 {"--output_bucket_url", $"s3://{bucketName}/"}
 }
};

var response = await _amazonGlue.StartJobRunAsync(request);
return response.JobRunId;
}
}
```

シナリオを実行するクラスを作成します。

```
global using Amazon.Glue;
global using GlueActions;
global using Microsoft.Extensions.Configuration;
global using Microsoft.Extensions.DependencyInjection;
global using Microsoft.Extensions.Hosting;
global using Microsoft.Extensions.Logging;
global using Microsoft.Extensions.Logging.Console;
global using Microsoft.Extensions.Logging.Debug;

using Amazon.Glue.Model;
using Amazon.S3;
using Amazon.S3.Model;

namespace GlueBasics;

public class GlueBasics
{
 private static ILogger logger = null!;
 private static IConfiguration _configuration = null!;
```

```
static async Task Main(string[] args)
{
 // Set up dependency injection for AWS Glue.
 using var host = Host.CreateDefaultBuilder(args)
 .ConfigureLogging(logging =>
 logging.AddFilter("System", LogLevel.Debug)
 .AddFilter<DebugLoggerProvider>("Microsoft",
LogLevel.Information)
 .AddFilter<ConsoleLoggerProvider>("Microsoft",
LogLevel.Trace))
 .ConfigureServices((_, services) =>
 services.AddAWSService<IAmazonGlue>()
 .AddTransient<GlueWrapper>()
 .AddTransient<UiWrapper>()
)
 .Build();

 logger = LoggerFactory.Create(builder => { builder.AddConsole(); })
 .CreateLogger<GlueBasics>();

 _configuration = new ConfigurationBuilder()
 .SetBasePath(Directory.GetCurrentDirectory())
 .AddJsonFile("settings.json") // Load settings from .json file.
 .AddJsonFile("settings.local.json",
 true) // Optionally load local settings.
 .Build();

 // These values are stored in settings.json
 // Once you have run the CDK script to deploy the resources,
 // edit the file to set "BucketName", "RoleName", and "ScriptURL"
 // to the appropriate values. Also set "CrawlerName" to the name
 // you want to give the crawler when it is created.
 string bucketName = _configuration["BucketName"]!;
 string bucketUrl = _configuration["BucketUrl"]!;
 string crawlerName = _configuration["CrawlerName"]!;
 string roleName = _configuration["RoleName"]!;
 string sourceData = _configuration["SourceData"]!;
 string dbName = _configuration["DbName"]!;
 string cron = _configuration["Cron"]!;
 string scriptUrl = _configuration["ScriptURL"]!;
 string jobName = _configuration["JobName"]!;

 var wrapper = host.Services.GetRequiredService<GlueWrapper>();
 var uiWrapper = host.Services.GetRequiredService<UiWrapper>();
}
```

```
uiWrapper.DisplayOverview();
uiWrapper.PressEnter();

// Create the crawler and wait for it to be ready.
uiWrapper.DisplayTitle("Create AWS Glue crawler");
Console.WriteLine("Let's begin by creating the AWS Glue crawler.");

var crawlerDescription = "Crawler created for the AWS Glue Basics
scenario.";
var crawlerCreated = await wrapper.CreateCrawlerAsync(crawlerName,
crawlerDescription, roleName, cron, sourceData, dbName);
if (crawlerCreated)
{
 Console.WriteLine($"The crawler: {crawlerName} has been created. Now
let's wait until it's ready.");
 CrawlerState crawlerState;
 do
 {
 crawlerState = await wrapper.GetCrawlerStateAsync(crawlerName);
 }
 while (crawlerState != "READY");
 Console.WriteLine($"The crawler {crawlerName} is now ready for
use.");
}
else
{
 Console.WriteLine($"Couldn't create crawler {crawlerName}.");
 return; // Exit the application.
}

uiWrapper.DisplayTitle("Start AWS Glue crawler");
Console.WriteLine("Now let's wait until the crawler has successfully
started.");
var crawlerStarted = await wrapper.StartCrawlerAsync(crawlerName);
if (crawlerStarted)
{
 CrawlerState crawlerState;
 do
 {
 crawlerState = await wrapper.GetCrawlerStateAsync(crawlerName);
 }
 while (crawlerState != "READY");
}
```

```
 Console.WriteLine($"The crawler {crawlerName} is now ready for
use.");
 }
 else
 {
 Console.WriteLine($"Couldn't start the crawler {crawlerName}.");
 return; // Exit the application.
 }

 uiWrapper.PressEnter();

 Console.WriteLine($"
Let's take a look at the database: {dbName}");
 var database = await wrapper.GetDatabaseAsync(dbName);

 if (database != null)
 {
 uiWrapper.DisplayTitle($"{database.Name} Details");
 Console.WriteLine($"{database.Name} created on
{database.CreateTime}");
 Console.WriteLine(database.Description);
 }

 uiWrapper.PressEnter();

 var tables = await wrapper.GetTablesAsync(dbName);
 if (tables.Count > 0)
 {
 tables.ForEach(table =>
 {
 Console.WriteLine($"{table.Name}\tCreated:
[table.CreateTime]\tUpdated: {table.UpdateTime}");
 });
 }

 uiWrapper.PressEnter();

 uiWrapper.DisplayTitle("Create AWS Glue job");
 Console.WriteLine("Creating a new AWS Glue job.");
 var description = "An AWS Glue job created using the AWS SDK for .NET";
 await wrapper.CreateJobAsync(dbName, tables[0].Name, bucketUrl, jobName,
roleName, description, scriptUrl);

 uiWrapper.PressEnter();
```

```
 uiWrapper.DisplayTitle("Starting AWS Glue job");
 Console.WriteLine("Starting the new AWS Glue job...");
 var jobRunId = await wrapper.StartJobRunAsync(jobName, dbName,
tables[0].Name, bucketName);
 var jobRunComplete = false;
 var jobRun = new JobRun();
 do
 {
 jobRun = await wrapper.GetJobRunAsync(jobName, jobRunId);
 if (jobRun.JobRunState == "SUCCEEDED" || jobRun.JobRunState ==
"STOPPED" ||
 jobRun.JobRunState == "FAILED" || jobRun.JobRunState ==
"TIMEOUT")
 {
 jobRunComplete = true;
 }
 } while (!jobRunComplete);

 uiWrapper.DisplayTitle($"Data in {bucketName}");

 // Get the list of data stored in the S3 bucket.
 var s3Client = new AmazonS3Client();

 var response = await s3Client.ListObjectsAsync(new ListObjectsRequest
{ BucketName = bucketName });
 response.S3Objects.ForEach(s3Object =>
 {
 Console.WriteLine(s3Object.Key);
 });

 uiWrapper.DisplayTitle("AWS Glue jobs");
 var jobNames = await wrapper.ListJobsAsync();
 jobNames.ForEach(jobName =>
 {
 Console.WriteLine(jobName);
 });

 uiWrapper.PressEnter();

 uiWrapper.DisplayTitle("Get AWS Glue job run information");
 Console.WriteLine("Getting information about the AWS Glue job.");
 var jobRuns = await wrapper.GetJobRunsAsync(jobName);

 jobRuns.ForEach(jobRun =>
```

```
 {
 Console.WriteLine($"{jobRun.JobName}\t{jobRun.JobRunState}\t{jobRun.CompletedOn}");
 });

 uiWrapper.PressEnter();

 uiWrapper.DisplayTitle("Deleting resources");
 Console.WriteLine("Deleting the AWS Glue job used by the example.");
 await wrapper.DeleteJobAsync(jobName);

 Console.WriteLine("Deleting the tables from the database.");
 tables.ForEach(async table =>
 {
 await wrapper.DeleteTableAsync(dbName, table.Name);
 });

 Console.WriteLine("Deleting the database.");
 await wrapper.DeleteDatabaseAsync(dbName);

 Console.WriteLine("Deleting the AWS Glue crawler.");
 await wrapper.DeleteCrawlerAsync(crawlerName);

 Console.WriteLine("The AWS Glue scenario has completed.");
 uiWrapper.PressEnter();
}
}

namespace GlueBasics;

public class UiWrapper
{
 public readonly string SepBar = new string('-', Console.WindowWidth);

 /// <summary>
 /// Show information about the scenario.
 /// </summary>
 public void DisplayOverview()
 {
 Console.Clear();
 DisplayTitle("Amazon Glue: get started with crawlers and jobs");

 Console.WriteLine("This example application does the following:");
 }
}
```

```

 Console.WriteLine("\t 1. Create a crawler, pass it the IAM role and the
URL to the public S3 bucket that contains the source data");
 Console.WriteLine("\t 2. Start the crawler.");
 Console.WriteLine("\t 3. Get the database created by the crawler and the
tables in the database.");
 Console.WriteLine("\t 4. Create a job.");
 Console.WriteLine("\t 5. Start a job run.");
 Console.WriteLine("\t 6. Wait for the job run to complete.");
 Console.WriteLine("\t 7. Show the data stored in the bucket.");
 Console.WriteLine("\t 8. List jobs for the account.");
 Console.WriteLine("\t 9. Get job run details for the job that was run.");
 Console.WriteLine("\t10. Delete the demo job.");
 Console.WriteLine("\t11. Delete the database and tables created for the
demo.");
 Console.WriteLine("\t12. Delete the crawler.");
 }

 /// <summary>
 /// Display a message and wait until the user presses enter.
 /// </summary>
 public void PressEnter()
 {
 Console.WriteLine("\nPlease press <Enter> to continue. ");
 _ = Console.ReadLine();
 }

 /// <summary>
 /// Pad a string with spaces to center it on the console display.
 /// </summary>
 /// <param name="strToCenter">The string to center on the screen.</param>
 /// <returns>The string padded to make it center on the screen.</returns>
 public string CenterString(string strToCenter)
 {
 var padAmount = (Console.WindowWidth - strToCenter.Length) / 2;
 var leftPad = new string(' ', padAmount);
 return $"{leftPad}{strToCenter}";
 }

 /// <summary>
 /// Display a line of hyphens, the centered text of the title and another
 /// line of hyphens.
 /// </summary>
 /// <param name="strTitle">The string to be displayed.</param>
 public void DisplayTitle(string strTitle)

```

```
{
 Console.WriteLine(SepBar);
 Console.WriteLine(CenterString(strTitle));
 Console.WriteLine(SepBar);
}
}
```

- APIの詳細については、「AWS SDK for .NET API リファレンス」の以下のトピックを参照してください。
  - [CreateCrawler](#)
  - [CreateJob](#)
  - [DeleteCrawler](#)
  - [DeleteDatabase](#)
  - [DeleteJob](#)
  - [DeleteTable](#)
  - [GetCrawler](#)
  - [GetDatabase](#)
  - [GetDatabases](#)
  - [GetJob](#)
  - [GetJobRun](#)
  - [GetJobRuns](#)
  - [GetTables](#)
  - [ListJobs](#)
  - [StartCrawler](#)
  - [StartJobRun](#)

## C++

## SDK for C++

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
//! Scenario which demonstrates using AWS Glue to add a crawler and run a job.
/*!
 \\sa runGettingStartedWithGlueScenario()
 \\param bucketName: An S3 bucket created in the setup.
 \\param roleName: An AWS Identity and Access Management (IAM) role created in the
 setup.
 \\param clientConfig: AWS client configuration.
 \\return bool: Successful completion.
*/

bool AwsDoc::Glue::runGettingStartedWithGlueScenario(const Aws::String
&bucketName,
 const Aws::String &roleName,
 const
Aws::Client::ClientConfiguration &clientConfig) {
 Aws::Glue::GlueClient client(clientConfig);

 Aws::String roleArn;
 if (!getRoleArn(roleName, roleArn, clientConfig)) {
 std::cerr << "Error getting role ARN for role." << std::endl;
 return false;
 }

 // 1. Upload the job script to the S3 bucket.
 {
 std::cout << "Uploading the job script '"
 << AwsDoc::Glue::PYTHON_SCRIPT
 << "'." << std::endl;

 if (!AwsDoc::Glue::uploadFile(bucketName,
 AwsDoc::Glue::PYTHON_SCRIPT_PATH,
 AwsDoc::Glue::PYTHON_SCRIPT,
```

```
 clientConfig)) {
 std::cerr << "Error uploading the job file." << std::endl;
 return false;
 }
 }

 // 2. Create a crawler.
 {
 Aws::Glue::Model::S3Target s3Target;
 s3Target.SetPath("s3://crawler-public-us-east-1/flight/2016/csv");
 Aws::Glue::Model::CrawlerTargets crawlerTargets;
 crawlerTargets.AddS3Targets(s3Target);

 Aws::Glue::Model::CreateCrawlerRequest request;
 request.SetTargets(crawlerTargets);
 request.SetName(CRAWLER_NAME);
 request.SetDatabaseName(CRAWLER_DATABASE_NAME);
 request.SetTablePrefix(CRAWLER_DATABASE_PREFIX);
 request.SetRole(roleArn);

 Aws::Glue::Model::CreateCrawlerOutcome outcome =
client.CreateCrawler(request);

 if (outcome.IsSuccess()) {
 std::cout << "Successfully created the crawler." << std::endl;
 }
 else {
 std::cerr << "Error creating a crawler. " <<
outcome.GetError().GetMessage()
 << std::endl;
 deleteAssets("", CRAWLER_DATABASE_NAME, "", bucketName,
clientConfig);
 return false;
 }
 }

 // 3. Get a crawler.
 {
 Aws::Glue::Model::GetCrawlerRequest request;
 request.SetName(CRAWLER_NAME);

 Aws::Glue::Model::GetCrawlerOutcome outcome = client.GetCrawler(request);

 if (outcome.IsSuccess()) {
```

```
 Aws::Glue::Model::CrawlerState crawlerState =
outcome.GetResult().GetCrawler().GetState();
 std::cout << "Retrieved crawler with state " <<

Aws::Glue::Model::CrawlerStateMapper::GetNameForCrawlerState(
 crawlerState)
 << "." << std::endl;
 }
 else {
 std::cerr << "Error retrieving a crawler. "
 << outcome.GetError().GetMessage() << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, "", bucketName,
 clientConfig);
 return false;
 }
}

// 4. Start a crawler.
{
 Aws::Glue::Model::StartCrawlerRequest request;
 request.SetName(CRAWLER_NAME);

 Aws::Glue::Model::StartCrawlerOutcome outcome =
client.StartCrawler(request);

 if (outcome.IsSuccess() || (Aws::Glue::GlueErrors::CRAWLER_RUNNING ==
 outcome.GetError().GetErrorType())) {
 if (!outcome.IsSuccess()) {
 std::cout << "Crawler was already started." << std::endl;
 }
 else {
 std::cout << "Successfully started crawler." << std::endl;
 }

 std::cout << "This may take a while to run." << std::endl;

 Aws::Glue::Model::CrawlerState crawlerState =
Aws::Glue::Model::CrawlerState::NOT_SET;
 int iterations = 0;
 while (Aws::Glue::Model::CrawlerState::READY != crawlerState) {
 std::this_thread::sleep_for(std::chrono::seconds(1));
 ++iterations;
 if ((iterations % 10) == 0) { // Log status every 10 seconds.
```

```
 std::cout << "Crawler status " <<

Aws::Glue::Model::CrawlerStateMapper::GetNameForCrawlerState(
 crawlerState)
 << ". After " << iterations
 << " seconds elapsed."
 << std::endl;
 }
 Aws::Glue::Model::GetCrawlerRequest getCrawlerRequest;
 getCrawlerRequest.SetName(CRAWLER_NAME);

 Aws::Glue::Model::GetCrawlerOutcome getCrawlerOutcome =
client.GetCrawler(
 getCrawlerRequest);

 if (getCrawlerOutcome.IsSuccess()) {
 crawlerState =
getCrawlerOutcome.GetResult().GetCrawler().GetState();
 }
 else {
 std::cerr << "Error getting crawler. "
 << getCrawlerOutcome.GetError().GetMessage() <<
std::endl;
 break;
 }
}

if (Aws::Glue::Model::CrawlerState::READY == crawlerState) {
 std::cout << "Crawler finished running after " << iterations
 << " seconds."
 << std::endl;
}
}
else {
 std::cerr << "Error starting a crawler. "
 << outcome.GetError().GetMessage()
 << std::endl;

 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, "", bucketName,
 clientConfig);
 return false;
}
}
```

```
// 5. Get a database.
{
 Aws::Glue::Model::GetDatabaseRequest request;
 request.SetName(CRAWLER_DATABASE_NAME);

 Aws::Glue::Model::GetDatabaseOutcome outcome =
client.GetDatabase(request);

 if (outcome.IsSuccess()) {
 const Aws::Glue::Model::Database &database =
outcome.GetResult().GetDatabase();

 std::cout << "Successfully retrieve the database\n" <<
 database.Jsonize().View().WriteReadable() << ". " <<
std::endl;
 }
 else {
 std::cerr << "Error getting the database. "
 << outcome.GetError().GetMessage() << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, "", bucketName,
 clientConfig);
 return false;
 }
}

// 6. Get tables.
Aws::String tableName;
{
 Aws::Glue::Model::GetTablesRequest request;
 request.SetDatabaseName(CRAWLER_DATABASE_NAME);
 std::vector<Aws::Glue::Model::Table> all_tables;
 Aws::String nextToken; // Used for pagination.
 do {
 Aws::Glue::Model::GetTablesOutcome outcome =
client.GetTables(request);

 if (outcome.IsSuccess()) {
 const std::vector<Aws::Glue::Model::Table> &tables =
outcome.GetResult().GetTableList();
 all_tables.insert(all_tables.end(), tables.begin(),
tables.end());
 nextToken = outcome.GetResult().GetNextToken();
 }
 else {
```

```
 std::cerr << "Error getting the tables. "
 << outcome.GetError().GetMessage()
 << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, "", bucketName,
 clientConfig);
 return false;
 }
} while (!nextToken.empty());

std::cout << "The database contains " << all_tables.size()
 << (all_tables.size() == 1 ?
 " table." : "tables.") << std::endl;
std::cout << "Here is a list of the tables in the database.";
for (size_t index = 0; index < all_tables.size(); ++index) {
 std::cout << " " << index + 1 << ": " <<
all_tables[index].GetName()
 << std::endl;
}

if (!all_tables.empty()) {
 int tableIndex = askQuestionForIntRange(
 "Enter an index to display the database detail ",
 1, static_cast<int>(all_tables.size()));
 std::cout << all_tables[tableIndex -
1].Jsonize().View().WriteReadable()
 << std::endl;

 tableName = all_tables[tableIndex - 1].GetName();
}
}

// 7. Create a job.
{
 Aws::Glue::Model::CreateJobRequest request;
 request.SetName(JOB_NAME);
 request.SetRole(roleArn);
 request.SetGlueVersion(GLUE_VERSION);

 Aws::Glue::Model::JobCommand command;
 command.SetName(JOB_COMMAND_NAME);
 command.SetPythonVersion(JOB_PYTHON_VERSION);
 command.SetScriptLocation(
 Aws::String("s3://") + bucketName + "/" + PYTHON_SCRIPT);
 request.SetCommand(command);
}
```

```
Aws::Glue::Model::CreateJobOutcome outcome = client.CreateJob(request);

if (outcome.IsSuccess()) {
 std::cout << "Successfully created the job." << std::endl;
}
else {
 std::cerr << "Error creating the job. " <<
outcome.GetError().GetMessage()
 << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, "", bucketName,
 clientConfig);
 return false;
}
}

// 8. Start a job run.
{
 Aws::Glue::Model::StartJobRunRequest request;
 request.SetJobName(JOB_NAME);

 Aws::Map<Aws::String, Aws::String> arguments;
 arguments["--input_database"] = CRAWLER_DATABASE_NAME;
 arguments["--input_table"] = tableName;
 arguments["--output_bucket_url"] = Aws::String("s3://") + bucketName +
"/";
 request.SetArguments(arguments);

 Aws::Glue::Model::StartJobRunOutcome outcome =
client.StartJobRun(request);

 if (outcome.IsSuccess()) {
 std::cout << "Successfully started the job." << std::endl;

 Aws::String jobRunId = outcome.GetResult().GetJobRunId();

 int iterator = 0;
 bool done = false;
 while (!done) {
 ++iterator;
 std::this_thread::sleep_for(std::chrono::seconds(1));
 Aws::Glue::Model::GetJobRunRequest jobRunRequest;
 jobRunRequest.SetJobName(JOB_NAME);
 jobRunRequest.SetRunId(jobRunId);
```

```

 Aws::Glue::Model::GetJobRunOutcome jobRunOutcome =
client.GetJobRun(
 jobRunRequest);

 if (jobRunOutcome.IsSuccess()) {
 const Aws::Glue::Model::JobRun &jobRun =
jobRunOutcome.GetResult().GetJobRun();
 Aws::Glue::Model::JobRunState jobRunState =
jobRun.GetJobRunState();

 if ((jobRunState == Aws::Glue::Model::JobRunState::STOPPED)
||
 (jobRunState == Aws::Glue::Model::JobRunState::FAILED) ||
 (jobRunState == Aws::Glue::Model::JobRunState::TIMEOUT))
{
 std::cerr << "Error running job. "
 << jobRun.GetErrorMessage()
 << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME,
JOB_NAME,
 bucketName,
 clientConfig);
 return false;
 }
 else if (jobRunState ==
 Aws::Glue::Model::JobRunState::SUCCEEDED) {
 std::cout << "Job run succeeded after " << iterator <<
 " seconds elapsed." << std::endl;
 done = true;
 }
 else if ((iterator % 10) == 0) { // Log status every 10
seconds.
 std::cout << "Job run status " <<

 Aws::Glue::Model::JobRunStateMapper::GetNameForJobRunState(
 jobRunState) <<
 ". " << iterator <<
 " seconds elapsed." << std::endl;
 }
 }
 else {
 std::cerr << "Error retrieving job run state. "
 << jobRunOutcome.GetError().GetMessage()

```

```
 << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, JOB_NAME,
 bucketName, clientConfig);
 return false;
 }
}
}
else {
 std::cerr << "Error starting a job. " <<
outcome.GetError().GetMessage()
 << std::endl;
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, JOB_NAME,
bucketName,
 clientConfig);
 return false;
}
}

// 9. List the output data stored in the S3 bucket.
{
 Aws::S3::S3Client s3Client;
 Aws::S3::Model::ListObjectsV2Request request;
 request.SetBucket(bucketName);
 request.SetPrefix(OUTPUT_FILE_PREFIX);

 Aws::String continuationToken; // Used for pagination.
 std::vector<Aws::S3::Model::Object> allObjects;
 do {
 if (!continuationToken.empty()) {
 request.SetContinuationToken(continuationToken);
 }
 Aws::S3::Model::ListObjectsV2Outcome outcome =
s3Client.ListObjectsV2(
 request);

 if (outcome.IsSuccess()) {
 const std::vector<Aws::S3::Model::Object> &objects =
 outcome.GetResult().GetContents();
 allObjects.insert(allObjects.end(), objects.begin(),
objects.end());
 continuationToken =
outcome.GetResult().GetNextContinuationToken();
 }
 else {
```

```
 std::cerr << "Error listing objects. "
 << outcome.GetError().GetMessage()
 << std::endl;
 break;
 }
} while (!continuationToken.empty());

std::cout << "Data from your job is in " << allObjects.size() <<
 " files in the S3 bucket, " << bucketName << "." << std::endl;

for (size_t i = 0; i < allObjects.size(); ++i) {
 std::cout << " " << i + 1 << ". " << allObjects[i].GetKey()
 << std::endl;
}

int objectIndex = askQuestionForIntRange(
 std::string(
 "Enter the number of a block to download it and see the
first ") +
 std::to_string(LINES_OF_RUN_FILE_TO_DISPLAY) +
 " lines of JSON output in the block: ", 1,
 static_cast<int>(allObjects.size()));

Aws::String objectKey = allObjects[objectIndex - 1].GetKey();

std::stringstream stringStream;
if (getObjectFromBucket(bucketName, objectKey, stringStream,
 clientConfig)) {
 for (int i = 0; i < LINES_OF_RUN_FILE_TO_DISPLAY && stringStream; +
+i) {
 std::string line;
 std::getline(stringStream, line);
 std::cout << " " << line << std::endl;
 }
}
else {
 deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, JOB_NAME,
bucketName,
 clientConfig);
 return false;
}
}

// 10. List all the jobs.
```

```
Aws::String jobName;
{
 Aws::Glue::Model::ListJobsRequest listJobsRequest;

 Aws::String nextToken;
 std::vector<Aws::String> allJobNames;

 do {
 if (!nextToken.empty()) {
 listJobsRequest.SetNextToken(nextToken);
 }
 Aws::Glue::Model::ListJobsOutcome listRunsOutcome = client.ListJobs(
 listJobsRequest);

 if (listRunsOutcome.IsSuccess()) {
 const std::vector<Aws::String> &jobNames =
listRunsOutcome.GetResult().GetJobNames();
 allJobNames.insert(allJobNames.end(), jobNames.begin(),
jobNames.end());
 nextToken = listRunsOutcome.GetResult().GetNextToken();
 }
 else {
 std::cerr << "Error listing jobs. "
 << listRunsOutcome.GetError().GetMessage()
 << std::endl;
 }
 } while (!nextToken.empty());
 std::cout << "Your account has " << allJobNames.size() << " jobs."
 << std::endl;
 for (size_t i = 0; i < allJobNames.size(); ++i) {
 std::cout << " " << i + 1 << ". " << allJobNames[i] << std::endl;
 }
 int jobIndex = askQuestionForIntRange(
 Aws::String("Enter a number between 1 and ") +
 std::to_string(allJobNames.size()) +
 " to see the list of runs for a job: ",
 1, static_cast<int>(allJobNames.size()));

 jobName = allJobNames[jobIndex - 1];
}

// 11. Get the job runs for a job.
Aws::String jobRunID;
if (!jobName.empty()) {
```

```
Aws::Glue::Model::GetJobRunsRequest getJobRunsRequest;
getJobRunsRequest.SetJobName(jobName);

Aws::String nextToken; // Used for pagination.
std::vector<Aws::Glue::Model::JobRun> allJobRuns;
do {
 if (!nextToken.empty()) {
 getJobRunsRequest.SetNextToken(nextToken);
 }
 Aws::Glue::Model::GetJobRunsOutcome jobRunsOutcome =
client.GetJobRuns(
 getJobRunsRequest);

 if (jobRunsOutcome.IsSuccess()) {
 const std::vector<Aws::Glue::Model::JobRun> &jobRuns =
jobRunsOutcome.GetResult().GetJobRuns();
 allJobRuns.insert(allJobRuns.end(), jobRuns.begin(),
jobRuns.end());

 nextToken = jobRunsOutcome.GetResult().GetNextToken();
 }
 else {
 std::cerr << "Error getting job runs. "
 << jobRunsOutcome.GetError().GetMessage()
 << std::endl;
 break;
 }
} while (!nextToken.empty());

std::cout << "There are " << allJobRuns.size() << " runs in the job '"
 <<
 jobName << "'." << std::endl;

for (size_t i = 0; i < allJobRuns.size(); ++i) {
 std::cout << " " << i + 1 << ". " << allJobRuns[i].GetJobName()
 << std::endl;
}

int runIndex = askQuestionForIntRange(
 Aws::String("Enter a number between 1 and ") +
 std::to_string(allJobRuns.size()) +
 " to see details for a run: ",
 1, static_cast<int>(allJobRuns.size()));
jobRunID = allJobRuns[runIndex - 1].GetId();
```

```
 }

 // 12. Get a single job run.
 if (!jobRunID.empty()) {
 Aws::Glue::Model::GetJobRunRequest jobRunRequest;
 jobRunRequest.SetJobName(jobName);
 jobRunRequest.SetRunId(jobRunID);

 Aws::Glue::Model::GetJobRunOutcome jobRunOutcome = client.GetJobRun(
 jobRunRequest);

 if (jobRunOutcome.IsSuccess()) {
 std::cout << "Displaying the job run JSON description." << std::endl;
 std::cout
 <<
jobRunOutcome.GetResult().GetJobRun().Jsonize().View().WriteReadable()
 << std::endl;
 }
 else {
 std::cerr << "Error get a job run. "
 << jobRunOutcome.GetError().GetMessage()
 << std::endl;
 }
 }

 return deleteAssets(CRAWLER_NAME, CRAWLER_DATABASE_NAME, JOB_NAME,
 bucketName,
 clientConfig);
}

//! Cleanup routine to delete created assets.
/*!
 \\sa deleteAssets()
 \\param crawler: Name of an AWS Glue crawler.
 \\param database: The name of an AWS Glue database.
 \\param job: The name of an AWS Glue job.
 \\param bucketName: The name of an S3 bucket.
 \\param clientConfig: AWS client configuration.
 \\return bool: Successful completion.
*/
bool AwsDoc::Glue::deleteAssets(const Aws::String &crawler, const Aws::String
 &database,
 const Aws::String &job, const Aws::String
 &bucketName,
```

```
const Aws::Client::ClientConfiguration
&clientConfig) {
 const Aws::Glue::GlueClient client(clientConfig);
 bool result = true;

 // 13. Delete a job.
 if (!job.empty()) {
 Aws::Glue::Model::DeleteJobRequest request;
 request.SetJobName(job);

 Aws::Glue::Model::DeleteJobOutcome outcome = client.DeleteJob(request);

 if (outcome.IsSuccess()) {
 std::cout << "Successfully deleted the job." << std::endl;
 }
 else {
 std::cerr << "Error deleting the job. " <<
outcome.GetError().GetMessage()
 << std::endl;
 result = false;
 }
 }

 // 14. Delete a database.
 if (!database.empty()) {
 Aws::Glue::Model::DeleteDatabaseRequest request;
 request.SetName(database);

 Aws::Glue::Model::DeleteDatabaseOutcome outcome = client.DeleteDatabase(
 request);

 if (outcome.IsSuccess()) {
 std::cout << "Successfully deleted the database." << std::endl;
 }
 else {
 std::cerr << "Error deleting database. " <<
outcome.GetError().GetMessage()
 << std::endl;
 result = false;
 }
 }

 // 15. Delete a crawler.
```

```

 if (!crawler.empty()) {
 Aws::Glue::Model::DeleteCrawlerRequest request;
 request.SetName(crawler);

 Aws::Glue::Model::DeleteCrawlerOutcome outcome =
client.DeleteCrawler(request);

 if (outcome.IsSuccess()) {
 std::cout << "Successfully deleted the crawler." << std::endl;
 }
 else {
 std::cerr << "Error deleting the crawler. "
 << outcome.GetError().GetMessage() << std::endl;
 result = false;
 }
 }

 // 16. Delete the job script and run data from the S3 bucket.
 result &= AwsDoc::Glue::deleteAllObjectsInS3Bucket(bucketName,
clientConfig);

 return result;
}

//! Routine which uploads a file to an S3 bucket.
/*!
 \\sa uploadFile()
 \\param bucketName: An S3 bucket created in the setup.
 \\param filePath: The path of the file to upload.
 \\param fileName The name for the uploaded file.
 \\param clientConfig: AWS client configuration.
 \\return bool: Successful completion.
*/
bool
AwsDoc::Glue::uploadFile(const Aws::String &bucketName,
 const Aws::String &filePath,
 const Aws::String &fileName,
 const Aws::Client::ClientConfiguration &clientConfig) {
 Aws::S3::S3Client s3_client(clientConfig);

 Aws::S3::Model::PutObjectRequest request;
 request.SetBucket(bucketName);
 request.SetKey(fileName);

 std::shared_ptr<Aws::IOStream> inputData =

```

```

 Aws::MakeShared<Aws::FStream>("SampleAllocationTag",
 filePath.c_str(),
 std::ios_base::in |
std::ios_base::binary);

 if (!*inputData) {
 std::cerr << "Error unable to read file " << filePath << std::endl;
 return false;
 }

 request.SetBody(inputData);

 Aws::S3::Model::PutObjectOutcome outcome =
 s3_client.PutObject(request);

 if (!outcome.IsSuccess()) {
 std::cerr << "Error: PutObject: " <<
 outcome.GetError().GetMessage() << std::endl;
 }
 else {
 std::cout << "Added object '" << filePath << "' to bucket '"
 << bucketName << "'." << std::endl;
 }

 return outcome.IsSuccess();
}

//! Routine which deletes all objects in an S3 bucket.
/*!
 \sa deleteAllObjectsInS3Bucket()
 \param bucketName: The S3 bucket name.
 \param clientConfig: AWS client configuration.
 \return bool: Successful completion.
 */
bool AwsDoc::Glue::deleteAllObjectsInS3Bucket(const Aws::String &bucketName,
 const
 Aws::Client::ClientConfiguration &clientConfig) {
 Aws::S3::S3Client client(clientConfig);
 Aws::S3::Model::ListObjectsV2Request listObjectsRequest;
 listObjectsRequest.SetBucket(bucketName);

 Aws::String continuationToken; // Used for pagination.
 bool result = true;
 do {

```

```
 if (!continuationToken.empty()) {
 listObjectsRequest.SetContinuationToken(continuationToken);
 }

 Aws::S3::Model::ListObjectsV2Outcome listObjectsOutcome =
client.ListObjectsV2(
 listObjectsRequest);

 if (listObjectsOutcome.IsSuccess()) {
 const std::vector<Aws::S3::Model::Object> &objects =
listObjectsOutcome.GetResult().GetContents();
 if (!objects.empty()) {
 Aws::S3::Model::DeleteObjectsRequest deleteObjectsRequest;
 deleteObjectsRequest.SetBucket(bucketName);

 std::vector<Aws::S3::Model::ObjectIdentifier> objectIdentifiers;
 for (const Aws::S3::Model::Object &object: objects) {
 objectIdentifiers.push_back(
 Aws::S3::Model::ObjectIdentifier().WithKey(
 object.GetKey()));
 }
 Aws::S3::Model::Delete objectsDelete;
 objectsDelete.SetObjects(objectIdentifiers);
 objectsDelete.SetQuiet(true);
 deleteObjectsRequest.SetDelete(objectsDelete);

 Aws::S3::Model::DeleteObjectsOutcome deleteObjectsOutcome =
 client.DeleteObjects(deleteObjectsRequest);

 if (!deleteObjectsOutcome.IsSuccess()) {
 std::cerr << "Error deleting objects. " <<
 deleteObjectsOutcome.GetError().GetMessage() <<
std::endl;

 result = false;
 break;
 }
 else {
 std::cout << "Successfully deleted the objects." <<
std::endl;

 }
 }
 }
 else {
 std::cout << "No objects to delete in '" << bucketName << "'."

```

```

 << std::endl;
 }

 continuationToken =
listObjectsOutcome.GetResult().GetNextContinuationToken();
 }
 else {
 std::cerr << "Error listing objects. "
 << listObjectsOutcome.GetError().GetMessage() << std::endl;
 result = false;
 break;
 }
} while (!continuationToken.empty());

return result;
}

//! Routine which retrieves an object from an S3 bucket.
/*!
 \sa getObjectFromBucket()
 \param bucketName: The S3 bucket name.
 \param objectKey: The object's name.
 \param objectStream: A stream to receive the retrieved data.
 \param clientConfig: AWS client configuration.
 \return bool: Successful completion.
 */
bool AwsDoc::Glue::getObjectFromBucket(const Aws::String &bucketName,
 const Aws::String &objectKey,
 std::ostream &objectStream,
 const Aws::Client::ClientConfiguration
&clientConfig) {
 Aws::S3::S3Client client(clientConfig);
 Aws::S3::Model::GetObjectRequest request;
 request.SetBucket(bucketName);
 request.SetKey(objectKey);

 Aws::S3::Model::GetObjectOutcome outcome = client.GetObject(request);

 if (outcome.IsSuccess()) {
 std::cout << "Successfully retrieved '" << objectKey << "'." <<
std::endl;
 auto &body = outcome.GetResult().GetBody();
 objectStream << body.rdbuf();
 }
}

```

```
 }
 else {
 std::cerr << "Error retrieving object. " <<
outcome.GetError().GetMessage()
 << std::endl;
 }

 return outcome.IsSuccess();
}
```

- APIの詳細については、『AWS SDK for C++ API リファレンス』の以下のトピックを参照してください。
  - [CreateCrawler](#)
  - [CreateJob](#)
  - [DeleteCrawler](#)
  - [DeleteDatabase](#)
  - [DeleteJob](#)
  - [DeleteTable](#)
  - [GetCrawler](#)
  - [GetDatabase](#)
  - [GetDatabases](#)
  - [GetJob](#)
  - [GetJobRun](#)
  - [GetJobRuns](#)
  - [GetTables](#)
  - [ListJobs](#)
  - [StartCrawler](#)
  - [StartJobRun](#)

## Java

## SDK for Java 2.x

 Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
/**
 *
 * Before running this Java V2 code example, set up your development
 * environment, including your credentials.
 *
 * For more information, see the following documentation topic:
 *
 * https://docs.aws.amazon.com/sdk-for-java/latest/developer-guide/get-
 * started.html
 *
 * To set up the resources, see this documentation topic:
 *
 * https://docs.aws.amazon.com/glue/latest/ug/tutorial-add-crawler.html
 *
 * This example performs the following tasks:
 *
 * 1. Create a database.
 * 2. Create a crawler.
 * 3. Get a crawler.
 * 4. Start a crawler.
 * 5. Get a database.
 * 6. Get tables.
 * 7. Create a job.
 * 8. Start a job run.
 * 9. List all jobs.
 * 10. Get job runs.
 * 11. Delete a job.
 * 12. Delete a database.
 * 13. Delete a crawler.
 */

public class GlueScenario {
```

```
public static final String DASHES = new String(new char[80]).replace("\0",
"-");

public static void main(String[] args) throws InterruptedException {
 final String usage = ""

 Usage:
 <iam> <s3Path> <cron> <dbName> <crawlerName> <jobName>\s

 Where:
 iam - The ARN of the IAM role that has AWS Glue and S3
permissions.\s
 s3Path - The Amazon Simple Storage Service (Amazon S3) target
that contains data (for example, CSV data).
 cron - A cron expression used to specify the schedule (i.e.,
cron(15 12 * * ? *).
 dbName - The database name.\s
 crawlerName - The name of the crawler.\s
 jobName - The name you assign to this job definition.
 scriptLocation - The Amazon S3 path to a script that runs a
job.
 locationUri - The location of the database
 bucketNameSc - The Amazon S3 bucket name used when creating a
job

 """;

 if (args.length != 9) {
 System.out.println(usage);
 System.exit(1);
 }

 String iam = args[0];
 String s3Path = args[1];
 String cron = args[2];
 String dbName = args[3];
 String crawlerName = args[4];
 String jobName = args[5];
 String scriptLocation = args[6];
 String locationUri = args[7];
 String bucketNameSc = args[8];

 Region region = Region.US_EAST_1;
 GlueClient glueClient = GlueClient.builder()
 .region(region)
```

```
 .build());
System.out.println(DASHES);
System.out.println("Welcome to the AWS Glue scenario.");
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("1. Create a database.");
createDatabase(glueClient, dbName, locationUri);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("2. Create a crawler.");
createGlueCrawler(glueClient, iam, s3Path, cron, dbName, crawlerName);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("3. Get a crawler.");
getSpecificCrawler(glueClient, crawlerName);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("4. Start a crawler.");
startSpecificCrawler(glueClient, crawlerName);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("5. Get a database.");
getSpecificDatabase(glueClient, dbName);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("**** Wait 5 min for the tables to become available");
TimeUnit.MINUTES.sleep(5);
System.out.println("6. Get tables.");
String myTableName = getGlueTables(glueClient, dbName);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("7. Create a job.");
createJob(glueClient, jobName, iam, scriptLocation);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("8. Start a Job run.");
```

```
startJob(glueClient, jobName, dbName, myTableName, bucketNameSc);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("9. List all jobs.");
getAllJobs(glueClient);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("10. Get job runs.");
getJobRuns(glueClient, jobName);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("11. Delete a job.");
deleteJob(glueClient, jobName);
System.out.println("*** Wait 5 MIN for the " + crawlerName + " to stop");
TimeUnit.MINUTES.sleep(5);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("12. Delete a database.");
deleteDatabase(glueClient, dbName);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("Delete a crawler.");
deleteSpecificCrawler(glueClient, crawlerName);
System.out.println(DASHES);

System.out.println(DASHES);
System.out.println("Successfully completed the AWS Glue Scenario");
System.out.println(DASHES);
}

public static void createDatabase(GlueClient glueClient, String dbName,
String locationUri) {
 try {
 DatabaseInput input = DatabaseInput.builder()
 .description("Built with the AWS SDK for Java V2")
 .name(dbName)
 .locationUri(locationUri)
 .build();
```

```
 CreateDatabaseRequest request = CreateDatabaseRequest.builder()
 .databaseInput(input)
 .build();

 glueClient.createDatabase(request);
 System.out.println(dbName + " was successfully created");

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
}

public static void createGlueCrawler(GlueClient glueClient,
 String iam,
 String s3Path,
 String cron,
 String dbName,
 String crawlerName) {

 try {
 S3Target s3Target = S3Target.builder()
 .path(s3Path)
 .build();

 List<S3Target> targetList = new ArrayList<>();
 targetList.add(s3Target);
 CrawlerTargets targets = CrawlerTargets.builder()
 .s3Targets(targetList)
 .build();

 CreateCrawlerRequest crawlerRequest = CreateCrawlerRequest.builder()
 .databaseName(dbName)
 .name(crawlerName)
 .description("Created by the AWS Glue Java API")
 .targets(targets)
 .role(iam)
 .schedule(cron)
 .build();

 glueClient.createCrawler(crawlerRequest);
 System.out.println(crawlerName + " was successfully created");

 } catch (GlueException e) {
```

```
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
}

public static void getSpecificCrawler(GlueClient glueClient, String
crawlerName) {
 try {
 GetCrawlerRequest crawlerRequest = GetCrawlerRequest.builder()
 .name(crawlerName)
 .build();

 boolean ready = false;
 while (!ready) {
 GetCrawlerResponse response =
glueClient.getCrawler(crawlerRequest);
 String status = response.crawler().stateAsString();
 if (status.compareTo("READY") == 0) {
 ready = true;
 }
 Thread.sleep(3000);
 }

 System.out.println("The crawler is now ready");

 } catch (GlueException | InterruptedException e) {
 System.err.println(e.getMessage());
 System.exit(1);
 }
}

public static void startSpecificCrawler(GlueClient glueClient, String
crawlerName) {
 try {
 StartCrawlerRequest crawlerRequest = StartCrawlerRequest.builder()
 .name(crawlerName)
 .build();

 glueClient.startCrawler(crawlerRequest);
 System.out.println(crawlerName + " was successfully started!");

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
}
```

```
 }
 }

 public static void getSpecificDatabase(GlueClient glueClient, String
databaseName) {
 try {
 GetDatabaseRequest databasesRequest = GetDatabaseRequest.builder()
 .name(databaseName)
 .build();

 GetDatabaseResponse response =
glueClient.getDatabase(databasesRequest);
 Instant createDate = response.database().createTime();

 // Convert the Instant to readable date.
 DateTimeFormatter formatter =
DateTimeFormatter.ofLocalizedDateTime(FormatStyle.SHORT)
 .withLocale(Locale.US)
 .withZone(ZoneId.systemDefault());

 formatter.format(createDate);
 System.out.println("The create date of the database is " +
createDate);

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
 }

 public static String getGlueTables(GlueClient glueClient, String dbName) {
 String myTableName = "";
 try {
 GetTablesRequest tableRequest = GetTablesRequest.builder()
 .databaseName(dbName)
 .build();

 GetTablesResponse response = glueClient.getTables(tableRequest);
 List<Table> tables = response.tableList();
 if (tables.isEmpty()) {
 System.out.println("No tables were returned");
 } else {
 for (Table table : tables) {
 myTableName = table.name();
 }
 }
 }
 }
}
```

```
 System.out.println("Table name is: " + myTableName);
 }
}

} catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
}
return myTableName;
}

public static void startJob(GlueClient glueClient, String jobName, String
inputDatabase, String inputTable,
 String outBucket) {
 try {
 Map<String, String> myMap = new HashMap<>();
 myMap.put("--input_database", inputDatabase);
 myMap.put("--input_table", inputTable);
 myMap.put("--output_bucket_url", outBucket);

 StartJobRunRequest runRequest = StartJobRunRequest.builder()
 .workerType(WorkerType.G_1_X)
 .numberOfWorkers(10)
 .arguments(myMap)
 .jobName(jobName)
 .build();

 StartJobRunResponse response = glueClient.startJobRun(runRequest);
 System.out.println("The request Id of the job is " +
response.responseMetadata().requestId());

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
}

public static void createJob(GlueClient glueClient, String jobName, String
iam, String scriptLocation) {
 try {
 JobCommand command = JobCommand.builder()
 .pythonVersion("3")
 .name("glueetl")
 .scriptLocation(scriptLocation)
```

```
 .build();

 CreateJobRequest jobRequest = CreateJobRequest.builder()
 .description("A Job created by using the AWS SDK for Java
V2")
 .glueVersion("2.0")
 .workerType(WorkerType.G_1_X)
 .numberOfWorkers(10)
 .name(jobName)
 .role(iam)
 .command(command)
 .build();

 glueClient.createJob(jobRequest);
 System.out.println(jobName + " was successfully created.");

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
}

public static void getAllJobs(GlueClient glueClient) {
 try {
 GetJobsRequest jobsRequest = GetJobsRequest.builder()
 .maxResults(10)
 .build();

 GetJobsResponse jobsResponse = glueClient.getJobs(jobsRequest);
 List<Job> jobs = jobsResponse.jobs();
 for (Job job : jobs) {
 System.out.println("Job name is : " + job.name());
 System.out.println("The job worker type is : " +
job.workerType().name());
 }

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
}

public static void getJobRuns(GlueClient glueClient, String jobName) {
 try {
```

```
GetJobRunsRequest runsRequest = GetJobRunsRequest.builder()
 .jobName(jobName)
 .maxResults(20)
 .build();

boolean jobDone = false;
while (!jobDone) {
 GetJobRunsResponse response = glueClient.getJobRuns(runsRequest);
 List<JobRun> jobRuns = response.jobRuns();
 for (JobRun jobRun : jobRuns) {
 String jobState = jobRun.jobRunState().name();
 if (jobState.compareTo("SUCCEEDED") == 0) {
 System.out.println(jobName + " has succeeded");
 jobDone = true;

 } else if (jobState.compareTo("STOPPED") == 0) {
 System.out.println("Job run has stopped");
 jobDone = true;

 } else if (jobState.compareTo("FAILED") == 0) {
 System.out.println("Job run has failed");
 jobDone = true;

 } else if (jobState.compareTo("TIMEOUT") == 0) {
 System.out.println("Job run has timed out");
 jobDone = true;

 } else {
 System.out.println("*** Job run state is " +
jobRun.jobRunState().name());
 System.out.println("Job run Id is " + jobRun.id());
 System.out.println("The Glue version is " +
jobRun.glueVersion());
 }
 TimeUnit.SECONDS.sleep(5);
 }
}

} catch (GlueException | InterruptedException e) {
 System.err.println(e.getMessage());
 System.exit(1);
}
}
```

```
public static void deleteJob(GlueClient glueClient, String jobName) {
 try {
 DeleteJobRequest jobRequest = DeleteJobRequest.builder()
 .jobName(jobName)
 .build();

 glueClient.deleteJob(jobRequest);
 System.out.println(jobName + " was successfully deleted");

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
}

public static void deleteDatabase(GlueClient glueClient, String databaseName)
{
 try {
 DeleteDatabaseRequest request = DeleteDatabaseRequest.builder()
 .name(databaseName)
 .build();

 glueClient.deleteDatabase(request);
 System.out.println(databaseName + " was successfully deleted");

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 System.exit(1);
 }
}

public static void deleteSpecificCrawler(GlueClient glueClient, String
crawlerName) {
 try {
 DeleteCrawlerRequest deleteCrawlerRequest =
DeleteCrawlerRequest.builder()
 .name(crawlerName)
 .build();

 glueClient.deleteCrawler(deleteCrawlerRequest);
 System.out.println(crawlerName + " was deleted");

 } catch (GlueException e) {
 System.err.println(e.awsErrorDetails().errorMessage());
 }
}
```

```
 System.exit(1);
 }
}
}
```

- APIの詳細については、『AWS SDK for Java 2.x API リファレンス』の以下のトピックを参照してください。

- [CreateCrawler](#)
- [CreateJob](#)
- [DeleteCrawler](#)
- [DeleteDatabase](#)
- [DeleteJob](#)
- [DeleteTable](#)
- [GetCrawler](#)
- [GetDatabase](#)
- [GetDatabases](#)
- [GetJob](#)
- [GetJobRun](#)
- [GetJobRuns](#)
- [GetTables](#)
- [ListJobs](#)
- [StartCrawler](#)
- [StartJobRun](#)

## JavaScript

### SDK for JavaScript (v3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

パブリックの Amazon Simple Storage Service (Amazon S3) バケットをクローलし、検出した CSV 形式データを記述するメタデータデータベースを生成するクローラーを作成して実行します。

```
const createCrawler = (name, role, dbName, tablePrefix, s3TargetPath) => {
 const client = new GlueClient({});

 const command = new CreateCrawlerCommand({
 Name: name,
 Role: role,
 DatabaseName: dbName,
 TablePrefix: tablePrefix,
 Targets: {
 S3Targets: [{ Path: s3TargetPath }],
 },
 });

 return client.send(command);
};

const getCrawler = (name) => {
 const client = new GlueClient({});

 const command = new GetCrawlerCommand({
 Name: name,
 });

 return client.send(command);
};

const startCrawler = (name) => {
 const client = new GlueClient({});

 const command = new StartCrawlerCommand({
 Name: name,
 });

 return client.send(command);
};

const crawlerExists = async ({ getCrawler }, crawlerName) => {
 try {
 await getCrawler(crawlerName);
 }
}
```

```
 return true;
 } catch {
 return false;
 }
};

/**
 * @param {{ createCrawler: import('../../../actions/create-crawler.js').createCrawler}} actions
 */
const makeCreateCrawlerStep = (actions) => async (context) => {
 if (await crawlerExists(actions, process.env.CRAWLER_NAME)) {
 log("Crawler already exists. Skipping creation.");
 } else {
 await actions.createCrawler(
 process.env.CRAWLER_NAME,
 process.env.ROLE_NAME,
 process.env.DATABASE_NAME,
 process.env.TABLE_PREFIX,
 process.env.S3_TARGET_PATH,
);

 log("Crawler created successfully.", { type: "success" });
 }

 return { ...context };
};

/**
 * @param {(name: string) => Promise<import('@aws-sdk/client-glue').GetCrawlerCommandOutput>} getCrawler
 * @param {string} crawlerName
 */
const waitForCrawler = async (getCrawler, crawlerName) => {
 const waitTimeInSeconds = 30;
 const { Crawler } = await getCrawler(crawlerName);

 if (!Crawler) {
 throw new Error(`Crawler with name ${crawlerName} not found.`);
 }

 if (Crawler.State === "READY") {
 return;
 }
}
```

```
log(`Crawler is ${Crawler.State}. Waiting ${waitTimeInSeconds} seconds...`);
await wait(waitTimeInSeconds);
return waitForCrawler(getCrawler, crawlerName);
};

const makeStartCrawlerStep =
 ({ startCrawler, getCrawler }) =>
 async (context) => {
 log("Starting crawler.");
 await startCrawler(process.env.CRAWLER_NAME);
 log("Crawler started.", { type: "success" });

 log("Waiting for crawler to finish running. This can take a while.");
 await waitForCrawler(getCrawler, process.env.CRAWLER_NAME);
 log("Crawler ready.", { type: "success" });

 return { ...context };
 };
```

のデータベースとテーブルに関する情報を一覧表示します AWS Glue Data Catalog。

```
const getDatabase = (name) => {
 const client = new GlueClient({});

 const command = new GetDatabaseCommand({
 Name: name,
 });

 return client.send(command);
};

const getTables = (databaseName) => {
 const client = new GlueClient({});

 const command = new GetTablesCommand({
 DatabaseName: databaseName,
 });

 return client.send(command);
};
```

```

const makeGetDatabaseStep =
 ({ getDatabase }) =>
 async (context) => {
 const {
 Database: { Name },
 } = await getDatabase(process.env.DATABASE_NAME);
 log(`Database: ${Name}`);
 return { ...context };
 };

/**
 * @param {{ getTables: () => Promise<import('@aws-sdk/client-glue').GetTablesCommandOutput>}} config
 */
const makeGetTablesStep =
 ({ getTables }) =>
 async (context) => {
 const { TableList } = await getTables(process.env.DATABASE_NAME);
 log("Tables:");
 log(TableList.map((table) => ` • ${table.Name}\n`));
 return { ...context };
 };

```

ソース Amazon S3 バケットから CSV 形式データを抽出し、フィールドを削除して名前を変更することで変換し、JSON 形式の出力を別の Amazon S3 バケットにロードするジョブを作成して実行します。

```

const createJob = (name, role, scriptBucketName, scriptKey) => {
 const client = new GlueClient({});

 const command = new CreateJobCommand({
 Name: name,
 Role: role,
 Command: {
 Name: "glueetl",
 PythonVersion: "3",
 ScriptLocation: `s3://${scriptBucketName}/${scriptKey}`,
 },
 GlueVersion: "3.0",
 });

 return client.send(command);
}

```

```
};

const startJobRun = (jobName, dbName, tableName, bucketName) => {
 const client = new GlueClient({});

 const command = new StartJobRunCommand({
 JobName: jobName,
 Arguments: {
 "--input_database": dbName,
 "--input_table": tableName,
 "--output_bucket_url": `s3://${bucketName}/`,
 },
 });

 return client.send(command);
};

const makeCreateJobStep =
 ({ createJob }) =>
 async (context) => {
 log("Creating Job.");
 await createJob(
 process.env.JOB_NAME,
 process.env.ROLE_NAME,
 process.env.BUCKET_NAME,
 process.env.PYTHON_SCRIPT_KEY,
);
 log("Job created.", { type: "success" });

 return { ...context };
 };

/**
 * @param {(name: string, runId: string) => Promise<import('@aws-sdk/client-glue').GetJobRunCommandOutput> } getJobRun
 * @param {string} jobName
 * @param {string} jobRunId
 */
const waitForJobRun = async (getJobRun, jobName, jobRunId) => {
 const waitTimeInSeconds = 30;
 const { JobRun } = await getJobRun(jobName, jobRunId);

 if (!JobRun) {
 throw new Error(`Job run with id ${jobRunId} not found.`);
 }
};
```

```
}

switch (JobRun.JobRunState) {
 case "FAILED":
 case "TIMEOUT":
 case "STOPPED":
 throw new Error(
 `Job ${JobRun.JobRunState}. Error: ${JobRun.ErrorMessage}`,
);
 case "RUNNING":
 break;
 case "SUCCEEDED":
 return;
 default:
 throw new Error(`Unknown job run state: ${JobRun.JobRunState}`);
}

log(
 `Job ${JobRun.JobRunState}. Waiting ${waitTimeInSeconds} more seconds...`,
);
await wait(waitTimeInSeconds);
return waitForJobRun(getJobRun, jobName, jobRunId);
};

/**
 * @param {{ prompter: { prompt: () => Promise<{ shouldOpen: boolean }>} }}
 * context
 */
const promptToOpen = async (context) => {
 const { shouldOpen } = await context.prompter.prompt({
 name: "shouldOpen",
 type: "confirm",
 message: "Open the output bucket in your browser?",
 });

 if (shouldOpen) {
 return open(
 `https://s3.console.aws.amazon.com/s3/buckets/${process.env.BUCKET_NAME} to
 view the output.`
);
 }
};

const makeStartJobRunStep =
```

```
({ startJobRun, getJobRun }) =>
async (context) => {
 log("Starting job.");
 const { JobRunId } = await startJobRun(
 process.env.JOB_NAME,
 process.env.DATABASE_NAME,
 process.env.TABLE_NAME,
 process.env.BUCKET_NAME,
);
 log("Job started.", { type: "success" });

 log("Waiting for job to finish running. This can take a while.");
 await waitForJobRun(getJobRun, process.env.JOB_NAME, JobRunId);
 log("Job run succeeded.", { type: "success" });

 await promptToOpen(context);

 return { ...context };
};
```

ジョブ実行に関する情報を一覧表示し、変換されたデータの一部を表示します。

```
const getJobRuns = (jobName) => {
 const client = new GlueClient({});
 const command = new GetJobRunsCommand({
 JobName: jobName,
 });

 return client.send(command);
};

const getJobRun = (jobName, jobRunId) => {
 const client = new GlueClient({});
 const command = new GetJobRunCommand({
 JobName: jobName,
 RunId: jobRunId,
 });

 return client.send(command);
};

/**
```

```
* @typedef {{ prompter: { prompt: () => Promise<{jobName: string}> } }} Context
*/

/**
 * @typedef {() => Promise<import('@aws-sdk/client-
glue').GetJobRunCommandOutput>} getJobRun
 */

/**
 * @typedef {() => Promise<import('@aws-sdk/client-
glue').GetJobRunsCommandOutput>} getJobRuns
 */

/**
 *
 * @param {getJobRun} getJobRun
 * @param {string} jobName
 * @param {string} jobRunId
 */
const logJobRunDetails = async (getJobRun, jobName, jobRunId) => {
 const { JobRun } = await getJobRun(jobName, jobRunId);
 log(JobRun, { type: "object" });
};

/**
 *
 * @param {{getJobRuns: getJobRuns, getJobRun: getJobRun }} funcs
 */
const makePickJobRunStep =
 ({ getJobRuns, getJobRun }) =>
 async (** @type { Context } */ context) => {
 if (context.selectedJobName) {
 const { JobRuns } = await getJobRuns(context.selectedJobName);

 const { jobRunId } = await context.prompter.prompt({
 name: "jobRunId",
 type: "list",
 message: "Select a job run to see details.",
 choices: JobRuns.map((run) => run.Id),
 });

 logJobRunDetails(getJobRun, context.selectedJobName, jobRunId);
 }
 }
}
```

```
 return { ...context };
};
```

デモによって作成されたすべてのリソースを削除します。

```
const deleteJob = (jobName) => {
 const client = new GlueClient({});

 const command = new DeleteJobCommand({
 JobName: jobName,
 });

 return client.send(command);
};

const deleteTable = (databaseName, tableName) => {
 const client = new GlueClient({});

 const command = new DeleteTableCommand({
 DatabaseName: databaseName,
 Name: tableName,
 });

 return client.send(command);
};

const deleteDatabase = (databaseName) => {
 const client = new GlueClient({});

 const command = new DeleteDatabaseCommand({
 Name: databaseName,
 });

 return client.send(command);
};

const deleteCrawler = (crawlerName) => {
 const client = new GlueClient({});

 const command = new DeleteCrawlerCommand({
 Name: crawlerName,
 });
};
```

```
 return client.send(command);
 };

/**
 *
 * @param {import('.././../actions/delete-job.js').deleteJob} deleteJobFn
 * @param {string[]} jobNames
 * @param {{ prompter: { prompt: () => Promise<any> }}} context
 */
const handleDeleteJobs = async (deleteJobFn, jobNames, context) => {
 /**
 * @type {{ selectedJobNames: string[] }}
 */
 const { selectedJobNames } = await context.prompter.prompt({
 name: "selectedJobNames",
 type: "checkbox",
 message: "Let's clean up jobs. Select jobs to delete.",
 choices: jobNames,
 });

 if (selectedJobNames.length === 0) {
 log("No jobs selected.");
 } else {
 log("Deleting jobs.");
 await Promise.all(
 selectedJobNames.map((n) => deleteJobFn(n).catch(console.error)),
);
 log("Jobs deleted.", { type: "success" });
 }
};

/**
 * @param {{
 * listJobs: import('.././../actions/list-jobs.js').listJobs,
 * deleteJob: import('.././../actions/delete-job.js').deleteJob
 * }} config
 */
const makeCleanUpJobsStep =
 ({ listJobs, deleteJob }) =>
 async (context) => {
 const { JobNames } = await listJobs();
 if (JobNames.length > 0) {
 await handleDeleteJobs(deleteJob, JobNames, context);
 }
 };
};
```

```
 }

 return { ...context };
 };

/**
 * @param {import('.././../actions/delete-table.js').deleteTable} deleteTable
 * @param {string} databaseName
 * @param {string[]} tableNames
 */
const deleteTables = (deleteTable, databaseName, tableNames) =>
 Promise.all(
 tableNames.map((tableName) =>
 deleteTable(databaseName, tableName).catch(console.error),
),
);

/**
 * @param {{
 * getTables: import('.././../actions/get-tables.js').getTables,
 * deleteTable: import('.././../actions/delete-table.js').deleteTable
 * }} config
 */
const makeCleanUpTablesStep =
 ({ getTables, deleteTable }) =>
 /**
 * @param {{ prompter: { prompt: () => Promise<any>}}} context
 */
 async (context) => {
 const { TableList } = await getTables(process.env.DATABASE_NAME).catch(
 () => ({ TableList: null }),
);

 if (TableList && TableList.length > 0) {
 /**
 * @type {{ tableNames: string[] }}
 */
 const { tableNames } = await context.prompter.prompt({
 name: "tableNames",
 type: "checkbox",
 message: "Let's clean up tables. Select tables to delete.",
 choices: TableList.map((t) => t.Name),
 });
 }
 }
};
```

```
 if (tableNames.length === 0) {
 log("No tables selected.");
 } else {
 log("Deleting tables.");
 await deleteTables(deleteTable, process.env.DATABASE_NAME, tableNames);
 log("Tables deleted.", { type: "success" });
 }
 }
}

return { ...context };
};

/**
 * @param {import('.././././actions/delete-database.js').deleteDatabase}
deleteDatabase
 * @param {string[]} databaseNames
 */
const deleteDatabases = (deleteDatabase, databaseNames) =>
 Promise.all(
 databaseNames.map((dbName) => deleteDatabase(dbName).catch(console.error)),
);

/**
 * @param {{
 * getDatabases: import('.././././actions/get-databases.js').getDatabases
 * deleteDatabase: import('.././././actions/delete-database.js').deleteDatabase
 * }} config
 */
const makeCleanUpDatabasesStep =
 ({ getDatabases, deleteDatabase }) =>
 /**
 * @param {{ prompter: { prompt: () => Promise<any> } } } context
 */
 async (context) => {
 const { DatabaseList } = await getDatabases();

 if (DatabaseList.length > 0) {
 /** @type {{ dbName: string[] }} */
 const { dbName } = await context.prompter.prompt({
 name: "dbNames",
 type: "checkbox",
 message: "Let's clean up databases. Select databases to delete.",
 choices: DatabaseList.map((db) => db.Name),
 });
 }
 });
```

```
 if (dbNames.length === 0) {
 log("No databases selected.");
 } else {
 log("Deleting databases.");
 await deleteDatabases(deleteDatabase, dbNames);
 log("Databases deleted.", { type: "success" });
 }
 }

 return { ...context };
};

const cleanUpCrawlerStep = async (context) => {
 log(`Deleting crawler.`);

 try {
 await deleteCrawler(process.env.CRAWLER_NAME);
 log("Crawler deleted.", { type: "success" });
 } catch (err) {
 if (err.name === "EntityNotFoundException") {
 log(`Crawler is already deleted.`);
 } else {
 throw err;
 }
 }

 return { ...context };
};
```

- APIの詳細については、「AWS SDK for JavaScript API リファレンス」の以下のトピックを参照してください。
  - [CreateCrawler](#)
  - [CreateJob](#)
  - [DeleteCrawler](#)
  - [DeleteDatabase](#)
  - [DeleteJob](#)
  - [DeleteTable](#)
  - [GetCrawler](#)

- [GetDatabase](#)
- [GetDatabases](#)
- [GetJob](#)
- [GetJobRun](#)
- [GetJobRuns](#)
- [GetTables](#)
- [ListJobs](#)
- [StartCrawler](#)
- [StartJobRun](#)

## Kotlin

### SDK for Kotlin

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
suspend fun main(args: Array<String>) {
 val usage = ""
 Usage:
 <iam> <s3Path> <cron> <dbName> <crawlerName> <jobName>
<scriptLocation> <locationUri>

 Where:
 iam - The Amazon Resource Name (ARN) of the AWS Identity and Access
Management (IAM) role that has AWS Glue and Amazon Simple Storage Service
(Amazon S3) permissions.
 s3Path - The Amazon Simple Storage Service (Amazon S3) target that
contains data (for example, CSV data).
 cron - A cron expression used to specify the schedule (for example,
cron(15 12 * * ? *)).
 dbName - The database name.
 crawlerName - The name of the crawler.
 jobName - The name you assign to this job definition.
```

```
 scriptLocation - Specifies the Amazon S3 path to a script that runs a
job.
 locationUri - Specifies the location of the database
 """"

 if (args.size != 8) {
 println(usage)
 exitProcess(1)
 }

 val iam = args[0]
 val s3Path = args[1]
 val cron = args[2]
 val dbName = args[3]
 val crawlerName = args[4]
 val jobName = args[5]
 val scriptLocation = args[6]
 val locationUri = args[7]

 println("About to start the AWS Glue Scenario")
 createDatabase(dbName, locationUri)
 createCrawler(iam, s3Path, cron, dbName, crawlerName)
 getCrawler(crawlerName)
 startCrawler(crawlerName)
 getDatabase(dbName)
 getGlueTables(dbName)
 createJob(jobName, iam, scriptLocation)
 startJob(jobName)
 getJobs()
 getJobRuns(jobName)
 deleteJob(jobName)
 println("**** Wait for 5 MIN so the $crawlerName is ready to be deleted")
 TimeUnit.MINUTES.sleep(5)
 deleteMyDatabase(dbName)
 deleteCrawler(crawlerName)
}

suspend fun createDatabase(
 dbName: String?,
 locationUriVal: String?,
) {
 val input =
 DatabaseInput {
 description = "Built with the AWS SDK for Kotlin"
 }
}
```

```
 name = dbName
 locationUri = locationUriVal
 }

 val request =
 CreateDatabaseRequest {
 databaseInput = input
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 glueClient.createDatabase(request)
 println("The database was successfully created")
 }
}

suspend fun createCrawler(
 iam: String?,
 s3Path: String?,
 cron: String?,
 dbName: String?,
 crawlerName: String,
) {
 val s3Target =
 S3Target {
 path = s3Path
 }

 val targetList = ArrayList<S3Target>()
 targetList.add(s3Target)

 val target0b =
 CrawlerTargets {
 s3Targets = targetList
 }

 val crawlerRequest =
 CreateCrawlerRequest {
 databaseName = dbName
 name = crawlerName
 description = "Created by the AWS Glue Java API"
 targets = target0b
 role = iam
 schedule = cron
 }
}
```

```
GlueClient { region = "us-east-1" }.use { glueClient ->
 glueClient.createCrawler(crawlerRequest)
 println("$crawlerName was successfully created")
}
}

suspend fun getCrawler(crawlerName: String?) {
 val request =
 GetCrawlerRequest {
 name = crawlerName
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 val response = glueClient.getCrawler(request)
 val role = response.crawler?.role
 println("The role associated with this crawler is $role")
 }
}

suspend fun startCrawler(crawlerName: String) {
 val crawlerRequest =
 StartCrawlerRequest {
 name = crawlerName
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 glueClient.startCrawler(crawlerRequest)
 println("$crawlerName was successfully started.")
 }
}

suspend fun getDatabase(databaseName: String?) {
 val request =
 GetDatabaseRequest {
 name = databaseName
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 val response = glueClient.getDatabase(request)
 val dbDesc = response.database?.description
 println("The database description is $dbDesc")
 }
}
```

```
suspend fun getGlueTables(dbName: String?) {
 val tableRequest =
 GetTablesRequest {
 databaseName = dbName
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 val response = glueClient.getTables(tableRequest)
 response.tableList?.forEach { tableName ->
 println("Table name is ${tableName.name}")
 }
 }
}

suspend fun startJob(jobNameVal: String?) {
 val runRequest =
 StartJobRunRequest {
 workerType = WorkerType.G1X
 numberOfWorkers = 10
 jobName = jobNameVal
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 val response = glueClient.startJobRun(runRequest)
 println("The job run Id is ${response.jobRunId}")
 }
}

suspend fun createJob(
 jobName: String,
 iam: String?,
 scriptLocationVal: String?,
) {
 val commandOb =
 JobCommand {
 pythonVersion = "3"
 name = "MyJob1"
 scriptLocation = scriptLocationVal
 }

 val jobRequest =
 CreateJobRequest {
 description = "A Job created by using the AWS SDK for Java V2"
 }
}
```

```
 glueVersion = "2.0"
 workerType = WorkerType.G1X
 numberOfWorkers = 10
 name = jobName
 role = iam
 command = commandOb
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 glueClient.createJob(jobRequest)
 println("$jobName was successfully created.")
 }
}

suspend fun getJobs() {
 val request =
 GetJobsRequest {
 maxResults = 10
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 val response = glueClient.getJobs(request)
 response.jobs?.forEach { job ->
 println("Job name is ${job.name}")
 }
 }
}

suspend fun getJobRuns(jobNameVal: String?) {
 val request =
 GetJobRunsRequest {
 jobName = jobNameVal
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 val response = glueClient.getJobRuns(request)
 response.jobRuns?.forEach { job ->
 println("Job name is ${job.jobName}")
 }
 }
}

suspend fun deleteJob(jobNameVal: String) {
 val jobRequest =
```

```
 DeleteJobRequest {
 jobName = jobNameVal
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 glueClient.deleteJob(jobRequest)
 println("$jobNameVal was successfully deleted")
 }
 }

suspend fun deleteMyDatabase(databaseName: String) {
 val request =
 DeleteDatabaseRequest {
 name = databaseName
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 glueClient.deleteDatabase(request)
 println("$databaseName was successfully deleted")
 }
}

suspend fun deleteCrawler(crawlerName: String) {
 val request =
 DeleteCrawlerRequest {
 name = crawlerName
 }

 GlueClient { region = "us-east-1" }.use { glueClient ->
 glueClient.deleteCrawler(request)
 println("$crawlerName was deleted")
 }
}
```

- APIの詳細については、『AWS SDK for Kotlin API リファレンス』の以下のトピックを参照してください。
  - [CreateCrawler](#)
  - [CreateJob](#)
  - [DeleteCrawler](#)
  - [DeleteDatabase](#)
  - [DeleteJob](#)

- [DeleteTable](#)
- [GetCrawler](#)
- [GetDatabase](#)
- [GetDatabases](#)
- [GetJob](#)
- [GetJobRun](#)
- [GetJobRuns](#)
- [GetTables](#)
- [ListJobs](#)
- [StartCrawler](#)
- [StartJobRun](#)

## PHP

### SDK for PHP

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

```
namespace Glue;

use Aws\Glue\GlueClient;
use Aws\S3\S3Client;
use AwsUtilities\AWSServiceClass;
use GuzzleHttp\Psr7\Stream;
use Iam\IAMService;

class GettingStartedWithGlue
{
 public function run()
 {
 echo("\n");
 echo("-----\n");
 print("Welcome to the AWS Glue getting started demo using PHP!\n");
 }
}
```

```
echo("-----\n");

$clientArgs = [
 'region' => 'us-west-2',
 'version' => 'latest',
 'profile' => 'default',
];
$uniqid = uniqid();

$glueClient = new GlueClient($clientArgs);
$glueService = new GlueService($glueClient);
$iamService = new IAMService();
$crawlerName = "example-crawler-test-" . $uniqid;

AWSServiceClass::$waitTime = 5;
AWSServiceClass::$maxWaitAttempts = 20;

$role = $iamService->getRole("AWSGlueServiceRole-DocExample");

$databaseName = "doc-example-database-$uniqid";
$path = 's3://crawler-public-us-east-1/flight/2016/csv';
$glueService->createCrawler($crawlerName, $role['Role']['Arn'],
$databaseName, $path);
$glueService->startCrawler($crawlerName);

echo "Waiting for crawler";
do {
 $crawler = $glueService->getCrawler($crawlerName);
 echo ".";
 sleep(10);
} while ($crawler['Crawler']['State'] != "READY");
echo "\n";

$database = $glueService->getDatabase($databaseName);
echo "Found a database named " . $database['Database']['Name'] . "\n";

//Upload job script
$s3client = new S3Client($clientArgs);
$bucketName = "test-glue-bucket-" . $uniqid;
$s3client->createBucket([
 'Bucket' => $bucketName,
 'CreateBucketConfiguration' => ['LocationConstraint' => 'us-west-2'],
]);
```

```
$s3client->putObject([
 'Bucket' => $bucketName,
 'Key' => 'run_job.py',
 'SourceFile' => __DIR__ . '/flight_etl_job_script.py'
]);
$s3client->putObject([
 'Bucket' => $bucketName,
 'Key' => 'setup_scenario_getting_started.yaml',
 'SourceFile' => __DIR__ . '/setup_scenario_getting_started.yaml'
]);

$tables = $glueService->getTables($databaseName);

$jobName = 'test-job-' . $uniqid;
$scriptLocation = "s3://$bucketName/run_job.py";
$job = $glueService->createJob($jobName, $role['Role']['Arn'],
$scriptLocation);

$outputBucketUrl = "s3://$bucketName";
$runId = $glueService->startJobRun($jobName, $databaseName, $tables,
$outputBucketUrl)['JobRunId'];

echo "waiting for job";
do {
 $jobRun = $glueService->getJobRun($jobName, $runId);
 echo ".";
 sleep(10);
} while (!array_intersect([$jobRun['JobRun']['JobRunState']],
['SUCCEEDED', 'STOPPED', 'FAILED', 'TIMEOUT']));
echo "\n";

$jobRuns = $glueService->getJobRuns($jobName);

$objects = $s3client->listObjects([
 'Bucket' => $bucketName,
])['Contents'];

foreach ($objects as $object) {
 echo $object['Key'] . "\n";
}

echo "Downloading " . $objects[1]['Key'] . "\n";
/** @var Stream $downloadObject */
$downloadObject = $s3client->getObject([
```

```
 'Bucket' => $bucketName,
 'Key' => $objects[1]['Key'],
)['Body']->getContents();
echo "Here is the first 1000 characters in the object.";
echo substr($downloadObject, 0, 1000);

$jobs = $glueService->listJobs();
echo "Current jobs:\n";
foreach ($jobs['JobNames'] as $jobsName) {
 echo "{$jobsName}\n";
}

echo "Delete the job.\n";
$glueClient->deleteJob([
 'JobName' => $job['Name'],
]);

echo "Delete the tables.\n";
foreach ($tables['TableList'] as $table) {
 $glueService->deleteTable($table['Name'], $databaseName);
}

echo "Delete the databases.\n";
$glueClient->deleteDatabase([
 'Name' => $databaseName,
]);

echo "Delete the crawler.\n";
$glueClient->deleteCrawler([
 'Name' => $crawlerName,
]);

$deleteObjects = $s3client->listObjectsV2([
 'Bucket' => $bucketName,
]);
echo "Delete all objects in the bucket.\n";
$deleteObjects = $s3client->deleteObjects([
 'Bucket' => $bucketName,
 'Delete' => [
 'Objects' => $deleteObjects['Contents'],
]
]);
echo "Delete the bucket.\n";
$s3client->deleteBucket(['Bucket' => $bucketName]);
```

```
 echo "This job was brought to you by the number $uniqid\n";
 }
}

namespace Glue;

use Aws\Glue\GlueClient;
use Aws\Result;

use function PHPUnit\Framework\isEmpty;

class GlueService extends \AwsUtilities\AWSServiceClass
{
 protected GlueClient $glueClient;

 public function __construct($glueClient)
 {
 $this->glueClient = $glueClient;
 }

 public function getCrawler($crawlerName)
 {
 return $this->customWaiter(function () use ($crawlerName) {
 return $this->glueClient->getCrawler([
 'Name' => $crawlerName,
]);
 });
 }

 public function createCrawler($crawlerName, $role, $databaseName, $path):
 Result
 {
 return $this->customWaiter(function () use ($crawlerName, $role,
 $databaseName, $path) {
 return $this->glueClient->createCrawler([
 'Name' => $crawlerName,
 'Role' => $role,
 'DatabaseName' => $databaseName,
 'Targets' => [
 'S3Targets' =>
 [[
 'Path' => $path,
]]
]
]);
 });
 }
}
```

```
],
 });
});
}

public function startCrawler($crawlerName): Result
{
 return $this->glueClient->startCrawler([
 'Name' => $crawlerName,
]);
}

public function getDatabase(string $databaseName): Result
{
 return $this->customWaiter(function () use ($databaseName) {
 return $this->glueClient->getDatabase([
 'Name' => $databaseName,
]);
 });
}

public function getTables($databaseName): Result
{
 return $this->glueClient->getTables([
 'DatabaseName' => $databaseName,
]);
}

public function createJob($jobName, $role, $scriptLocation, $pythonVersion =
'3', $glueVersion = '3.0'): Result
{
 return $this->glueClient->createJob([
 'Name' => $jobName,
 'Role' => $role,
 'Command' => [
 'Name' => 'glueetl',
 'ScriptLocation' => $scriptLocation,
 'PythonVersion' => $pythonVersion,
],
 'GlueVersion' => $glueVersion,
]);
}
```

```
public function startJobRun($jobName, $databaseName, $tables,
$outputBucketUrl): Result
{
 return $this->glueClient->startJobRun([
 'JobName' => $jobName,
 'Arguments' => [
 'input_database' => $databaseName,
 'input_table' => $tables['TableList'][0]['Name'],
 'output_bucket_url' => $outputBucketUrl,
 '--input_database' => $databaseName,
 '--input_table' => $tables['TableList'][0]['Name'],
 '--output_bucket_url' => $outputBucketUrl,
],
]);
}

public function listJobs($maxResults = null, $nextToken = null, $tags = []):
Result
{
 $arguments = [];
 if ($maxResults) {
 $arguments['MaxResults'] = $maxResults;
 }
 if ($nextToken) {
 $arguments['NextToken'] = $nextToken;
 }
 if (!empty($tags)) {
 $arguments['Tags'] = $tags;
 }
 return $this->glueClient->listJobs($arguments);
}

public function getJobRuns($jobName, $maxResults = 0, $nextToken = ''):
Result
{
 $arguments = ['JobName' => $jobName];
 if ($maxResults) {
 $arguments['MaxResults'] = $maxResults;
 }
 if ($nextToken) {
 $arguments['NextToken'] = $nextToken;
 }
 return $this->glueClient->getJobRuns($arguments);
}
```

```
public function getJobRun($jobName, $runId, $predecessorsIncluded = false):
Result
{
 return $this->glueClient->getJobRun([
 'JobName' => $jobName,
 'RunId' => $runId,
 'PredecessorsIncluded' => $predecessorsIncluded,
]);
}

public function deleteJob($jobName)
{
 return $this->glueClient->deleteJob([
 'JobName' => $jobName,
]);
}

public function deleteTable($tableName, $databaseName)
{
 return $this->glueClient->deleteTable([
 'DatabaseName' => $databaseName,
 'Name' => $tableName,
]);
}

public function deleteDatabase($databaseName)
{
 return $this->glueClient->deleteDatabase([
 'Name' => $databaseName,
]);
}

public function deleteCrawler($crawlerName)
{
 return $this->glueClient->deleteCrawler([
 'Name' => $crawlerName,
]);
}
}
```

- APIの詳細については、『AWS SDK for PHP API リファレンス』の以下のトピックを参照してください。
  - [CreateCrawler](#)
  - [CreateJob](#)
  - [DeleteCrawler](#)
  - [DeleteDatabase](#)
  - [DeleteJob](#)
  - [DeleteTable](#)
  - [GetCrawler](#)
  - [GetDatabase](#)
  - [GetDatabases](#)
  - [GetJob](#)
  - [GetJobRun](#)
  - [GetJobRuns](#)
  - [GetTables](#)
  - [ListJobs](#)
  - [StartCrawler](#)
  - [StartJobRun](#)

## Python

### SDK for Python (Boto3)

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

シナリオで使用される AWS Glue 関数をラップするクラスを作成します。

```
class GlueWrapper:
 """Encapsulates AWS Glue actions."""

 def __init__(self, glue_client):
```

```
 """
 :param glue_client: A Boto3 Glue client.
 """
 self.glue_client = glue_client

def get_crawler(self, name):
 """
 Gets information about a crawler.

 :param name: The name of the crawler to look up.
 :return: Data about the crawler.
 """
 crawler = None
 try:
 response = self.glue_client.get_crawler(Name=name)
 crawler = response["Crawler"]
 except ClientError as err:
 if err.response["Error"]["Code"] == "EntityNotFoundException":
 logger.info("Crawler %s doesn't exist.", name)
 else:
 logger.error(
 "Couldn't get crawler %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
 return crawler

def create_crawler(self, name, role_arn, db_name, db_prefix, s3_target):
 """
 Creates a crawler that can crawl the specified target and populate a
 database in your AWS Glue Data Catalog with metadata that describes the
 data
 in the target.

 :param name: The name of the crawler.
 :param role_arn: The Amazon Resource Name (ARN) of an AWS Identity and
 Access
 Management (IAM) role that grants permission to let AWS
 Glue
 access the resources it needs.
```

```
 :param db_name: The name to give the database that is created by the
crawler.
 :param db_prefix: The prefix to give any database tables that are created
by
 the crawler.
 :param s3_target: The URL to an S3 bucket that contains data that is
 the target of the crawler.
 """
 try:
 self.glue_client.create_crawler(
 Name=name,
 Role=role_arn,
 DatabaseName=db_name,
 TablePrefix=db_prefix,
 Targets={"S3Targets": [{"Path": s3_target}]},
)
 except ClientError as err:
 logger.error(
 "Couldn't create crawler. Here's why: %s: %s",
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise

def start_crawler(self, name):
 """
 Starts a crawler. The crawler crawls its configured target and creates
 metadata that describes the data it finds in the target data source.

 :param name: The name of the crawler to start.
 """
 try:
 self.glue_client.start_crawler(Name=name)
 except ClientError as err:
 logger.error(
 "Couldn't start crawler %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
```

```
def get_database(self, name):
 """
 Gets information about a database in your Data Catalog.

 :param name: The name of the database to look up.
 :return: Information about the database.
 """
 try:
 response = self.glue_client.get_database(Name=name)
 except ClientError as err:
 logger.error(
 "Couldn't get database %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
 else:
 return response["Database"]

def get_tables(self, db_name):
 """
 Gets a list of tables in a Data Catalog database.

 :param db_name: The name of the database to query.
 :return: The list of tables in the database.
 """
 try:
 response = self.glue_client.get_tables(DatabaseName=db_name)
 except ClientError as err:
 logger.error(
 "Couldn't get tables %s. Here's why: %s: %s",
 db_name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
 else:
 return response["TableList"]

def create_job(self, name, description, role_arn, script_location):
 """
```

```
Creates a job definition for an extract, transform, and load (ETL) job
that can
 be run by AWS Glue.

:param name: The name of the job definition.
:param description: The description of the job definition.
:param role_arn: The ARN of an IAM role that grants AWS Glue the
permissions
 it requires to run the job.
:param script_location: The Amazon S3 URL of a Python ETL script that is
run as
 part of the job. The script defines how the data
is
 transformed.
"""
try:
 self.glue_client.create_job(
 Name=name,
 Description=description,
 Role=role_arn,
 Command={
 "Name": "glueetl",
 "ScriptLocation": script_location,
 "PythonVersion": "3",
 },
 GlueVersion="3.0",
)
except ClientError as err:
 logger.error(
 "Couldn't create job %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise

def start_job_run(self, name, input_database, input_table,
output_bucket_name):
 """
 Starts a job run. A job run extracts data from the source, transforms it,
 and loads it to the output bucket.

 :param name: The name of the job definition.
```

```
 :param input_database: The name of the metadata database that contains
tables
 that describe the source data. This is typically
created
 by a crawler.
 :param input_table: The name of the table in the metadata database that
 describes the source data.
 :param output_bucket_name: The S3 bucket where the output is written.
 :return: The ID of the job run.
 """
 try:
 # The custom Arguments that are passed to this function are used by
the
 # Python ETL script to determine the location of input and output
data.

 response = self.glue_client.start_job_run(
 JobName=name,
 Arguments={
 "--input_database": input_database,
 "--input_table": input_table,
 "--output_bucket_url": f"s3://{output_bucket_name}/",
 },
)
 except ClientError as err:
 logger.error(
 "Couldn't start job run %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
 else:
 return response["JobRunId"]

def list_jobs(self):
 """
 Lists the names of job definitions in your account.

 :return: The list of job definition names.
 """
 try:
 response = self.glue_client.list_jobs()
 except ClientError as err:
```

```
 logger.error(
 "Couldn't list jobs. Here's why: %s: %s",
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
 else:
 return response["JobNames"]

def get_job_runs(self, job_name):
 """
 Gets information about runs that have been performed for a specific job
 definition.

 :param job_name: The name of the job definition to look up.
 :return: The list of job runs.
 """
 try:
 response = self.glue_client.get_job_runs(JobName=job_name)
 except ClientError as err:
 logger.error(
 "Couldn't get job runs for %s. Here's why: %s: %s",
 job_name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
 else:
 return response["JobRuns"]

def get_job_run(self, name, run_id):
 """
 Gets information about a single job run.

 :param name: The name of the job definition for the run.
 :param run_id: The ID of the run.
 :return: Information about the run.
 """
 try:
 response = self.glue_client.get_job_run(JobName=name, RunId=run_id)
 except ClientError as err:
 logger.error(
```

```
 "Couldn't get job run %s/%s. Here's why: %s: %s",
 name,
 run_id,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
else:
 return response["JobRun"]

def delete_job(self, job_name):
 """
 Deletes a job definition. This also deletes data about all runs that are
 associated with this job definition.

 :param job_name: The name of the job definition to delete.
 """
 try:
 self.glue_client.delete_job(JobName=job_name)
 except ClientError as err:
 logger.error(
 "Couldn't delete job %s. Here's why: %s: %s",
 job_name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise

def delete_table(self, db_name, table_name):
 """
 Deletes a table from a metadata database.

 :param db_name: The name of the database that contains the table.
 :param table_name: The name of the table to delete.
 """
 try:
 self.glue_client.delete_table(DatabaseName=db_name, Name=table_name)
 except ClientError as err:
 logger.error(
 "Couldn't delete table %s. Here's why: %s: %s",
 table_name,
 err.response["Error"]["Code"],
```

```
 err.response["Error"]["Message"],
)
 raise

def delete_database(self, name):
 """
 Deletes a metadata database from your Data Catalog.

 :param name: The name of the database to delete.
 """
 try:
 self.glue_client.delete_database(Name=name)
 except ClientError as err:
 logger.error(
 "Couldn't delete database %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise

def delete_crawler(self, name):
 """
 Deletes a crawler.

 :param name: The name of the crawler to delete.
 """
 try:
 self.glue_client.delete_crawler(Name=name)
 except ClientError as err:
 logger.error(
 "Couldn't delete crawler %s. Here's why: %s: %s",
 name,
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
```

シナリオを実行するクラスを作成します。

```
class GlueCrawlerJobScenario:
 """
 Encapsulates a scenario that shows how to create an AWS Glue crawler and job
 and use
 them to transform data from CSV to JSON format.
 """

 def __init__(self, glue_client, glue_service_role, glue_bucket):
 """
 :param glue_client: A Boto3 AWS Glue client.
 :param glue_service_role: An AWS Identity and Access Management (IAM)
role
 that AWS Glue can assume to gain access to the
 resources it requires.
 :param glue_bucket: An S3 bucket that can hold a job script and output
data
 from AWS Glue job runs.
 """
 self.glue_client = glue_client
 self.glue_service_role = glue_service_role
 self.glue_bucket = glue_bucket

 @staticmethod
 def wait(seconds, tick=12):
 """
 Waits for a specified number of seconds, while also displaying an
animated
 spinner.

 :param seconds: The number of seconds to wait.
 :param tick: The number of frames per second used to animate the spinner.
 """
 progress = "|/-\\"
 waited = 0
 while waited < seconds:
 for frame in range(tick):
 sys.stdout.write(f"\r{progress[frame % len(progress)]}")
 sys.stdout.flush()
 time.sleep(1 / tick)
 waited += 1

 def upload_job_script(self, job_script):
```

```

 """
 Uploads a Python ETL script to an S3 bucket. The script is used by the
 AWS Glue
 job to transform data.

 :param job_script: The relative path to the job script.
 """
 try:
 self.glue_bucket.upload_file(Filename=job_script, Key=job_script)
 print(f"Uploaded job script '{job_script}' to the example bucket.")
 except S3UploadFailedError as err:
 logger.error("Couldn't upload job script. Here's why: %s", err)
 raise

 def run(self, crawler_name, db_name, db_prefix, data_source, job_script,
 job_name):
 """
 Runs the scenario. This is an interactive experience that runs at a
 command
 prompt and asks you for input throughout.

 :param crawler_name: The name of the crawler used in the scenario. If the
 crawler does not exist, it is created.
 :param db_name: The name to give the metadata database created by the
 crawler.
 :param db_prefix: The prefix to give tables added to the database by the
 crawler.
 :param data_source: The location of the data source that is targeted by
 the
 crawler and extracted during job runs.
 :param job_script: The job script that is used to transform data during
 job
 runs.
 :param job_name: The name to give the job definition that is created
 during the
 scenario.
 """
 wrapper = GlueWrapper(self.glue_client)
 print(f"Checking for crawler {crawler_name}.")
 crawler = wrapper.get_crawler(crawler_name)
 if crawler is None:
 print(f"Creating crawler {crawler_name}.")
 wrapper.create_crawler(
 crawler_name,

```

```
 self.glue_service_role.arn,
 db_name,
 db_prefix,
 data_source,
)
 print(f"Created crawler {crawler_name}.")
 crawler = wrapper.get_crawler(crawler_name)
pprint(crawler)
print("-" * 88)

print(
 f"When you run the crawler, it crawls data stored in {data_source}
and "
 f"creates a metadata database in the AWS Glue Data Catalog that
describes "
 f"the data in the data source."
)
print("In this example, the source data is in CSV format.")
ready = False
while not ready:
 ready = Question.ask_question(
 "Ready to start the crawler? (y/n) ", Question.is_yesno
)
wrapper.start_crawler(crawler_name)
print("Let's wait for the crawler to run. This typically takes a few
minutes.")
crawler_state = None
while crawler_state != "READY":
 self.wait(10)
 crawler = wrapper.get_crawler(crawler_name)
 crawler_state = crawler["State"]
 print(f"Crawler is {crawler['State']}.")
print("-" * 88)

database = wrapper.get_database(db_name)
print(f"The crawler created database {db_name}:")
pprint(database)
print(f"The database contains these tables:")
tables = wrapper.get_tables(db_name)
for index, table in enumerate(tables):
 print(f"\t{index + 1}. {table['Name']}")
table_index = Question.ask_question(
 f"Enter the number of a table to see more detail: ",
 Question.is_int,
```

```
 Question.in_range(1, len(tables)),
)
 pprint(tables[table_index - 1])
 print("-" * 88)

 print(f"Creating job definition {job_name}.")
 wrapper.create_job(
 job_name,
 "Getting started example job.",
 self.glue_service_role.arn,
 f"s3://{self.glue_bucket.name}/{job_script}",
)
 print("Created job definition.")
 print(
 f"When you run the job, it extracts data from {data_source},
transforms it "
 f"by using the {job_script} script, and loads the output into "
 f"S3 bucket {self.glue_bucket.name}."
)
 print(
 "In this example, the data is transformed from CSV to JSON, and only
a few "
 "fields are included in the output."
)
 job_run_status = None
 if Question.ask_question(f"Ready to run? (y/n) ", Question.is_yesno):
 job_run_id = wrapper.start_job_run(
 job_name, db_name, tables[0]["Name"], self.glue_bucket.name
)
 print(f"Job {job_name} started. Let's wait for it to run.")
 while job_run_status not in ["SUCCEEDED", "STOPPED", "FAILED",
"TIMEOUT"]:
 self.wait(10)
 job_run = wrapper.get_job_run(job_name, job_run_id)
 job_run_status = job_run["JobRunState"]
 print(f"Job {job_name}/{job_run_id} is {job_run_status}.")
 print("-" * 88)

 if job_run_status == "SUCCEEDED":
 print(
 f"Data from your job run is stored in your S3 bucket
'{self.glue_bucket.name}':"
)
 try:
```

```

 keys = [
 obj.key for obj in
self.glue_bucket.objects.filter(Prefix="run-")
]
 for index, key in enumerate(keys):
 print(f"\t{index + 1}: {key}")
 lines = 4
 key_index = Question.ask_question(
 f"Enter the number of a block to download it and see the
first {lines} "
 f"lines of JSON output in the block: ",
 Question.is_int,
 Question.in_range(1, len(keys)),
)
 job_data = io.BytesIO()
 self.glue_bucket.download_fileobj(keys[key_index - 1], job_data)
 job_data.seek(0)
 for _ in range(lines):
 print(job_data.readline().decode("utf-8"))
 except ClientError as err:
 logger.error(
 "Couldn't get job run data. Here's why: %s: %s",
 err.response["Error"]["Code"],
 err.response["Error"]["Message"],
)
 raise
 print("-" * 88)

 job_names = wrapper.list_jobs()
 if job_names:
 print(f"Your account has {len(job_names)} jobs defined:")
 for index, job_name in enumerate(job_names):
 print(f"\t{index + 1}. {job_name}")
 job_index = Question.ask_question(
 f"Enter a number between 1 and {len(job_names)} to see the list
of runs for "
 f"a job: ",
 Question.is_int,
 Question.in_range(1, len(job_names)),
)
 job_runs = wrapper.get_job_runs(job_names[job_index - 1])
 if job_runs:
 print(f"Found {len(job_runs)} runs for job {job_names[job_index -
1]}:")

```

```

 for index, job_run in enumerate(job_runs):
 print(
 f"\t{index + 1}. {job_run['JobRunState']} on "
 f"{job_run['CompletedOn']:%Y-%m-%d %H:%M:%S}"
)
 run_index = Question.ask_question(
 f"Enter a number between 1 and {len(job_runs)} to see details
for a run: ",
 Question.is_int,
 Question.in_range(1, len(job_runs)),
)
 pprint(job_runs[run_index - 1])
 else:
 print(f"No runs found for job {job_names[job_index - 1]}")
else:
 print("Your account doesn't have any jobs defined.")
print("-" * 88)

print(
 f"Let's clean up. During this example we created job definition
'{job_name}'."
)
if Question.ask_question(
 "Do you want to delete the definition and all runs? (y/n) ",
 Question.is_yesno,
):
 wrapper.delete_job(job_name)
 print(f"Job definition '{job_name}' deleted.")
tables = wrapper.get_tables(db_name)
print(f"We also created database '{db_name}' that contains these
tables:")
for table in tables:
 print(f"\t{table['Name']}")
if Question.ask_question(
 "Do you want to delete the tables and the database? (y/n) ",
 Question.is_yesno,
):
 for table in tables:
 wrapper.delete_table(db_name, table["Name"])
 print(f"Deleted table {table['Name']}")
 wrapper.delete_database(db_name)
 print(f"Deleted database {db_name}.")
print(f"We also created crawler '{crawler_name}'.")
if Question.ask_question(

```

```
 "Do you want to delete the crawler? (y/n) ", Question.is_yesno
):
 wrapper.delete_crawler(crawler_name)
 print(f"Deleted crawler {crawler_name}.")
 print("-" * 88)

def parse_args(args):
 """
 Parse command line arguments.

 :param args: The command line arguments.
 :return: The parsed arguments.
 """
 parser = argparse.ArgumentParser(
 description="Runs the AWS Glue getting started with crawlers and jobs
 scenario. "
 "Before you run this scenario, set up scaffold resources by running "
 "'python scaffold.py deploy'."
)
 parser.add_argument(
 "role_name",
 help="The name of an IAM role that AWS Glue can assume. This role must
 grant access "
 "to Amazon S3 and to the permissions granted by the AWSGlueServiceRole "
 "managed policy.",
)
 parser.add_argument(
 "bucket_name",
 help="The name of an S3 bucket that AWS Glue can access to get the job
 script and "
 "put job results.",
)
 parser.add_argument(
 "--job_script",
 default="flight_etl_job_script.py",
 help="The name of the job script file that is used in the scenario.",
)
 return parser.parse_args(args)

def main():
 args = parse_args(sys.argv[1:])
 try:
```

```
print("-" * 88)
print(
 "Welcome to the AWS Glue getting started with crawlers and jobs
scenario."
)
print("-" * 88)
scenario = GlueCrawlerJobScenario(
 boto3.client("glue"),
 boto3.resource("iam").Role(args.role_name),
 boto3.resource("s3").Bucket(args.bucket_name),
)
scenario.upload_job_script(args.job_script)
scenario.run(
 "doc-example-crawler",
 "doc-example-database",
 "doc-example-",
 "s3://crawler-public-us-east-1/flight/2016/csv",
 args.job_script,
 "doc-example-job",
)
print("-" * 88)
print(
 "To destroy scaffold resources, including the IAM role and S3 bucket
"
 "used in this scenario, run 'python scaffold.py destroy'."
)
print("\nThanks for watching!")
print("-" * 88)
except Exception:
 logging.exception("Something went wrong with the example.")
```

ジョブの実行中にデータを抽出、変換、ロード AWS Glue するためにで使用される ETL スクリプトを作成します。

```
import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job
```

```
"""
These custom arguments must be passed as Arguments to the StartJobRun request.
 --input_database The name of a metadata database that is contained in
your
 AWS Glue Data Catalog and that contains tables that
describe
 the data to be processed.
 --input_table The name of a table in the database that describes the
data to
 be processed.
 --output_bucket_url An S3 bucket that receives the transformed output data.
"""
args = getResolvedOptions(
 sys.argv, ["JOB_NAME", "input_database", "input_table", "output_bucket_url"]
)
sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session
job = Job(glueContext)
job.init(args["JOB_NAME"], args)

Script generated for node S3 Flight Data.
S3FlightData_node1 = glueContext.create_dynamic_frame.from_catalog(
 database=args["input_database"],
 table_name=args["input_table"],
 transformation_ctx="S3FlightData_node1",
)

This mapping performs two main functions:
1. It simplifies the output by removing most of the fields from the data.
2. It renames some fields. For example, `fl_date` is renamed to `flight_date`.
ApplyMapping_node2 = ApplyMapping.apply(
 frame=S3FlightData_node1,
 mappings=[
 ("year", "long", "year", "long"),
 ("month", "long", "month", "tinyint"),
 ("day_of_month", "long", "day", "tinyint"),
 ("fl_date", "string", "flight_date", "string"),
 ("carrier", "string", "carrier", "string"),
 ("fl_num", "long", "flight_num", "long"),
 ("origin_city_name", "string", "origin_city_name", "string"),
 ("origin_state_abr", "string", "origin_state_abr", "string"),
 ("dest_city_name", "string", "dest_city_name", "string"),
```

```
 ("dest_state_abr", "string", "dest_state_abr", "string"),
 ("dep_time", "long", "departure_time", "long"),
 ("wheels_off", "long", "wheels_off", "long"),
 ("wheels_on", "long", "wheels_on", "long"),
 ("arr_time", "long", "arrival_time", "long"),
 ("mon", "string", "mon", "string"),
],
 transformation_ctx="ApplyMapping_node2",
)

Script generated for node Revised Flight Data.
RevisedFlightData_node3 = glueContext.write_dynamic_frame.from_options(
 frame=ApplyMapping_node2,
 connection_type="s3",
 format="json",
 connection_options={"path": args["output_bucket_url"], "partitionKeys": []},
 transformation_ctx="RevisedFlightData_node3",
)

job.commit()
```

- APIの詳細については、『AWS SDK for Python (Boto3) API リファレンス』の以下のトピックを参照してください。

- [CreateCrawler](#)
- [CreateJob](#)
- [DeleteCrawler](#)
- [DeleteDatabase](#)
- [DeleteJob](#)
- [DeleteTable](#)
- [GetCrawler](#)
- [GetDatabase](#)
- [GetDatabases](#)
- [GetJob](#)
- [GetJobRun](#)
- [GetJobRuns](#)
- [GetTables](#)

- [ListJobs](#)
- [StartCrawler](#)
- [StartJobRun](#)

## Ruby

### SDK for Ruby

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

シナリオで使用される AWS Glue 関数をラップするクラスを作成します。

```
The `GlueWrapper` class serves as a wrapper around the AWS Glue API, providing
a simplified interface for common operations.
It encapsulates the functionality of the AWS SDK for Glue and provides methods
for interacting with Glue crawlers, databases, tables, jobs, and S3 resources.
The class initializes with a Glue client and a logger, allowing it to make API
calls and log any errors or informational messages.
class GlueWrapper
 def initialize(glue_client, logger)
 @glue_client = glue_client
 @logger = logger
 end

 # Retrieves information about a specific crawler.
 #
 # @param name [String] The name of the crawler to retrieve information about.
 # @return [Aws::Glue::Types::Crawler, nil] The crawler object if found, or nil
 # if not found.
 def get_crawler(name)
 @glue_client.get_crawler(name: name)
 rescue Aws::Glue::Errors::EntityNotFoundException
 @logger.info("Crawler #{name} doesn't exist.")
 false
 rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not get crawler #{name}: \n#{e.message}")
 end
end
```

```
 raise
 end

 # Creates a new crawler with the specified configuration.
 #
 # @param name [String] The name of the crawler.
 # @param role_arn [String] The ARN of the IAM role to be used by the crawler.
 # @param db_name [String] The name of the database where the crawler stores its
 metadata.
 # @param db_prefix [String] The prefix to be added to the names of tables that
 the crawler creates.
 # @param s3_target [String] The S3 path that the crawler will crawl.
 # @return [void]
 def create_crawler(name, role_arn, db_name, db_prefix, s3_target)
 @glue_client.create_crawler(
 name: name,
 role: role_arn,
 database_name: db_name,
 targets: {
 s3_targets: [
 {
 path: s3_target
 }
]
 }
)
 rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not create crawler: \n#{e.message}")
 raise
 end

 # Starts a crawler with the specified name.
 #
 # @param name [String] The name of the crawler to start.
 # @return [void]
 def start_crawler(name)
 @glue_client.start_crawler(name: name)
 rescue Aws::Glue::Errors::ServiceError => e
 @logger.error("Glue could not start crawler #{name}: \n#{e.message}")
 raise
 end

 # Deletes a crawler with the specified name.
 #
```

```
@param name [String] The name of the crawler to delete.
@return [void]
def delete_crawler(name)
 @glue_client.delete_crawler(name: name)
rescue Aws::Glue::Errors::ServiceError => e
 @logger.error("Glue could not delete crawler #{name}: \n#{e.message}")
 raise
end

Retrieves information about a specific database.
#
@param name [String] The name of the database to retrieve information about.
@return [Aws::Glue::Types::Database, nil] The database object if found, or
nil if not found.
def get_database(name)
 response = @glue_client.get_database(name: name)
 response.database
rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not get database #{name}: \n#{e.message}")
 raise
end

Retrieves a list of tables in the specified database.
#
@param db_name [String] The name of the database to retrieve tables from.
@return [Array<Aws::Glue::Types::Table>]
def get_tables(db_name)
 response = @glue_client.get_tables(database_name: db_name)
 response.table_list
rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not get tables #{db_name}: \n#{e.message}")
 raise
end

Creates a new job with the specified configuration.
#
@param name [String] The name of the job.
@param description [String] The description of the job.
@param role_arn [String] The ARN of the IAM role to be used by the job.
@param script_location [String] The location of the ETL script for the job.
@return [void]
def create_job(name, description, role_arn, script_location)
 @glue_client.create_job(
 name: name,
```

```
description: description,
role: role_arn,
command: {
 name: "glueetl",
 script_location: script_location,
 python_version: "3"
},
glue_version: "3.0"
)
rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not create job #{name}: \n#{e.message}")
 raise
end

Starts a job run for the specified job.
#
@param name [String] The name of the job to start the run for.
@param input_database [String] The name of the input database for the job.
@param input_table [String] The name of the input table for the job.
@param output_bucket_name [String] The name of the output S3 bucket for the
job.
@return [String] The ID of the started job run.
def start_job_run(name, input_database, input_table, output_bucket_name)
 response = @glue_client.start_job_run(
 job_name: name,
 arguments: {
 '--input_database': input_database,
 '--input_table': input_table,
 '--output_bucket_url': "s3://#{output_bucket_name}/"
 }
)
 response.job_run_id
rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not start job run #{name}: \n#{e.message}")
 raise
end

Retrieves a list of jobs in AWS Glue.
#
@return [Aws::Glue::Types::ListJobsResponse]
def list_jobs
 @glue_client.list_jobs
rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not list jobs: \n#{e.message}")
end
```

```
 raise
 end

 # Retrieves a list of job runs for the specified job.
 #
 # @param job_name [String] The name of the job to retrieve job runs for.
 # @return [Array<Aws::Glue::Types::JobRun>]
 def get_job_runs(job_name)
 response = @glue_client.get_job_runs(job_name: job_name)
 response.job_runs
 rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not get job runs: \n#{e.message}")
 end

 # Retrieves data for a specific job run.
 #
 # @param job_name [String] The name of the job run to retrieve data for.
 # @return [Glue::Types::GetJobRunResponse]
 def get_job_run(job_name, run_id)
 @glue_client.get_job_run(job_name: job_name, run_id: run_id)
 rescue Aws::Glue::Errors::GlueException => e
 @logger.error("Glue could not get job runs: \n#{e.message}")
 end

 # Deletes a job with the specified name.
 #
 # @param job_name [String] The name of the job to delete.
 # @return [void]
 def delete_job(job_name)
 @glue_client.delete_job(job_name: job_name)
 rescue Aws::Glue::Errors::ServiceError => e
 @logger.error("Glue could not delete job: \n#{e.message}")
 end

 # Deletes a table with the specified name.
 #
 # @param database_name [String] The name of the catalog database in which the
 table resides.
 # @param table_name [String] The name of the table to be deleted.
 # @return [void]
 def delete_table(database_name, table_name)
 @glue_client.delete_table(database_name: database_name, name: table_name)
 rescue Aws::Glue::Errors::ServiceError => e
 @logger.error("Glue could not delete job: \n#{e.message}")
 end
end
```

```
end

Removes a specified database from a Data Catalog.
#
@param database_name [String] The name of the database to delete.
@return [void]
def delete_database(database_name)
 @glue_client.delete_database(name: database_name)
rescue Aws::Glue::Errors::ServiceError => e
 @logger.error("Glue could not delete database: \n#{e.message}")
end

Uploads a job script file to an S3 bucket.
#
@param file_path [String] The local path of the job script file.
@param bucket_resource [Aws::S3::Bucket] The S3 bucket resource to upload the
file to.
@return [void]
def upload_job_script(file_path, bucket_resource)
 File.open(file_path) do |file|
 bucket_resource.client.put_object({
 body: file,
 bucket: bucket_resource.name,
 key: file_path
 })
 end
rescue Aws::S3::Errors::S3UploadFailedError => e
 @logger.error("S3 could not upload job script: \n#{e.message}")
 raise
end

end
```

シナリオを実行するクラスを作成します。

```
class GlueCrawlerJobScenario
 def initialize(glue_client, glue_service_role, glue_bucket, logger)
 @glue_client = glue_client
 @glue_service_role = glue_service_role
 @glue_bucket = glue_bucket
 @logger = logger
 end
end
```

```
def run(crawler_name, db_name, db_prefix, data_source, job_script, job_name)
 wrapper = GlueWrapper.new(@glue_client, @logger)

 new_step(1, "Create a crawler")
 puts "Checking for crawler #{crawler_name}."
 crawler = wrapper.get_crawler(crawler_name)
 if crawler == false
 puts "Creating crawler #{crawler_name}."
 wrapper.create_crawler(crawler_name, @glue_service_role.arn, db_name,
db_prefix, data_source)
 puts "Successfully created #{crawler_name}:"
 crawler = wrapper.get_crawler(crawler_name)
 puts JSON.pretty_generate(crawler).yellow
 end
 print "\nDone!\n".green

 new_step(2, "Run a crawler to output a database.")
 puts "Location of input data analyzed by crawler: #{data_source}"
 puts "Outputs: a Data Catalog database in CSV format containing metadata on
input."
 wrapper.start_crawler(crawler_name)
 puts "Starting crawler... (this typically takes a few minutes)"
 crawler_state = nil
 while crawler_state != "READY"
 custom_wait(15)
 crawler = wrapper.get_crawler(crawler_name)
 crawler_state = crawler[0]["state"]
 print "Status check: #{crawler_state}.".yellow
 end
 print "\nDone!\n".green

 new_step(3, "Query the database.")
 database = wrapper.get_database(db_name)
 puts "The crawler created database #{db_name}:"
 print "#{database}.".yellow
 puts "\nThe database contains these tables:"
 tables = wrapper.get_tables(db_name)
 tables.each_with_index do |table, index|
 print "\t#{index + 1}. #{table['name']}".yellow
 end
 print "\nDone!\n".green

 new_step(4, "Create a job definition that runs an ETL script.")
```

```
puts "Uploading Python ETL script to S3..."
wrapper.upload_job_script(job_script, @glue_bucket)
puts "Creating job definition #{job_name}:\n"
response = wrapper.create_job(job_name, "Getting started example job.",
@glue_service_role.arn, "s3://#{@glue_bucket.name}/#{@job_script}")
puts JSON.pretty_generate(response).yellow
print "\nDone!\n".green

new_step(5, "Start a new job")
job_run_status = nil
job_run_id = wrapper.start_job_run(
 job_name,
 db_name,
 tables[0]["name"],
 @glue_bucket.name
)
puts "Job #{job_name} started. Let's wait for it to run."
until ["SUCCEEDED", "STOPPED", "FAILED", "TIMEOUT"].include?(job_run_status)
 custom_wait(10)
 job_run = wrapper.get_job_runs(job_name)
 job_run_status = job_run[0]["job_run_state"]
 print "Status check: #{job_name}/#{job_run_id} - #{job_run_status}.".yellow
end
print "\nDone!\n".green

new_step(6, "View results from a successful job run.")
if job_run_status == "SUCCEEDED"
 puts "Data from your job run is stored in your S3 bucket
'#{@glue_bucket.name}'. Files include:"
 begin

 # Print the key name of each object in the bucket.
 @glue_bucket.objects.each do |object_summary|
 if object_summary.key.include?("run-")
 print "#{object_summary.key}".yellow
 end
 end

 # Print the first 256 bytes of a run file
 desired_sample_objects = 1
 @glue_bucket.objects.each do |object_summary|
 if object_summary.key.include?("run-")
 if desired_sample_objects > 0
 sample_object = @glue_bucket.object(object_summary.key)
 end
 end
 end
 end
end
```

```

 sample = sample_object.get(range: "bytes=0-255").body.read
 puts "\nSample run file contents:"
 print "#{sample}".yellow
 desired_sample_objects -= 1
 end
 end
 end
end
rescue Aws::S3::Errors::ServiceError => e
 logger.error(
 "Couldn't get job run data. Here's why: %s: %s",
 e.response.error.code, e.response.error.message
)
 raise
end
end
print "\nDone!\n".green

new_step(7, "Delete job definition and crawler.")
wrapper.delete_job(job_name)
puts "Job deleted: #{job_name}."
wrapper.delete_crawler(crawler_name)
puts "Crawler deleted: #{crawler_name}."
wrapper.delete_table(db_name, tables[0]["name"])
puts "Table deleted: #{tables[0]["name"]} in #{db_name}."
wrapper.delete_database(db_name)
puts "Database deleted: #{db_name}."
print "\nDone!\n".green
end
end

def main

 banner("../helpers/banner.txt")
 puts
 "#####"
 puts "#
 #".yellow
 puts "# EXAMPLE CODE DEMO:
 #".yellow
 puts "# AWS Glue
 #".yellow
 puts "#
 #".yellow

```

```
puts
#####
puts ""
puts "You have launched a demo of AWS Glue using the AWS for Ruby v3 SDK. Over
the next 60 seconds, it will"
puts "do the following:"
puts " 1. Create a crawler."
puts " 2. Run a crawler to output a database."
puts " 3. Query the database."
puts " 4. Create a job definition that runs an ETL script."
puts " 5. Start a new job."
puts " 6. View results from a successful job run."
puts " 7. Delete job definition and crawler."
puts ""

confirm_begin
billing
security
puts "\e[H\e[2J"

Set input file names
job_script_filepath = "job_script.py"
resource_names = YAML.load_file("resource_names.yaml")

Instantiate existing IAM role.
iam = Aws::IAM::Resource.new(region: "us-east-1")
iam_role_name = resource_names["glue_service_role"]
iam_role = iam.role(iam_role_name)

Instantiate existing S3 bucket.
s3 = Aws::S3::Resource.new(region: "us-east-1")
s3_bucket_name = resource_names["glue_bucket"]
s3_bucket = s3.bucket(s3_bucket_name)

scenario = GlueCrawlerJobScenario.new(
 Aws::Glue::Client.new(region: "us-east-1"),
 iam_role,
 s3_bucket,
 @logger
)

random_int = rand(10 ** 4)
scenario.run(
 "doc-example-crawler-#{random_int}",
```

```

 "doc-example-database-#{random_int}",
 "doc-example-#{random_int}-",
 "s3://crawler-public-us-east-1/flight/2016/csv",
 job_script_filepath,
 "doc-example-job-#{random_int}"
)

 puts "-" * 88
 puts "You have reached the end of this tour of AWS Glue."
 puts "To destroy CDK-created resources, run:\n cdk destroy"
 puts "-" * 88

end

```

ジョブの実行中にデータを抽出、変換、ロード AWS Glue するためにで使用される ETL スクリプトを作成します。

```

import sys
from awsglue.transforms import *
from awsglue.utils import getResolvedOptions
from pyspark.context import SparkContext
from awsglue.context import GlueContext
from awsglue.job import Job

"""
These custom arguments must be passed as Arguments to the StartJobRun request.
 --input_database The name of a metadata database that is contained in
your
 AWS Glue Data Catalog and that contains tables that
describe
 the data to be processed.
 --input_table The name of a table in the database that describes the
data to
 be processed.
 --output_bucket_url An S3 bucket that receives the transformed output data.
"""
args = getResolvedOptions(
 sys.argv, ["JOB_NAME", "input_database", "input_table", "output_bucket_url"]
)
sc = SparkContext()
glueContext = GlueContext(sc)
spark = glueContext.spark_session

```

```
job = Job(glueContext)
job.init(args["JOB_NAME"], args)

Script generated for node S3 Flight Data.
S3FlightData_node1 = glueContext.create_dynamic_frame.from_catalog(
 database=args["input_database"],
 table_name=args["input_table"],
 transformation_ctx="S3FlightData_node1",
)

This mapping performs two main functions:
1. It simplifies the output by removing most of the fields from the data.
2. It renames some fields. For example, `fl_date` is renamed to `flight_date`.
ApplyMapping_node2 = ApplyMapping.apply(
 frame=S3FlightData_node1,
 mappings=[
 ("year", "long", "year", "long"),
 ("month", "long", "month", "tinyint"),
 ("day_of_month", "long", "day", "tinyint"),
 ("fl_date", "string", "flight_date", "string"),
 ("carrier", "string", "carrier", "string"),
 ("fl_num", "long", "flight_num", "long"),
 ("origin_city_name", "string", "origin_city_name", "string"),
 ("origin_state_abr", "string", "origin_state_abr", "string"),
 ("dest_city_name", "string", "dest_city_name", "string"),
 ("dest_state_abr", "string", "dest_state_abr", "string"),
 ("dep_time", "long", "departure_time", "long"),
 ("wheels_off", "long", "wheels_off", "long"),
 ("wheels_on", "long", "wheels_on", "long"),
 ("arr_time", "long", "arrival_time", "long"),
 ("mon", "string", "mon", "string"),
],
 transformation_ctx="ApplyMapping_node2",
)

Script generated for node Revised Flight Data.
RevisedFlightData_node3 = glueContext.write_dynamic_frame.from_options(
 frame=ApplyMapping_node2,
 connection_type="s3",
 format="json",
 connection_options={"path": args["output_bucket_url"], "partitionKeys": []},
 transformation_ctx="RevisedFlightData_node3",
)
```

```
job.commit()
```

- APIの詳細については、「AWS SDK for Ruby API リファレンス」の以下のトピックを参照してください。
  - [CreateCrawler](#)
  - [CreateJob](#)
  - [DeleteCrawler](#)
  - [DeleteDatabase](#)
  - [DeleteJob](#)
  - [DeleteTable](#)
  - [GetCrawler](#)
  - [GetDatabase](#)
  - [GetDatabases](#)
  - [GetJob](#)
  - [GetJobRun](#)
  - [GetJobRuns](#)
  - [GetTables](#)
  - [ListJobs](#)
  - [StartCrawler](#)
  - [StartJobRun](#)

## Rust

### SDK for Rust

#### Note

については、「」を参照してください GitHub。 [AWS コード例リポジトリ](#) で全く同じ例を見つけて、設定と実行の方法を確認してください。

パブリックの Amazon Simple Storage Service (Amazon S3) バケットをクローलし、検出した CSV 形式データを記述するメタデータデータベースを生成するクローラーを作成して実行します。

```
let create_crawler = glue
 .create_crawler()
 .name(self.crawler())
 .database_name(self.database())
 .role(self.iam_role.expose_secret())
 .targets(
 CrawlerTargets::builder()
 .s3_targets(S3Target::builder().path(CRAWLER_TARGET).build())
 .build(),
)
 .send()
 .await;

match create_crawler {
 Err(err) => {
 let glue_err: aws_sdk_glue::Error = err.into();
 match glue_err {
 aws_sdk_glue::Error::AlreadyExistsException(_) => {
 info!("Using existing crawler");
 Ok(())
 }
 _ => Err(GlueMvpError::GlueSdk(glue_err)),
 }
 }
 Ok(_) => Ok(()),
}??;

let start_crawler =
glue.start_crawler().name(self.crawler()).send().await;

match start_crawler {
 Ok(_) => Ok(()),
 Err(err) => {
 let glue_err: aws_sdk_glue::Error = err.into();
 match glue_err {
 aws_sdk_glue::Error::CrawlerRunningException(_) => Ok(()),
 _ => Err(GlueMvpError::GlueSdk(glue_err)),
 }
 }
}
```

```
}?;
```

のデータベースとテーブルに関する情報を一覧表示します AWS Glue Data Catalog。

```
let database = glue
 .get_database()
 .name(self.database())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?
 .to_owned();
let database = database
 .database()
 .ok_or_else(|| GlueMvpError::Unknown("Could not find
database".into()))?;

let tables = glue
 .get_tables()
 .database_name(self.database())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;

let tables = tables.table_list();
```

ソース Amazon S3 バケットから CSV 形式データを抽出し、フィールドを削除して名前を変更することで変換し、JSON 形式の出力を別の Amazon S3 バケットにロードするジョブを作成して実行します。

```
let create_job = glue
 .create_job()
 .name(self.job())
 .role(self.iam_role.expose_secret())
 .command(
 JobCommand::builder()
 .name("glueetl")
 .python_version("3")
 .script_location(format!("s3://{}/job.py", self.bucket()))
 .build(),
)
```

```
.glue_version("3.0")
.send()
.await
.map_err(GlueMvpError::from_glue_sdk)?;

let job_name = create_job.name().ok_or_else(|| {
 GlueMvpError::Unknown("Did not get job name after creating
job".into())
})?;

let job_run_output = glue
.start_job_run()
.job_name(self.job())
.arguments("--input_database", self.database())
.arguments(
 "--input_table",
 self.tables
 .first()
 .ok_or_else(|| GlueMvpError::Unknown("Missing crawler
table".into()))?
 .name(),
)
.arguments("--output_bucket_url", self.bucket())
.send()
.await
.map_err(GlueMvpError::from_glue_sdk)?;

let job = job_run_output
.job_run_id()
.ok_or_else(|| GlueMvpError::Unknown("Missing run id from just
started job".into()))?
.to_string();
```

デモによって作成されたすべてのリソースを削除します。

```
glue.delete_job()
.job_name(self.job())
.send()
.await
.map_err(GlueMvpError::from_glue_sdk)?;

for t in &self.tables {
```

```
 glue.delete_table()
 .name(t.name())
 .database_name(self.database())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;
 }

 glue.delete_database()
 .name(self.database())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;

 glue.delete_crawler()
 .name(self.crawler())
 .send()
 .await
 .map_err(GlueMvpError::from_glue_sdk)?;
```

- APIの詳細については、「AWS SDK for Rust API リファレンス」の以下のトピックを参照してください。

- [CreateCrawler](#)
- [CreateJob](#)
- [DeleteCrawler](#)
- [DeleteDatabase](#)
- [DeleteJob](#)
- [DeleteTable](#)
- [GetCrawler](#)
- [GetDatabase](#)
- [GetDatabases](#)
- [GetJob](#)
- [GetJobRun](#)
- [GetJobRuns](#)
- [GetTables](#)
- [ListJobs](#)

- [StartCrawler](#)
- [StartJobRun](#)

AWS SDK デベロッパーガイドとコード例の完全なリストについては、「」を参照してください [AWS SDK でこのサービスを使用する](#)。このトピックには、使用開始方法に関する情報と、以前の SDK バージョンの詳細も含まれています。

# AWS Glue のセキュリティ

AWS では、クラウドのセキュリティが最優先事項です。セキュリティを最も重視する組織の要件を満たすために構築された AWS のデータセンターとネットワークアーキテクチャは、お客様に大きく貢献します。

セキュリティは、AWS と顧客の間の責任共有です。[責任共有モデル](#)では、この責任がクラウドのセキュリティおよびクラウド内のセキュリティとして説明されています。

- クラウドのセキュリティ - AWS は、AWS クラウドで AWS のサービスを実行するインフラストラクチャを保護する責任を負います。また、AWS は、使用するサービスを安全に提供します。[AWS コンプライアンスプログラム](#)の一環として、サードパーティーの監査が定期的にセキュリティの有効性をテストおよび検証しています。AWS Glue に適用されるコンプライアンスプログラムの詳細については、「[コンプライアンスプログラムによる対象範囲内の AWS のサービス](#)」を参照してください。
- クラウド内のセキュリティ - ユーザーの責任は、使用する AWS のサービスに応じて異なります。また、お客様は、データの機密性、お客様の会社の要件、および適用される法律および規制など、その他の要因についても責任を負います。

このドキュメントは、AWS Glue を使用する際に責任共有モデルを適用する方法を理解するのに役立ちます。以下のトピックでは、セキュリティおよびコンプライアンスの目的を達成するために AWS Glue を設定する方法を示します。また、AWS Glue リソースのモニタリングや保護に役立つ、その他 AWS サービスの使用方法についても説明します。

## トピック

- [AWS Glue でのデータ保護](#)
- [AWS Glue の Identity and Access Management](#)
- [AWS Glue でのログ記録とモニタリング](#)
- [のコンプライアンス検証 AWS Glue](#)
- [のレジリエンス AWS Glue](#)
- [AWS Glue 内のインフラストラクチャセキュリティ](#)

## AWS Glue でのデータ保護

AWS Glue には、データを保護するために設計された機能がいくつか用意されています。

## トピック

- [保管中の暗号化](#)
- [送信中の暗号化](#)
- [FIPS 準拠](#)
- [キーの管理](#)
- [AWS の他のサービスへの AWS Glue の依存関係](#)
- [開発エンドポイント](#)

## 保管中の暗号化

AWS Glue は、[AWS Glue Studio でビジュアル ETL ジョブを作成する](#) と [開発エンドポイントを使用してスクリプトを開発する](#) のために保管時のデータの暗号化をサポートしています。[AWS Key Management Service \(AWS KMS\)](#) キーを使用して保管時の暗号化されたデータを書き込むように ETL (抽出、変換、ロード) ジョブと開発エンドポイントを設定できます。また、AWS KMS で管理するキーを使用して、[AWS Glue Data Catalog](#) に格納されているメタデータを暗号化することもできます。さらに、AWS KMS キーを使用して、ジョブのブックマークや、[クローラー](#)および ETL ジョブで生成されたログを暗号化することもできます。

ジョブ、クローラー、開発エンドポイントによって Amazon Simple Storage Service (Amazon S3) と Amazon CloudWatch Logs に書き込まれたデータに加えて、AWS Glue Data Catalog でメタデータオブジェクトを暗号化できます。ジョブ、クローラーそして開発エンドポイントを AWS Glue で作成する場合、セキュリティ設定をアタッチすることで暗号化設定を指定します。セキュリティ設定には、Amazon S3 で管理されるサーバー側の暗号化キー (SSE-S3)、または AWS KMS (SSE-KMS) で保管されるカスタマーマスターキー (CMK) があります。AWS Glue コンソールを使用してセキュリティ設定を作成できます。

アカウント内の Data Catalog 全体の暗号化を有効にすることもできます。そのためには、AWS KMS に格納されている CMK を指定します。

### Important

AWS Glue は、対称カスタマーマネージドキーのみをサポートしています。詳細については、「AWS Key Management Service デベロッパーガイド」の「[カスタマーマネージドキー \(CMK\)](#)」を参照してください。

暗号化を有効にすると、データカタログオブジェクトを追加するとき、クローラーやジョブを実行するとき、または開発エンドポイントを開始するとき、SSE-S3 または SSE-KMS キーを使用して、保管中のデータの書き込みが行われます。さらに、信頼された Transport Layer Security (TLS) プロトコルを使用して、Java Database Connectivity (JDBC) データストアのみにアクセスするように AWS Glue を設定できます。

AWS Glue では、以下の場所で暗号化設定をコントロールします。

- Data Catalog の設定。
- 作成したセキュリティ設定。
- AWS Glue ETL (抽出、変換、ロード) ジョブにパラメータとして渡されるサーバー側の暗号化設定 (SSE-S3 または SSE-KMS)。

暗号化を設定する方法の詳細については、「[AWS Glue での暗号化のセットアップ](#)」を参照してください。

## トピック

- [Data Catalog の暗号化](#)
- [接続パスワードの暗号化](#)
- [AWS Glue によって書き込まれたデータの暗号化](#)

## Data Catalog の暗号化

AWS Glue Data Catalog 暗号化により、機密データのセキュリティが強化されます。AWS Glue は AWS Key Management Service (AWS KMS) と統合し、Data Catalog に保存されているメタデータを暗号化します。AWS Glue コンソールまたは AWS CLI を使用し、Data Catalog でリソースの暗号化設定を有効または無効にできます。

Data Catalog の暗号化を有効にすると、新しく作成するオブジェクトはすべて暗号化されます。暗号化を無効にすると、新しく作成するオブジェクトは暗号化されませんが、既存の暗号化オブジェクトは暗号化された状態が維持されます。

AWS マネージド暗号化キーまたはカスタマーマネージド暗号化キーを使用し、Data Catalog 全体を暗号化できます。キーの種類と状態の詳細については、「AWS Key Management Service デベロッパーガイド」の「[AWS Key Management Service の概念](#)」を参照してください。

## AWS マネージドキー

AWS マネージドキーはお客様のアカウントにある KMS キーであり、AWS KMS と統合された AWS サービスがお客様に代わって作成、管理、使用します。アカウントの AWS マネージドキーの表示、キーポリシーの表示し、AWS CloudTrail ログでその使用の監査を行うことができます。ただし、キーを管理したり、許可を変更したりすることはできません。

保管時の暗号化は、メタデータの暗号化に使用される AWS Glue の AWS マネージドキーを管理するため、AWS KMS と自動的に統合されます。メタデータの暗号化を有効にしたときに AWS マネージドキーが存在しない場合、AWS KMS は自動的に新しいキーを作成します。

詳細については、「[AWS マネージドキー](#)」を参照してください。

### カスタマーマネージドキー

カスタマーマネージドキーは AWS アカウント内の KMS キーで、ユーザーが作成、所有、および管理します。ユーザーは、この KMS キーに関する完全なコントロール権を持ちます。次のようにできます。

- キーポリシー、IAM ポリシー、許可付与の確立と維持
- 有効化と無効化
- 暗号化マテリアルのローテーション
- タグを追加する
- 参照するエイリアスの作成
- 削除のスケジュール設定

カスタマーマネージドキーの許可を管理する方法の詳細については、「[カスタマーマネージドキー](#)」を参照してください。

#### Important

AWS Glue は、対称カスタマーマネージドキーのみをサポートしています。KMS キー リストには、対称キーのみが表示されます。ただし、[KMS キー ARN の選択] を選択した場合、任意のキータイプの ARN をコンソールで入力します。対称キーには ARN だけを入力するようにしてください。

対称カスタマーマネージドキーを作成するには、「AWS Key Management Service デベロッパーガイド」の「[対称カスタマーマネージドキーの作成](#)」の手順に従ってください。

保存時に Data Catalog の暗号化を有効にすると、KMS キーを使用して次のリソースタイプが暗号化されます。

- データベース
- テーブル
- パーティション
- テーブルのバージョン
- 列統計
- ユーザー定義関数
- Data Catalog のビュー

### AWS Glue の暗号化コンテキスト

[暗号化コンテキスト](#)とは、データに関する追加のコンテキスト情報を含むために、使用する (オプションの) キーと値のペアのセットです。AWS KMS は、暗号化コンテキストを[追加の認証済みデータ](#)として使用し、[暗号化の認証](#)をサポートします。データの暗号化リクエストに暗号化コンテキストを組み込むと、AWS KMS は暗号化コンテキストを暗号化後のデータにバインドします。データを復号化するには、リクエストに同じ暗号化コンテキストを含めます。AWS Glue はすべての AWS KMS の暗号化処理に同じ暗号化コンテキストを使用する際、キーは `glue_catalog_id` で値は `catalogId` を指定します。

```
"encryptionContext": {
 "glue_catalog_id": "111122223333"
}
```

AWS マネージドキーまたは対称カスターマネージドキーを使用して Data Catalog を暗号化するとき、監査レコードとログで暗号化コンテキストも使用することができ、キーがどのように使用されているか特定することができます。暗号化コンテキストは、AWS CloudTrail または Amazon CloudWatch ログが生成したログにも表示されます。

### 暗号化の有効化

AWS Glue コンソールの [データカタログの設定] または AWS CLI を使用し、AWS Glue Data Catalog オブジェクトの暗号化を有効にできます。

## Console

コンソールを使用して暗号化を有効にするには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインで、[データカタログ] を選択します。
3. [データカタログの設定] ページで、[メタデータの暗号化] チェックボックスを選択し、AWS KMS キーを選択します。

暗号化を有効にするとき、カスターマネージドキーを指定しない場合、暗号化設定では AWS マネージド KMS キーが使用されます。

4. (オプション) カスターマネージドキーを使用して Data Catalog を暗号化すると、リソースを暗号化および復号化するため、Data Catalog は IAM ロールを登録するオプションを提供します。お客様に代わって AWS Glue が引き受ける IAM ロールの許可を付与する必要があります。データを暗号化および復号する AWS KMS アクセス許可が含まれます。

Data Catalog に新しいリソースを作成すると、AWS Glue はデータを暗号化するために提供される IAM ロールを引き受けます。同様に、コンシューマがリソースにアクセスすると、AWS Glue はデータを復号化する IAM ロールを引き受けます。必要な許可を持つ IAM ロールを登録した場合、呼び出し元のプリンシパルはキーにアクセスしてデータを復号化する許可の必要がなくなります。

### Important

カスターマネージドキーを使用して Data Catalog リソースを暗号化するときに限り、KMS オペレーションを IAM ロールに委任できます。KMS ロール委任の機能は、現時点では Data Catalog リソースの暗号化に AWS マネージドキーの使用をサポートしていません。

### Warning

IAM ロールを有効にして KMS の操作を委任すると、以前に AWS マネージドキーで暗号化された Data Catalog リソースにアクセスできなくなります。

- a. お客様に代わって AWS Glue が引き受けられる IAM ロールを有効にし、データの暗号化および復号化できるようにするには、[KMS 操作を IAM ロールに委任する] オプションを選択します。
- b. 次に、IAM ロールを選択します。

IAM ロールを作成するには、「[AWS Glue の IAM ロールを作成する](#)」を参照してください。

Data Catalog にアクセスするために AWS Glue が引き受ける IAM ロールには、Data Catalog のメタデータを暗号化および復号化する許可が必要です。IAM ロールを作成し、次のインラインポリシーをアタッチできます。

- 次のポリシーを追加し、Data Catalog を暗号化および復号化する AWS KMS 許可を含めてください。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "kms:Decrypt",
 "kms:Encrypt",
 "kms:GenerateDataKey"
],
 "Resource": "arn:aws:kms:<region>:<account-id>:key/<key-id>"
 }
]
}
```

- 次に、AWS Glue サービスが IAM ロールを引き受けるには、ロールに次の信頼ポリシーを追加します。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "",
 "Effect": "Allow",
```

```
 "Principal": {
 "Service": "glue.amazonaws.com"
 },
 "Action": "sts:AssumeRole"
 }
]
```

- 次に、IAM ロールに `iam:PassRole` 許可を追加します。

```
 {
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "iam:PassRole"
],
 "Resource": [
 "arn:aws:iam::<account-id>:role/<encryption-role-name>"
]
 }
]
 }
]
```

暗号化を有効にするとき、AWS Glue が引き受ける IAM ロールを指定していない場合、Data Catalog にアクセスするプリンシパルには、次の API 操作を実行する許可が必要です。

- `kms:Decrypt`
- `kms:Encrypt`
- `kms:GenerateDataKey`

## AWS CLI

SDK または AWS CLI を使用して暗号化を有効にするには

- PutDataCatalogEncryptionSettings API オペレーションを使用します。キーが指定されていない場合、AWS Glue はカスタマーアカウントに AWS マネージド暗号化キーを使用し、Data Catalog を暗号化します。

```
aws glue put-data-catalog-encryption-settings \
 --data-catalog-encryption-settings '{
 "EncryptionAtRest": {
 "CatalogEncryptionMode": "SSE-KMS-WITH-SERVICE-ROLE",
 "SseAwsKmsKeyId": "arn:aws:kms:<region>:<account-id>:key/<key-id>",
 "CatalogEncryptionServiceRole": "arn:aws:iam::<account-
id>:role/<encryption-role-name>"
 }
 }'
'
```

暗号化を有効にすると、Data Catalog オブジェクトで作成するすべてのオブジェクトが暗号化されます。この設定をオフにすると、Data Catalog で作成するオブジェクトが暗号化されなくなります。必要な KMS 許可を使用して、Data Catalog の既存の暗号化オブジェクトに引き続きアクセスできます。

### Important

AWS KMS キーは、Data Catalog でそのキーを使用して暗号化されたすべてのオブジェクトの AWS KMS キーストアに、使用可能な状態で保持する必要があります。キーを削除すると、オブジェクトは復号化できなくなります。ある種のシナリオでは、このようにして、Data Catalog メタデータにアクセスできないようにする場合があります。

## AWS Glue 用に KMS キーの監視

Data Catalog リソースに KMS キーを使用すると、AWS CloudTrail または Amazon CloudWatch ログを使用して AWS Glue が AWS KMS に送信するリクエストを追跡できます。AWS CloudTrail

は、KMS キーで暗号化されたデータにアクセスするために AWS Glue が呼び出す KMS 操作を監視および記録します。

次の例は、Decrypt および GenerateDataKey オペレーションの AWS CloudTrail イベントです。

## Decrypt

```
{
 "eventVersion": "1.08",
 "userIdentity": {
 "type": "AssumedRole",
 "principalId": "AROAXPHTESTANDEXAMPLE:Sampleuser01",
 "arn": "arn:aws:sts::111122223333:assumed-role/Admin/Sampleuser01",
 "accountId": "111122223333",
 "accessKeyId": "AKIAIOSFODNN7EXAMPLE",
 "sessionContext": {
 "sessionIssuer": {
 "type": "Role",
 "principalId": "AROAXPHTESTANDEXAMPLE",
 "arn": "arn:aws:iam::111122223333:role/Admin",
 "accountId": "111122223333",
 "userName": "Admin"
 },
 "webIdFederationData": {},
 "attributes": {
 "creationDate": "2024-01-10T14:33:56Z",
 "mfaAuthenticated": "false"
 }
 },
 "invokedBy": "glue.amazonaws.com"
 },
 "eventTime": "2024-01-10T15:18:11Z",
 "eventSource": "kms.amazonaws.com",
 "eventName": "Decrypt",
 "awsRegion": "eu-west-2",
 "sourceIPAddress": "glue.amazonaws.com",
 "userAgent": "glue.amazonaws.com",
 "requestParameters": {
 "encryptionContext": {
 "glue_catalog_id": "111122223333"
 },
 "encryptionAlgorithm": "SYMMETRIC_DEFAULT"
 },
 "responseElements": null,
```

```

"requestID": "43b019aa-34b8-4798-9b98-ee968b2d63df",
"eventID": "d7614763-d3fe-4f84-a1e1-3ca4d2a5bbd5",
"readOnly": true,
"resources": [
 {
 "accountId": "111122223333",
 "type": "AWS::KMS::Key",
 "ARN": "arn:aws:kms:<region>:111122223333:key/<key-id>"
 }
],
"eventType": "AwsApiCall",
"managementEvent": true,
"recipientAccountId": "111122223333",
"eventCategory": "Management",
"sessionCredentialFromConsole": "true"
}

```

## GenerateDataKey

```

{
 "eventVersion": "1.08",
 "userIdentity": {
 "type": "AssumedRole",
 "principalId":
"AROAXPHTESTANDEXAMPLE:V_00_GLUE_KMS_GENERATE_DATA_KEY_111122223333",
 "arn": "arn:aws:sts::111122223333:assumed-role/Admin/
V_00_GLUE_KMS_GENERATE_DATA_KEY_111122223333",
 "accountId": "111122223333",
 "accessKeyId": "AKIAIOSFODNN7EXAMPLE",
 "sessionContext": {
 "sessionIssuer": {
 "type": "Role",
 "principalId": "AROAXPHTESTANDEXAMPLE",
 "arn": "arn:aws:iam::111122223333:role/Admin",
 "accountId": "AKIAIOSFODNN7EXAMPLE",
 "userName": "Admin"
 },
 "webIdFederationData": {},
 "attributes": {
 "creationDate": "2024-01-05T21:15:47Z",
 "mfaAuthenticated": "false"
 }
 }
 }
}

```

```
 },
 "invokedBy": "glue.amazonaws.com"
 },
 "eventTime": "2024-01-05T21:15:47Z",
 "eventSource": "kms.amazonaws.com",
 "eventName": "GenerateDataKey",
 "awsRegion": "eu-west-2",
 "sourceIPAddress": "glue.amazonaws.com",
 "userAgent": "glue.amazonaws.com",
 "requestParameters": {
 "keyId": "arn:aws:kms:eu-west-2:AKIAIOSFODNN7EXAMPLE:key/
AKIAIOSFODNN7EXAMPLE",
 "encryptionContext": {
 "glue_catalog_id": "111122223333"
 },
 "keySpec": "AES_256"
 },
 "responseElements": null,
 "requestID": "64d1783a-4b62-44ba-b0ab-388b50188070",
 "eventID": "1c73689b-2ef2-443b-aed7-8c126585ca5e",
 "readOnly": true,
 "resources": [
 {
 "accountId": "111122223333",
 "type": "AWS::KMS::Key",
 "ARN": "arn:aws:kms:eu-west-2:111122223333:key/AKIAIOSFODNN7EXAMPLE"
 }
],
 "eventType": "AwsApiCall",
 "managementEvent": true,
 "recipientAccountId": "111122223333",
 "eventCategory": "Management"
}
```

## 接続パスワードの暗号化

AWS Glue Data Catalog の接続パスワードを取得するには、GetConnection および GetConnections の API オペレーションを使用します。これらのパスワードは Data Catalog 接続に保存され、AWS Glue から Java Database Connectivity (JDBC) データストアに接続される際に使用されます。接続が作成または更新された際、Data Catalog 設定のオプションによって、パスワードが暗号化されたかどうか判断され、暗号化されている場合には、指定された AWS Key Management Service (AWS KMS) キーが判断されます。

AWS Glue コンソールの [Data catalog settings] (データカタログ設定) ページで、このオプションを有効にできます。

接続パスワードを暗号化するには

1. AWS Management Console にサインインし、AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) を開きます。
2. ナビゲーションペインで [設定] を選択します。
3. [Data catalog settings] (データカタログ設定) ページで、[Encrypt connection passwords] (接続パスワードの暗号化) を選択し、AWS KMS を選択します。

### Important

AWS Glue は、対称カスタマーマスターキー (CMK) のみをサポートしています。[AWS KMS key] (KMS キー) リストには、対称キーのみが表示されます。ただし、[Choose a AWS KMS key ARN] (KMS キー ARN を選択します) を選択した場合、任意のキータイプの ARN をコンソールで入力します。対称キーには ARN だけを入力するようにしてください。

詳細については、「[データカタログの設定](#)」を参照してください。

## AWS Glue によって書き込まれたデータの暗号化

セキュリティ設定は、AWS Glue が使用できるセキュリティプロパティのセットです。セキュリティ構成を使用して、保管時にデータを暗号化できます。以下のシナリオでは、セキュリティ設定を使用できるいくつかの方法を示します。

- 暗号化された Amazon CloudWatch Logs を書き込む AWS Glue クローラーにセキュリティ設定をアタッチします。クローラーにセキュリティ設定のアタッチの詳細については、「[the section called “ステップ 3: セキュリティ設定を構成する”](#)」を参照してください。
- 抽出、変換、ロード (ETL) ジョブにセキュリティ設定をアタッチして、暗号化された Amazon Simple Storage Service (Amazon S3) ターゲット、および暗号化された CloudWatch Logs を書き込みます。
- ETL ジョブにセキュリティ設定をアタッチして、暗号化された Amazon S3 データとしてジョブブックマークを書き込みます。
- 開発エンドポイントにセキュリティ設定をアタッチして、暗号化された Amazon S3 ターゲットを書き込みます。

#### Important

現時点では、セキュリティ設定は、ETL ジョブのパラメータとして渡されるサーバー側暗号化 (SSE-S3) 設定をすべてオーバーライドします。したがって、ジョブにセキュリティ設定と SSE-S3 の両方が関連付けられている場合、SSE-S3 パラメータは無視されます。

セキュリティ設定の詳細については、「[AWS Glue コンソールでのセキュリティ設定の使用](#)」を参照してください。

#### トピック

- [セキュリティ設定を使用するための AWS Glue のセットアップ](#)
- [VPC ジョブとクローラーの AWS KMS へのルートを作成する](#)
- [AWS Glue コンソールでのセキュリティ設定の使用](#)

#### セキュリティ設定を使用するための AWS Glue のセットアップ

以下の手順を実行して、セキュリティ設定を使用するように AWS Glue 環境をセットアップします。

1. AWS Glue クローラーに渡される IAM ロールおよび CloudWatch Logs を暗号化するジョブへの AWS KMS アクセス許可を付与するため、AWS Key Management Service (AWS KMS) キーを作成または更新します。詳しくは、Amazon CloudWatch Logs ユーザーガイドの「[AWS KMS を使用して CloudWatch Logs のログデータを暗号化する](#)」を参照してください。

次の例では、`"role1"`、`"role2"`、および `"role3"` が、クローラーおよびジョブに渡される IAM ロールです。

```
{
 "Effect": "Allow",
 "Principal": { "Service": "logs.region.amazonaws.com",
 "AWS": [
 "role1",
 "role2",
 "role3"
] },
 "Action": [
 "kms:Encrypt*",
 "kms:Decrypt*",
 "kms:ReEncrypt*",
 "kms:GenerateDataKey*",
 "kms:Describe*"
],
 "Resource": "*"
}
```

キーを使用して CloudWatch Logs を暗号化する場合、`"Service": "logs.region.amazonaws.com"` と表示される Service ステートメントが必要です。

2. 使用する前に、AWS KMS キーが ENABLED であることを確認します。

#### Note

データレイクフレームワークとして Iceberg を使用している場合、Iceberg テーブルにはサーバー側の暗号化を有効にする独自のメカニズムがあります。AWS Glue のセキュリティ設定に加えて、これらの設定を有効にする必要があります。Iceberg テーブルでサーバー側の暗号化を有効にするには、「[Iceberg ドキュメント](#)」のガイダンスを確認してください。

## VPC ジョブとクローラーの AWS KMS へのルートを作成する

インターネットを介さずに、仮想プライベートクラウド (VPC) のプライベートエンドポイントから AWS KMS に直接接続できます。VPC エンドポイントを使用すると、VPC と AWS KMS の間の通信は完全に AWS ネットワーク内で実施されます。

VPC 内に AWS KMS VPC エンドポイントを作成することができます。この手順を実行しないと、ジョブまたはクローラーがジョブの kms timeout またはクローラーの internal service exception で失敗する可能性があります。詳細な手順については、AWS Key Management Service デベロッパーガイドの「[VPC エンドポイントを介した AWS KMS への接続](#)」を参照してください。

これらの手順に従って、[VPC コンソール](#)で次の操作を行う必要があります。

- [Enable Private DNS name] (プライベート DNS 名を有効にする) を選択します。
- Java Database Connectivity (JDBC) にアクセスするジョブまたはクローラーに使用する [Security group] (セキュリティグループ) (自己参照ルールを持つ) を選択します。AWS Glue 接続の詳細については、「[データへの接続](#)」を参照してください。

JDBC データストアにアクセスするクローラーまたはジョブにセキュリティ設定を追加する場合、AWS Glue には AWS KMS エンドポイントへのルートが必要です。ネットワークアドレス変換 (NAT) ゲートウェイまたは AWS KMS VPC エンドポイントを使用してルートを指定できます。NAT ゲートウェイを作成するには、Amazon VPC ユーザーガイドの「[NAT ゲートウェイ](#)」を参照してください。

AWS Glue コンソールでのセキュリティ設定の使用

#### Warning

AWS Glue セキュリティ設定は、現在、Ray ジョブではサポートされていません。

AWS Glue の [security configuration] (セキュリティ設定) には、暗号化されたデータを書き込むときに必要なプロパティが含まれています。AWS Glue コンソールでセキュリティ設定を作成し、クローラー、ジョブ、および開発エンドポイントによって使用される暗号化プロパティを指定します。

作成したすべてのセキュリティ設定を表示するには、<https://console.aws.amazon.com/glue/> で AWS Glue コンソールを開き、ナビゲーションペインで [Security configurations] (セキュリティ設定) をクリックします。

[セキュリティ設定] には、設定に関する次のプロパティが表示されます。

#### 名前

設定を作成したときに指定した一意の名前です。名前に使用できるのは、文字 (A~Z)、数字 (0~9)、ハイフン (-)、またはアンダースコア (\_) で、長さは 255 文字までです。

## Amazon S3 暗号化を有効にする

オンの場合、データカタログのメタデータ保存に Amazon Simple Storage Service (Amazon S3) 暗号化モード (SSE-KMS、SSE-S3 など) が有効になります。

## Amazon CloudWatch Logs 暗号化を有効にする

オンの場合、Amazon CloudWatch にログを書き込む際に Amazon S3 暗号化モード (SSE-KMS など) が有効になります。

## 詳細設定: ジョブのブックマーク暗号化を有効にする

オンの場合、ジョブがブックマークされた際に Amazon S3 暗号化モード (SSE-KMS など) が有効になります。

コンソールの [セキュリティ設定] セクションで設定を追加または削除できます。設定の詳細を表示するには、リスト内の設定名を選択します。詳細には、設定の作成時に定義した情報が含まれます。

## セキュリティ設定の追加

セキュリティ設定を追加するには、AWS Glue コンソールの [セキュリティ設定] ページで [セキュリティ設定の追加] を選択します。

## Add security configuration

Choose encryption and permission options for your accounts data catalog.

### Security configuration properties

**Name**

Name may contain letters (A-Z), numbers (0-9), hyphens (-), or underscores (\_), and can be up to 255 characters long.

### Encryption settings

Enable and choose options for at-rest encryption.

**Enable S3 encryption**  
Enable at-rest encryption for metadata stored in the data catalog.

**Enable CloudWatch logs encryption**  
Enable at-rest encryption when writing logs to Amazon CloudWatch.

▼ **Advanced settings**

**Enable job bookmark encryption**  
Enable at-rest encryption of job bookmark.

**Cancel** **Save**

### セキュリティ設定プロパティ

一意のセキュリティ設定名を入力します。名前に使用できるのは、文字 (A~Z)、数字 (0~9)、ハイフン (-)、またはアンダースコア (\_) で、長さは 255 文字までです。

### 暗号化設定

Amazon CloudWatch のログと Amazon S3 のデータカタログに保存されているメタデータの保管時の暗号化を有効にできます。AWS Glue コンソールで AWS Key Management Service (AWS KMS) キーを使用したデータとメタデータの暗号化を設定するには、コンソールユーザーにポリシーを追加します。このポリシーでは次の例のように、許可されているリソースを、Amazon S3 データストア

を暗号化するために使用されるキーの Amazon リソースネーム (ARN) として、指定する必要があります。

```
{
 "Version": "2012-10-17",
 "Statement": {
 "Effect": "Allow",
 "Action": [
 "kms:GenerateDataKey",
 "kms:Decrypt",
 "kms:Encrypt"],
 "Resource": "arn:aws:kms:region:account-id:key/key-id"
 }
}
```

#### Important

セキュリティ設定がクローラまたはジョブにアタッチされている場合、渡される IAM ロールには AWS KMS アクセス許可が必要です。詳細については、[AWS Glue によって書き込まれたデータの暗号化](#) を参照してください。

設定を定義する場合、次のプロパティの値を指定できます。

#### S3 暗号化を有効にする

Amazon S3 データを書き込む場合は、Amazon S3 管理キー (SSE-S3) によるサーバー側の暗号化、または AWS KMS 管理キー (SSE-KMS) によるサーバー側の暗号化を使用します。このフィールドはオプションです。Amazon S3 へのアクセス許可を有効にするには、AWS KMS キーを選択するか、または [Enter a key ARN] (キーの ARN を入力) を選択してキーの ARN を指定します。arn:aws:kms:*region*:*account-id*:key/*key-id* の形式で ARN を入力します。arn:aws:kms:*region*:*account-id*:alias/*alias-name* など、キーのエイリアスで ARN を指定することもできます。

ジョブで Spark UI を有効にすると、Amazon S3 にアップロードされた Spark UI ログファイルに同じ暗号化が適用されます。

#### Important

AWS Glue は、対称カスタマーマスターキー (CMK) のみをサポートしています。[AWS KMS key] (KMS キー) リストには、対称キーのみが表示されます。ただし、[Choose a

AWS KMS key ARN] (KMS キー ARN を選択します) を選択した場合、任意のキータイプの ARN をコンソールで入力します。対称キーには ARN だけを入力するようにしてください。

## CloudWatch Logs 暗号化を有効にする

サーバー側 (SSE-KMS) 暗号化が CloudWatch Logs の暗号化に使用されます。このフィールドはオプションです。有効にするには、AWS KMS キーを選択するか、または [Enter a key ARN] (キーの ARN を入力) を選択してキーの ARN を指定します。arn:aws:kms:region:account-id:key/key-id の形式で ARN を入力します。arn:aws:kms:region:account-id:alias/alias-name など、キーのエイリアスで ARN を指定することもできます。

### 詳細設定: ジョブブックマークの暗号化

クライアント側 (CSE-KMS) 暗号化は、ジョブのブックマークを暗号化するために使用されます。このフィールドはオプションです。ブックマークデータは、Amazon S3 に保存するために送信される前に暗号化されます。有効にするには、AWS KMS キーを選択するか、または [Enter a key ARN] (キーの ARN を入力) を選択してキーの ARN を指定します。arn:aws:kms:region:account-id:key/key-id の形式で ARN を入力します。arn:aws:kms:region:account-id:alias/alias-name など、キーのエイリアスで ARN を指定することもできます。

詳細については、Amazon Simple Storage Service コンソールユーザーガイドの以下のトピックを参照してください。

- SSE-S3 の詳細については、「[Amazon S3 で管理された暗号化キーによるサーバー側の暗号化 \(SSE-S3\) を使用したデータの保護](#)」を参照してください。
- SSE-KMS の詳細については、「[AWS KMS keys によるサーバー側の暗号化を使用したデータの保護](#)」を参照してください。
- CSE-KMS の詳細については、「[AWS KMS に保存されている KMS キーの使用](#)」を参照してください。

## 送信中の暗号化

AWS は転送中のデータに対して Transport Layer Security (TLS) 暗号化を提供します。AWS Glue の [\[security configurations\]](#) (セキュリティ設定) を使用して、クローラー、ETL ジョブ、開発エンドポイ

ントの暗号化を設定できます。Data Catalog の設定を介して、AWS Glue Data Catalog 暗号化を有効にできます。

2018 年 9 月 4 日現在、(AWS KMSbring your own keyと server-side encryption) が AWS Glue ETLおよび AWS Glue Data Catalog がサポートされています。

## FIPS 準拠

コマンドラインインターフェイスまたは API により AWS にアクセスするときに FIPS 140-2 検証済み暗号化モジュールが必要な場合は、FIPS エンドポイントを使用します。利用可能な FIPS エンドポイントの詳細については、「[連邦情報処理規格 \(FIPS\) 140-2](#)」を参照してください。

## キーの管理

AWS Glue で AWS Identity and Access Management (IAM) を使用して、ユーザー、AWS リソース、グループ、ロール、およびアクセスや拒否などに関する詳細なポリシーを定義できます。

組織のニーズに応じて、リソースベースのポリシーとアイデンティティベースのポリシーの両方を使用してメタデータへのアクセスを定義できます。リソースベースのポリシーでは、リソースへのアクセスを許可または拒否するプリンシパルを一覧表示し、クロスアカウントアクセスなどのポリシーを設定できます。アイデンティティポリシーは、IAM 内の特定のユーザー、グループ、およびロールにアタッチします。

詳細な例については、AWS ビッグデータブログ「[リソースレベルの IAM アクセス許可とリソースベースのポリシーで、AWS Glue Data Catalog へのアクセスを制限する](#)」を参照してください。

ポリシーの詳細なアクセス部分については、Resource 句内で定義します。この部分では、アクションを実行する対象の AWS Glue Data Catalog オブジェクトと、このオペレーションで返される結果のオブジェクトの両方を定義します。

開発エンドポイントは、AWS Glue スクリプトの開発およびテストに使用できる環境です。開発エンドポイントの SSH キーを追加、削除、または変更できます。

2018 年 9 月 4 日現在、AWS KMS ETLおよび AWS Glue の (AWS Glue Data Catalogbring your own key と server-side encryption) がサポートされています。

## AWS の他のサービスへの AWS Glue の依存関係

ユーザーが AWS Glue コンソールを使用するには、AWS アカウントで AWS Glue リソースの使用を許可する最小限のアクセス許可セットが必要です。これらの AWS Glue アクセス許可に加えて、コンソールでは次のサービスからのアクセス許可が必要になります。

- ログを表示するための Amazon CloudWatch Logs のアクセス許可。
- ロールをリスト化して渡すための AWS Identity and Access Management (IAM) のアクセス許可。
- スタックを操作する AWS CloudFormation のアクセス許可。
- Virtual Private Cloud (VPC)、サブネット、セキュリティグループ、インスタンスなどのオブジェクトをリストする Amazon Elastic Compute Cloud (Amazon EC2) のアクセス許可 (ジョブ、クローラーの実行時や、開発エンドポイントの作成時の VPC などの Amazon EC2 項目をセットアップするため)。
- バケットとオブジェクトをリストし、スクリプトを取得して保存する Amazon Simple Storage Service (Amazon S3) のアクセス許可。
- クラスターでの作業に必要な Amazon Redshift のアクセス許可。
- インスタンスをリスト化するための Amazon Relational Database Service (Amazon RDS) のアクセス許可。

## 開発エンドポイント

開発エンドポイントは、AWS Glue スクリプトの作成およびテストに使用できる環境です。AWS Glue を使用して、開発エンドポイントを作成、編集、削除できます。作成されたすべての開発エンドポイントを一覧表示できます。開発エンドポイントの SSH キーを追加、削除、または変更できます。開発エンドポイントで使用するノートブックを作成することもできます。

開発環境をプロビジョニングするための設定値を提供します。これらの値は、開発エンドポイントに安全にアクセスし、エンドポイントからデータストアにアクセスできるようにネットワークを設定する方法を AWS Glue に指定します。次に、開発エンドポイントに接続するノートブックを作成できます。ノートブックを使用して ETL スクリプトを作成し、テストします。

AWS Glue の ETL ジョブの実行に使用する IAM ロールと同様のアクセス許可を持つ AWS Identity and Access Management (IAM) ロールを使用します。Virtual Private Cloud (VPC)、サブネット、セキュリティグループを使用して、データリソースに安全に接続できる開発エンドポイントを作成します。SSH を使用して開発環境に接続するための SSH キーペアを生成します。

JDBC を通じてデータセットにアクセスするために使用できる VPC 内で、Amazon S3 のデータ用の開発エンドポイントを作成できます。

ローカルマシンに Jupyter Notebook クライアントをインストールし、これを使用して開発エンドポイントで ETL スクリプトをデバッグおよびテストできます。または Sagemaker ノートブックを使用して、AWS 上の JupyterLab で ETL スクリプトを記述できます。「[開発エンドポイントで Amazon SageMaker ノートブックを使用する](#)」を参照してください。

AWS Glue は、aws-glue-dev-endpoint をプレフィックスとする名前を Amazon EC2 インスタンスにタグ付けします。

開発エンドポイントにノートブックサーバーをセットアップして、AWS Glue 拡張機能で PySpark を実行できます。

## AWS Glue の Identity and Access Management

AWS Identity and Access Management (IAM) は、管理者が AWS リソースへのアクセスを安全に制御 AWS のサービス するのに役立つです。IAM 管理者は、誰を認証 (サインイン) し、誰に AWS Glue リソースの使用を承認する (アクセス許可を付与する) かを制御します。IAM は、追加料金なしで AWS のサービス 使用できる です。

### Note

AWS Glue Data Catalog 内のデータへのアクセスは、AWS Glue メソッドまたは AWS Lake Formation 許可を使用して許可できます。AWS Identity and Access Management (IAM) ポリシーを使用して、AWS Glue メソッドできめ細かなアクセスコントロールを設定できます。Lake Formation では、リレーショナルデータベースシステムの GRANT/REVOKE コマンドに似た、よりシンプルな GRANT/REVOKE アクセス許可モデルを使用します。

このセクションでは、AWS Glue メソッドの使用方法を説明します Lake Formation の付与の詳細については、AWS Lake Formation デベロッパーガイドの「[Granting Lake Formation Permissions](#)」(Lake Formation 許可の付与) を参照してください。

### トピック

- [対象者](#)
- [アイデンティティを使用した認証](#)
- [ポリシーを使用したアクセスの管理](#)
- [AWS Glue と IAM の連携方法](#)
- [AWS Glue の IAM アクセス許可の設定](#)
- [AWS Glue のアクセスコントロールポリシーの例](#)
- [AWS Glue の管理ポリシー](#)
- [AWS Glue リソース ARN の指定](#)

- [クロスアカウントアクセス許可の付与](#)
- [AWS Glue のアイデンティティとアクセスのトラブルシューティング](#)

## 対象者

AWS Identity and Access Management (IAM) の使用方法は、AWS Glue で行う作業によって異なります。

サービスユーザー – ジョブを実行するために AWS Glue サービスを使用する場合、管理者から必要な認証情報とアクセス許可が与えられます。さらに多くの AWS Glue 機能を使用して作業を行う場合は、追加のアクセス許可が必要になることがあります。アクセスの管理方法を理解しておく、管理者に適切な許可をリクエストするうえで役立ちます。AWS Glue の機能にアクセスできないときは、「[AWS Glue のアイデンティティとアクセスのトラブルシューティング](#)」を参照してください。

サービス管理者 – 社内の AWS Glue リソースを担当している場合は、通常、AWS Glue へのフルアクセスがあります。サービスユーザーがどの AWS Glue 機能やリソースにアクセスするかを決めるのは管理者の仕事です。その後、IAM 管理者にリクエストを送信して、サービスユーザーの権限を変更する必要があります。このページの情報を点検して、IAM の基本概念を理解してください。会社で AWS Glue で IAM を使用方法の詳細については、「」を参照してください [AWS Glue と IAM の連携方法](#)。

IAM 管理者 – IAM 管理者は、AWS Glue へのアクセスを管理するポリシーの作成方法の詳細について確認する場合があります。IAM で使用できる AWS Glue アイデンティティベースのポリシーの例を表示するには、「」を参照してください [AWS Glue のアイデンティティベースポリシーの例](#)。

## アイデンティティを使用した認証

認証とは、ID 認証情報 AWS を使用して にサインインする方法です。として、IAM ユーザーとして AWS アカウントのルートユーザー、または IAM ロールを引き受けて認証 ( にサインイン AWS) される必要があります。

ID ソースを介して提供された認証情報を使用して、フェデレーテッド ID AWS として にサインインできます。AWS IAM Identity Center ( IAM Identity Center) ユーザー、会社のシングルサインオン認証、Google または Facebook の認証情報は、フェデレーション ID の例です。フェデレーテッド ID としてサインインする場合、IAM ロールを使用して、前もって管理者により ID フェデレーションが設定されています。フェデレーション AWS を使用して にアクセスすると、間接的にロールを引き受けることとなります。

ユーザーのタイプに応じて、AWS Management Console または AWS アクセスポータルにサインインできます。へのサインインの詳細については AWS、「ユーザーガイド」の「[へのサインイン AWS アカウント](#)方法AWS サインイン」を参照してください。

AWS プログラムでにアクセスする場合、は Software Development Kit (SDK) とコマンドラインインターフェイス (CLI) AWS を提供し、認証情報を使用してリクエストに暗号で署名します。AWS ツールを使用しない場合は、リクエストに自分で署名する必要があります。推奨される方法を使用してリクエストを自分で署名する方法の詳細については、IAM [ユーザーガイドの API AWS リクエスト](#)の署名を参照してください。

使用する認証方法を問わず、追加セキュリティ情報の提供をリクエストされる場合もあります。例えば、AWS では、多要素認証 (MFA) を使用してアカウントのセキュリティを向上させることをお勧めします。詳細については、「AWS IAM Identity Center ユーザーガイド」の「[Multi-factor authentication](#)」(多要素認証) および「IAM ユーザーガイド」の「[AWSでの多要素認証 \(MFA\) の使用](#)」を参照してください。

## AWS アカウント ルートユーザー

を作成するときは AWS アカウント、アカウント内のすべての AWS のサービス およびリソースへの完全なアクセス権を持つ1つのサインインアイデンティティから始めます。この ID は AWS アカウント ルートユーザーと呼ばれ、アカウントの作成に使用した E メールアドレスとパスワードでサインインすることでアクセスできます。日常的なタスクには、ルートユーザーを使用しないことを強くお勧めします。ルートユーザーの認証情報は保護し、ルートユーザーでしか実行できないタスクを実行するときに使用します。ルートユーザーとしてサインインする必要があるタスクの完全なリストについては、IAM ユーザーガイドの[ルートユーザー認証情報が必要なタスク](#)を参照してください。

## フェデレーティッドアイデンティティ

ベストプラクティスとして、管理者アクセスを必要とするユーザーを含む人間のユーザーに、一時的な認証情報を使用してにアクセスするための ID プロバイダーとのフェデレーションの使用を要求 AWS のサービスします。

フェデレーティッド ID は、エンタープライズユーザーディレクトリ、ウェブ ID プロバイダー、AWS Directory Service、アイデンティティセンターディレクトリのユーザー、または ID ソースを通じて提供された認証情報 AWS のサービス を使用してにアクセスするユーザーです。フェデレーティッド ID がにアクセスすると AWS アカウント、ロールを引き受け、ロールは一時的な認証情報を提供します。

アクセスを一元管理する場合は、AWS IAM Identity Centerを使用することをお勧めします。IAM Identity Center でユーザーとグループを作成することも、独自の ID ソース内のユーザーとグルー

プのセットに接続して同期して、すべての AWS アカウント とアプリケーションで使用することもできます。IAM Identity Center の詳細については、「AWS IAM Identity Center ユーザーガイド」の「[What is IAM Identity Center?](#)」(IAM Identity Center とは) を参照してください。

## IAM ユーザーとグループ

[IAM ユーザー](#)は、単一のユーザーまたはアプリケーションに対して特定のアクセス許可 AWS アカウント を持つ 内のアイデンティティです。可能であれば、パスワードやアクセスキーなどの長期的な認証情報を保有する IAM ユーザーを作成する代わりに、一時的な認証情報を使用することをお勧めします。ただし、IAM ユーザーでの長期的な認証情報が必要な特定のユースケースがある場合は、アクセスキーをローテーションすることをお勧めします。詳細については、IAM ユーザーガイドの[長期的な認証情報を必要とするユースケースのためにアクセスキーを定期的にローテーションする](#)を参照してください。

[IAM グループ](#)は、IAM ユーザーの集団を指定するアイデンティティです。グループとしてサインインすることはできません。グループを使用して、複数のユーザーに対して一度に権限を指定できます。多数のユーザーグループがある場合、グループを使用することで権限の管理が容易になります。例えば、IAMAdmins という名前のグループを設定して、そのグループに IAM リソースを管理する許可を与えることができます。

ユーザーは、ロールとは異なります。ユーザーは 1 人の人または 1 つのアプリケーションに一意に関連付けられますが、ロールはそれを必要とする任意の人が引き受けるようになっています。ユーザーには永続的な長期の認証情報がありますが、ロールでは一時的な認証情報が提供されます。詳細については、「IAM ユーザーガイド」の「[IAM ユーザー \(ロールではなく\) の作成が適している場合](#)」を参照してください。

## IAM ロール

[IAM ロール](#)は、特定のアクセス許可 AWS アカウント を持つ 内のアイデンティティです。これは IAM ユーザーに似ていますが、特定のユーザーには関連付けられていません。ロール を切り替える AWS Management Console ことで、[で IAM ロール](#)を一時的に引き受けることができます。ロールを引き受けるには、または AWS API AWS CLI オペレーションを呼び出すか、カスタム URL を使用します。ロールを使用する方法の詳細については、「IAM ユーザーガイド」の「[IAM ロールの使用](#)」を参照してください。

IAM ロールと一時的な認証情報は、次の状況で役立ちます:

- フェデレーションユーザーアクセス – フェデレーティッド ID に許可を割り当てるには、ロールを作成してそのロールの許可を定義します。フェデレーティッド ID が認証されると、その ID はロールに関連付けられ、ロールで定義されている許可が付与されます。フェデレーションの詳細

については、「IAM ユーザーガイド」の「[Creating a role for a third-party Identity Provider](#)」(サードパーティーアイデンティティプロバイダー向けロールの作成)を参照してください。IAM Identity Center を使用する場合は、許可セットを設定します。アイデンティティが認証後にアクセスできるものを制御するため、IAM Identity Center は、権限セットを IAM のロールに関連付けます。アクセス許可セットの詳細については、「AWS IAM Identity Center ユーザーガイド」の「[アクセス許可セット](#)」を参照してください。

- 一時的な IAM ユーザー権限 - IAM ユーザーまたはロールは、特定のタスクに対して複数の異なる権限を一時的に IAM ロールで引き受けることができます。
- クロスアカウントアクセス - IAM ロールを使用して、自分のアカウントのリソースにアクセスすることを、別のアカウントの人物(信頼済みプリンシパル)に許可できます。クロスアカウントアクセス権を付与する主な方法は、ロールを使用することです。ただし、一部の AWS サービス、(ロールをプロキシとして使用する代わりに)ポリシーをリソースに直接アタッチできます。クロスアカウントアクセスにおけるロールとリソースベースのポリシーの違いについては、「IAM ユーザーガイド」の「[IAM でのクロスアカウントのリソースへのアクセス](#)」を参照してください。
- クロスサービスアクセス — 一部の AWS サービスは、他の AWS サービスを使用します。例えば、あるサービスで呼び出しを行うと、通常そのサービスによって Amazon EC2 でアプリケーションが実行されたり、Amazon S3 にオブジェクトが保存されたりします。サービスでは、呼び出し元プリンシパルの許可、サービスロール、またはサービスリンクロールを使用してこれを行う場合があります。
- 転送アクセスセッション (FAS) - IAM ユーザーまたはロールを使用してアクションを実行する場合 AWS、ユーザーはプリンシパルと見なされます。一部のサービスを使用する際に、アクションを実行することで、別のサービスの別のアクションがトリガーされることがあります。FAS は、呼び出すプリンシパルのアクセス許可を AWS のサービス、ダウンストリームサービス AWS のサービスへのリクエストのリクエストと組み合わせて使用します。FAS リクエストは、サービスが他の AWS のサービスまたはリソースとのやり取りを完了する必要があるリクエストを受け取った場合にのみ行われます。この場合、両方のアクションを実行するためのアクセス許可が必要です。FAS リクエストを行う際のポリシーの詳細については、「[転送アクセスセッション](#)」を参照してください。
- サービスロール - サービスがユーザーに代わってアクションを実行するために引き受ける [IAM ロール](#)です。IAM 管理者は、IAM 内からサービスロールを作成、変更、削除できます。詳細については、「IAM ユーザーガイド」の「[AWS のサービスにアクセス許可を委任するロールの作成](#)」を参照してください。
- サービスにリンクされたロール - サービスにリンクされたロールは、にリンクされたサービスロールの一種です AWS のサービス。サービスは、ユーザーに代わってアクションを実行する

ロールを引き受けることができます。サービスにリンクされたロールは [IAM ユーザーガイド](#) に表示され、AWS アカウント、サービスによって所有されます。IAM 管理者は、サービスにリンクされたロールのアクセス許可を表示できますが、編集することはできません。

- Amazon EC2 で実行されているアプリケーション – IAM ロールを使用して、EC2 インスタンスで実行され、AWS CLI または AWS API リクエストを行うアプリケーションの一時的な認証情報を管理できます。これは、EC2 インスタンス内でのアクセスキーの保存に推奨されます。AWS ロールを EC2 インスタンスに割り当て、そのすべてのアプリケーションで使用できるようにするには、インスタンスにアタッチされたインスタンスプロファイルを作成します。インスタンスプロファイルにはロールが含まれ、EC2 インスタンスで実行されるプログラムは一時的な認証情報を取得できます。詳細については、IAM ユーザーガイドの [Amazon EC2 インスタンスで実行されるアプリケーションに IAM ロールを使用して許可を付与する](#) を参照してください。

IAM ロールと IAM ユーザーのどちらを使用するかについては、IAM ユーザーガイドの [\(IAM ユーザーではなく\) IAM ロールをいつ作成したら良いのか?](#) を参照してください。

## ポリシーを使用したアクセスの管理

でアクセスを制御する AWS には、ポリシーを作成し、AWS ID またはリソースにアタッチします。ポリシーは AWS、アイデンティティまたはリソースに関連付けられているときにアクセス許可を定義するオブジェクトです。は、プリンシパル (ユーザー、ルートユーザー、またはロールセッション) がリクエストを行うときに、これらのポリシー AWS を評価します。ポリシーでの権限により、リクエストが許可されるか拒否されるかが決まります。ほとんどのポリシーは JSON ドキュメント AWS として保存されます。JSON ポリシードキュメントの構造と内容の詳細については、IAM ユーザーガイドの [JSON ポリシー概要](#) を参照してください。

管理者は AWS JSON ポリシーを使用して、誰が何にアクセスできるかを指定できます。つまり、どのプリンシパルがどんなリソースにどんな条件でアクションを実行できるかということです。

デフォルトでは、ユーザーやロールに権限はありません。IAM 管理者は、リソースで必要なアクションを実行するための権限をユーザーに付与する IAM ポリシーを作成できます。その後、管理者はロールに IAM ポリシーを追加し、ユーザーはロールを引き継ぐことができます。

IAM ポリシーは、オペレーションの実行方法を問わず、アクションの許可を定義します。例えば、iam:GetRole アクションを許可するポリシーがあるとします。そのポリシーを持つユーザーは、AWS Management Console、AWS CLI または AWS API からロール情報を取得できます。

## アイデンティティベースのポリシー

アイデンティティベースポリシーは、IAM ユーザーグループ、ユーザーのグループ、ロールなど、アイデンティティにアタッチできる JSON 許可ポリシードキュメントです。これらのポリシーは、ユーザーとロールが実行できるアクション、リソース、および条件をコントロールします。アイデンティティベースのポリシーを作成する方法については、IAM ユーザーガイドの[IAM ポリシーの作成](#)を参照してください。

アイデンティティベースのポリシーは、さらにインラインポリシーまたはマネージドポリシーに分類できます。インラインポリシーは、単一のユーザー、グループ、またはロールに直接埋め込まれています。管理ポリシーは、内の複数のユーザー、グループ、ロールにアタッチできるスタンドアロンポリシーです AWS アカウント。管理ポリシーには、AWS 管理ポリシーとカスタマー管理ポリシーが含まれます。マネージドポリシーまたはインラインポリシーのいずれかを選択する方法については、IAM ユーザーガイドの[マネージドポリシーとインラインポリシーの比較](#)を参照してください。

## リソースベースのポリシー

リソースベースのポリシーは、リソースに添付する JSON ポリシードキュメントです。リソースベースのポリシーには例として、IAM ロールの信頼ポリシー や Amazon S3 バケットポリシー があげられます。リソースベースのポリシーをサポートするサービスでは、サービス管理者はポリシーを使用して特定のリソースへのアクセスを制御できます。ポリシーがアタッチされているリソースの場合、指定されたプリンシパルがそのリソースに対して実行できるアクションと条件は、ポリシーによって定義されます。リソースベースのポリシーでは、[プリンシパルを指定する](#)必要があります。プリンシパルには、アカウント、ユーザー、ロール、フェデレーテッドユーザー、またはを含めることができます AWS のサービス。

リソースベースのポリシーは、そのサービス内にあるインラインポリシーです。リソースベースのポリシーでは、IAM の AWS マネージドポリシーを使用できません。

## アクセスコントロールリスト (ACL)

アクセスコントロールリスト (ACL) は、どのプリンシパル (アカウントメンバー、ユーザー、またはロール) がリソースにアクセスするための許可を持つかを制御します。ACL はリソースベースのポリシーに似ていますが、JSON ポリシードキュメント形式は使用しません。

Amazon S3、AWS WAF、および Amazon VPC は、ACLs。ACL の詳細については、Amazon Simple Storage Service デベロッパーガイドの[アクセスコントロールリスト \(ACL\) の概要](#)を参照してください。

## その他のポリシータイプ

AWS は、一般的ではない追加のポリシータイプをサポートします。これらのポリシータイプでは、より一般的なポリシータイプで付与された最大の権限を設定できます。

- **アクセス許可の境界** - アクセス許可の境界は、アイデンティティベースのポリシーによって IAM エンティティ (IAM ユーザーまたはロール) に付与できる権限の上限を設定する高度な機能です。エンティティにアクセス許可の境界を設定できます。結果として得られる権限は、エンティティのアイデンティティベースポリシーとそのアクセス許可の境界の共通部分になります。Principal フィールドでユーザーまたはロールを指定するリソースベースのポリシーでは、アクセス許可の境界は制限されません。これらのポリシーのいずれかを明示的に拒否した場合、権限は無効になります。アクセス許可の境界の詳細については、IAM ユーザーガイドの[IAM エンティティのアクセス許可の境界](#)を参照してください。
- **サービスコントロールポリシー (SCPs)** - SCPs は、 の組織または組織単位 (OU) に対する最大アクセス許可を指定する JSON ポリシーです AWS Organizations。AWS Organizations は、AWS アカウント ビジネスが所有する複数の をグループ化して一元管理するためのサービスです。組織内のすべての機能を有効にすると、サービスコントロールポリシー (SCP) を一部またはすべてのアカウントに適用できます。SCP は、各 を含むメンバーアカウントのエンティティのアクセス許可を制限します AWS アカウントのルートユーザー。Organizations と SCP の詳細については、AWS Organizations ユーザーガイドの「[SCP の仕組み](#)」を参照してください。
- **セッションポリシー** - セッションポリシーは、ロールまたはフェデレーションユーザーの一時的なセッションをプログラムで作成する際にパラメータとして渡す高度なポリシーです。結果としてセッションの権限は、ユーザーまたはロールのアイデンティティベースポリシーとセッションポリシーの共通部分になります。また、リソースベースのポリシーから権限が派生する場合もあります。これらのポリシーのいずれかを明示的に拒否した場合、権限は無効になります。詳細については、IAM ユーザーガイドの[セッションポリシー](#)を参照してください。

## 複数のポリシータイプ

1 つのリクエストに複数のタイプのポリシーが適用されると、結果として作成される権限を理解するのがさらに難しくなります。複数のポリシータイプが関与する場合にリクエストを許可するかどうか AWS を決定する方法については、IAM ユーザーガイドの「[ポリシー評価ロジック](#)」を参照してください。

## AWS Glue と IAM の連携方法

IAM を使用して AWS Glue へのアクセスを管理する前に、AWS Glue で使用できる IAM 機能について学びます。

### AWS Glue で使用できる IAM の機能

IAM 機能	AWS Glue のサポート
<a href="#">アイデンティティベースのポリシー</a>	あり
<a href="#">リソースベースのポリシー</a>	部分的
<a href="#">ポリシーアクション</a>	あり
<a href="#">ポリシーリソース</a>	はい
<a href="#">ポリシー条件キー (サービス固有)</a>	はい
<a href="#">ACL</a>	なし
<a href="#">ABAC (ポリシー内のタグ)</a>	部分的
<a href="#">一時的な認証情報</a>	あり
<a href="#">プリンシパル権限</a>	いいえ
<a href="#">サービスロール</a>	あり
<a href="#">サービスリンクロール</a>	なし

AWS Glue およびその他の AWS のサービスがほとんどの IAM 機能と連携する方法の概要を把握するには、「IAM ユーザーガイド」の[AWS 「IAM と連携する のサービス」](#)を参照してください。

### AWS Glue のアイデンティティベースのポリシー

アイデンティティベースポリシーをサポートする	あり
------------------------	----

アイデンティティベースポリシーは、IAM ユーザー、ユーザーグループ、ロールなど、アイデンティティにアタッチできる JSON 許可ポリシードキュメントです。これらのポリシーは、ユーザーとロールが実行できるアクション、リソース、および条件をコントロールします。アイデンティティベースのポリシーを作成する方法については、IAM ユーザーガイドの[IAM ポリシーの作成](#)を参照してください。

IAM アイデンティティベースのポリシーでは、許可または拒否するアクションとリソース、およびアクションを許可または拒否する条件を指定できます。プリンシパルは、それが添付されているユーザーまたはロールに適用されるため、アイデンティティベースのポリシーでは指定できません。JSON ポリシーで使用できるすべての要素について学ぶには、IAM ユーザーガイドの[IAM JSON ポリシーの要素のリファレンス](#)を参照してください。

AWS Glue では、すべての AWS Glue の操作についてアイデンティティベースのポリシー (IAM ポリシー) がサポートされています。ポリシーをアタッチすることで、AWS Glue Data Catalog 内のテーブルなどの AWS Glue リソースを作成し、そのリソースにアクセスして変更するためのアクセス許可を付与できます。

## AWS Glue のアイデンティティベースのポリシーの例

AWS Glue アイデンティティベースのポリシーの例を表示するには、「」を参照してください[AWS Glue のアイデンティティベースポリシーの例](#)。

## AWS Glue 内のリソースベースのポリシー

リソースベースのポリシーのサポート

部分的

リソースベースのポリシーは、リソースに添付する JSON 許可ポリシードキュメントです。リソースベースのポリシーには例として、IAM ロールの信頼ポリシー や Amazon S3 バケットポリシー があげられます。リソースベースのポリシーをサポートするサービスでは、サービス管理者はポリシーを使用して特定のリソースへのアクセスを制御できます。ポリシーがアタッチされているリソースの場合、指定されたプリンシパルがそのリソースに対して実行できるアクションと条件は、ポリシーによって定義されます。リソースベースのポリシーでは、[プリンシパルを指定する](#)必要があります。プリンシパルには、アカウント、ユーザー、ロール、フェデレーテッドユーザー、またはを含めることができます AWS のサービス。

クロスアカウントアクセスを有効にするには、アカウント全体、または別のアカウントの IAM エンティティをリソースベースのポリシーのプリンシパルとして指定します。リソースベースのポリシー

にクロスアカウントのプリンシパルを追加しても、信頼関係は半分しか確立されない点に注意してください。プリンシパルとリソースが異なる がある場合 AWS アカウント、信頼されたアカウントの IAM 管理者は、プリンシパルエンティティ (ユーザーまたはロール) にリソースへのアクセス許可も付与する必要があります。IAM 管理者は、アイデンティティベースのポリシーをエンティティにアタッチすることで権限を付与します。ただし、リソースベースのポリシーで、同じアカウントのプリンシパルへのアクセス権が付与されている場合は、アイデンティティベースのポリシーをさらに付与する必要はありません。詳細については、[「IAM ユーザーガイド」の「IAM でのクロスアカウントリソースアクセス」](#)を参照してください。

### Note

AWS Glue リソースポリシーは、Data Catalog リソースのアクセス許可を管理する場合にのみ使用できます。ジョブ、トリガー、開発エンドポイント、クローラ、または分類子など、他の AWS Glue リソースにアタッチすることはできません。リソースポリシーは 1 カタログにつき 1 つのみ許可されます。サイズの上限は 10 KB です。

AWS Glue では、リソースポリシーがカタログ にアタッチされます。カタログ は、前述のすべての種類の Data Catalog リソースの仮想コンテナです。各 AWS アカウントは、カタログ ID が AWS アカウント ID と同じ AWS リージョンで 1 つのカタログを所有しています。カタログを削除したり変更したりすることはできません。

リソースポリシーは、カタログへのすべての API コールで評価されます。この場合、呼び出し元のプリンシパルはポリシードキュメントの "Principal" ブロックに含まれているものとします。

AWS Glue リソースベースのポリシーの例を表示するには、「」を参照してください[AWS Glue のリソースベースのポリシーの例](#)。

## AWS Glue のポリシーアクション

ポリシーアクションに対するサポート	あり
-------------------	----

管理者は AWS JSON ポリシーを使用して、誰が何にアクセスできるかを指定できます。つまり、どのプリンシパルがどんなリソースにどんな条件でアクションを実行できるかということです。

JSON ポリシーの Action 要素には、ポリシー内のアクセスを許可または拒否するために使用できるアクションが記述されます。ポリシーアクションの名前は通常、関連付けられた AWS API オペレー

ションと同じです。一致する API オペレーションのない許可のみのアクションなど、いくつかの例外があります。また、ポリシーに複数のアクションが必要なオペレーションもあります。これらの追加アクションは、依存アクションと呼ばれます。

このアクションは、関連付けられたオペレーションを実行するための権限を付与するポリシーで使用されます。

AWS Glue アクションのリストを確認するには、「サービス認証リファレンス」の「[AWS Glue で定義されるアクション](#)」を参照してください。

AWS Glue のポリシーアクションは、アクションの前に次のプレフィックスを使用します。

```
glue
```

単一のステートメントで複数のアクションを指定するには、アクションをカンマで区切ります。

```
"Action": [
 "glue:action1",
 "glue:action2"
]
```

ワイルドカード (\*) を使用して複数アクションを指定できます。例えば、Get という単語で始まるすべてのアクションを指定するには、次のアクションを含めます。

```
"Action": "glue:Get*"
```

ポリシーの例を表示するには、「[AWS Glue のアクセスコントロールポリシーの例](#)」を参照してください。

## AWS Glue のポリシーリソース

ポリシーリソースに対するサポート	あり
------------------	----

管理者は AWS JSON ポリシーを使用して、誰が何にアクセスできるかを指定できます。つまり、どのプリンシパルがどのリソースにどのような条件でアクションを実行できるかということです。

Resource JSON ポリシー要素は、アクションが適用されるオブジェクトを指定します。ステートメントには、Resource または NotResource 要素を含める必要があります。ベストプラクティスとして、[Amazon リソースネーム \(ARN\)](#) を使用してリソースを指定します。これは、リソースレベルの許可と呼ばれる特定のリソースタイプをサポートするアクションに対して実行できます。

オペレーションのリスト化など、リソースレベルの権限をサポートしないアクションの場合は、ステートメントがすべてのリソースに適用されることを示すために、ワイルドカード (\*) を使用します。

```
"Resource": "*"
```

ARN を使用して AWS Glue リソースへのアクセスを制御する方法の詳細については、「」を参照してください。[AWS Glue リソース ARN の指定](#)。ARNs

AWS Glue リソースタイプとその ARNs」の「[AWS Glue で定義されるリソース](#)」を参照してください。各リソースの ARN を指定するために使用できるアクションについては、「[AWS Glue で定義されるアクション](#)」を参照してください。

## AWS Glue のポリシー条件キー

サービス固有のポリシー条件キーのサポート	あり
----------------------	----

管理者は AWS JSON ポリシーを使用して、誰が何にアクセスできるかを指定できます。つまり、どのプリンシパルがどんなリソースにどんな条件でアクションを実行できるかということです。

Condition 要素 (または Condition ブロック) を使用すると、ステートメントが有効な条件を指定できます。Condition 要素はオプションです。イコールや未満などの [条件演算子](#) を使用して条件式を作成することで、ポリシーの条件とリクエスト内の値を一致させることができます。

1 つのステートメントに複数の Condition 要素を指定する場合、または 1 つの Condition 要素に複数のキーを指定する場合、AWS では AND 論理演算子を使用してそれらを評価します。1 つの条件キーに複数の値を指定すると、は論理 OR オペレーションを使用して条件 AWS を評価します。ステートメントの権限が付与される前にすべての条件が満たされる必要があります。

条件を指定する際にプレースホルダー変数も使用できます。例えば IAM ユーザーに、IAM ユーザー名がタグ付けされている場合のみリソースにアクセスできる権限を付与することができます。詳細については、IAM ユーザーガイドの [IAM ポリシーの要素: 変数およびタグ](#) を参照してください。

AWS は、グローバル条件キーとサービス固有の条件キーをサポートします。すべての AWS グローバル条件キーを確認するには、「IAM ユーザーガイド」の [AWS 「グローバル条件コンテキストキー」](#) を参照してください。

AWS Glue の条件キーのリストを確認するには、「サービス認証リファレンス」の「[AWS Glue の条件キー](#)」を参照してください。条件キーを使用できるアクションとリソースについては、「[AWS Glue で定義されるアクション](#)」を参照してください。

ポリシーの例を表示するには、「[条件キーまたはコンテキストキーを使用して設定を制御する](#)」を参照してください。

## AWS Glue ACLs

ACL のサポート

なし

アクセスコントロールリスト (ACL) は、どのプリンシパル (アカウントメンバー、ユーザー、またはロール) がリソースにアクセスするための許可を持つかをコントロールします。ACL はリソーススペースのポリシーに似ていますが、JSON ポリシードキュメント形式は使用しません。

## AWS Glue での ABAC

ABAC (ポリシー内のタグ) のサポート

部分的

属性ベースのアクセス制御 (ABAC) は、属性に基づいてアクセス許可を定義する認可戦略です。では AWS、これらの属性はタグと呼ばれます。タグは、IAM エンティティ (ユーザーまたはロール) および多くの AWS リソースにアタッチできます。エンティティとリソースのタグ付けは、ABAC の最初の手順です。その後、プリンシパルのタグがアクセスしようとしているリソースのタグと一致した場合にオペレーションを許可するように ABAC ポリシーをします。

ABAC は、急成長する環境やポリシー管理が煩雑になる状況で役立ちます。

タグに基づいてアクセスを管理するには、`aws:ResourceTag/key-name`、`aws:RequestTag/key-name`、または `aws:TagKeys` の条件キーを使用して、ポリシーの [条件要素](#) でタグ情報を提供します。

サービスがすべてのリソースタイプに対して 3 つの条件キーすべてをサポートする場合、そのサービスの値はありです。サービスが一部のリソースタイプに対してのみ 3 つの条件キーのすべてをサポートする場合、値は「部分的」になります。

ABAC の詳細については、IAM ユーザーガイドの[ABAC とは?](#)を参照してください。ABAC をセットアップするステップを説明するチュートリアルについては、IAM ユーザーガイドの[属性に基づくアクセスコントロール \(ABAC\) を使用する](#)を参照してください。

#### Important

条件コンテキストキーは、クローラー、ジョブ、トリガー、開発エンドポイントに対する AWS Glue API アクションにのみ適用されます。影響を受ける API オペレーションの詳細については、「[AWS Glue の条件キー](#)」を参照してください。

AWS Glue データカタログ API オペレーションは、現在、aws:referrer および aws:UserAgent のグローバルな条件コンテキストキーをサポートしていません。

リソースのタグに基づいてリソースへのアクセスを制限するためのアイデンティティベースポリシーの例を表示するには、「[タグを使用してアクセスを許可する](#)」を参照してください。

## AWS Glue での一時的な認証情報の使用

一時的な認証情報のサポート

あり

一部の は、一時的な認証情報を使用してサインインすると機能 AWS のサービスしません。一時的な認証情報 AWS のサービス を使用する などの詳細については、IAM ユーザーガイドの[AWS のサービス 「IAM と連携する](#)」を参照してください。

ユーザー名とパスワード以外の AWS Management Console 方法で にサインインする場合、一時的な認証情報を使用します。例えば、会社の Single Sign-On (SSO) リンク AWS を使用して にアクセスすると、そのプロセスによって一時的な認証情報が自動的に作成されます。また、ユーザーとしてコンソールにサインインしてからロールを切り替える場合も、一時的な認証情報が自動的に作成されます。ロールの切り替えに関する詳細については、IAM ユーザーガイドの[ロールへの切り替え \(コンソール\)](#)を参照してください。

一時的な認証情報は、AWS CLI または AWS API を使用して手動で作成できます。その後、これらの一時的な認証情報を使用して .AWS recommends にアクセスできます AWS。これは、長期的なアクセスキーを使用する代わりに、一時的な認証情報を動的に生成することを推奨しています。詳細については、[IAM の一時的セキュリティ認証情報](#)を参照してください。

## AWS Glue のクロスサービスプリンシパル許可

転送アクセスセッション (FAS) をサポート	なし
-------------------------	----

IAM ユーザーまたはロールを使用してアクションを実行すると AWS、プリンシパルと見なされます。一部のサービスを使用する際に、アクションを実行することで、別のサービスの別のアクションがトリガーされることがあります。FAS は、 を呼び出すプリンシパルのアクセス許可を AWS のサービス、ダウンストリームサービス AWS のサービス へのリクエストのリクエストと組み合わせて使用します。FAS リクエストは、サービスが他の AWS のサービス またはリソースとのやり取りを完了する必要があるリクエストを受け取った場合にのみ行われます。この場合、両方のアクションを実行するためのアクセス許可が必要です。FAS リクエストを行う際のポリシーの詳細については、「[転送アクセスセッション](#)」を参照してください。

## AWS Glue のサービスロール

サービスロールに対するサポート	あり
-----------------	----

サービスロールとは、サービスがユーザーに代わってアクションを実行するために引き受ける [IAM ロール](#) です。IAM 管理者は、IAM 内からサービスロールを作成、変更、削除できます。詳細については、「IAM ユーザーガイド」の「[AWS のサービスにアクセス許可を委任するロールの作成](#)」を参照してください。

### Warning

サービスロールのアクセス許可を変更すると、AWS Glue の機能が破損する可能性があります。AWS Glue が指示する場合以外は、サービスロールを編集しないでください。

AWS Glue のサービスロールを作成する詳細な手順については、[ステップ 1: AWS Glue サービスの IAM ポリシーを作成する](#)「」および「」を参照してください。[ステップ 2: AWS Glue 用の IAM ロールを作成する](#)。

## AWS Glue のサービスにリンクされたロール

サービスにリンクされたロールのサポート	なし
---------------------	----

サービスにリンクされたロールは、にリンクされたサービスロールの一種です AWS のサービス。サービスは、ユーザーに代わってアクションを実行するロールを引き受けることができます。サービスにリンクされたロールは に表示され AWS アカウント、サービスによって所有されます。IAM 管理者は、サービスにリンクされたロールのアクセス許可を表示できますが、編集することはできません。

サービスリンクロールの作成または管理の詳細については、[IAM と提携するAWS のサービス](#)を参照してください。表の中から、[Service-linked role] (サービスにリンクされたロール) 列に Yes と記載されたサービスを見つけます。サービスリンクロールに関するドキュメントをサービスで表示するには、はい リンクを選択します。

## AWS Glue の IAM アクセス許可の設定

AWS Identity and Access Management (IAM) を使用して、リソースにアクセスするために AWS Glue が使用するポリシーとロールを定義します。次の手順では、AWS Glue のアクセス許可を設定するためのさまざまなオプションについて説明します。ビジネスニーズに応じて、リソースへのアクセスを追加または削減できます。

### Note

オプションを設定せず、基本的な IAM アクセス許可で AWS Glue の使用を開始する場合は、「[AWS Glue 用の IAM アクセス許可のセットアップ](#)」を参照してください。

1. [AWS Glue サービスの IAM ポリシーを作成する](#): AWS Glue リソースへのアクセスを許可するサービスポリシーを作成します。
2. [AWS Glue の IAM ロールを作成する](#): IAM ロールを作成して、AWS Glue が使用する Amazon Simple Storage Service (Amazon S3) リソース用の AWS Glue サービスポリシーおよびポリシーをアタッチします。
3. [AWS Glue にアクセスするユーザーまたはグループにポリシーをアタッチする](#): AWS Glue コンソールにサインインするユーザーまたはグループにポリシーをアタッチします。
4. [ノートブックの IAM ポリシーを作成する](#): 開発エンドポイント上のノートブックサーバーの作成に使用する、ノートブックサーバーポリシーを作成します。
5. [ノートブックの IAM ロールを作成する](#): IAM ロールを作成し、ノートブックサーバーポリシーをアタッチします。
6. [Amazon SageMaker ノートブックの IAM ポリシーを作成する](#): 開発エンドポイント上に Amazon SageMaker ノートブックを作成する際に使用する IAM ポリシーを作成します。

7. [Amazon SageMaker ノートブックの IAM ロールを作成する](#): IAM ロールを作成し、開発エンドポイント上に Amazon SageMaker ノートブックを作成する際にアクセスを許可するために、ポリシーをアタッチします。

## ステップ 1: AWS Glue サービスの IAM ポリシーを作成する

Amazon S3 のオブジェクトにアクセスするなど、別の AWS リソース上のデータにアクセスするオペレーションの場合、AWS Glue には、ユーザーの代わりにリソースにアクセスするためのアクセス許可が必要です。AWS Identity and Access Management (IAM) を使用してアクセス権限を提供できます。

### Note

AWS 管理ポリシー **AWSGlueServiceRole** を使用する場合は、このステップをスキップできます。

このステップでは、**AWSGlueServiceRole** に似たポリシーを作成します。**AWSGlueServiceRole** の最新バージョンは IAM コンソールで取得できます。

AWS Glue の IAM ポリシーを作成するには

このポリシーは、AWS Glue がこのポリシーを使用してロールを引き受ける際に必要となる、アカウント内のリソースを管理する Amazon S3 アクションの一部を許可します。このポリシーで指定されているリソースの一部は、Amazon S3 バケット、Amazon S3 ETL スクリプト、CloudWatch Logs、および Amazon EC2 リソースのために AWS Glue で使用される、デフォルト名を参照しています。簡素化のため AWS Glue のデフォルトでは、aws-glue-\* のプレフィックスが付いた Amazon S3 オブジェクトをアカウント内のバケットに書き込みます。

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. 左のナビゲーションペインの [ポリシー] を選択します。
3. [ポリシーの作成] を選択します。
4. [ポリシーの作成] 画面で、JSON 編集のためのタブに移動します。次の JSON ステートメントを使用してポリシードキュメントを作成して、[ポリシーの確認] を選択します。

**Note**

Amazon S3 リソースに必要なアクセス許可を追加します。アクセスポリシーのリソースセクションを必要なリソースだけに絞り込みたい場合があるかもしれません。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "glue:*",
 "s3:GetBucketLocation",
 "s3:ListBucket",
 "s3:ListAllMyBuckets",
 "s3:GetBucketAcl",
 "ec2:DescribeVpcEndpoints",
 "ec2:DescribeRouteTables",
 "ec2:CreateNetworkInterface",
 "ec2>DeleteNetworkInterface",
 "ec2:DescribeNetworkInterfaces",
 "ec2:DescribeSecurityGroups",
 "ec2:DescribeSubnets",
 "ec2:DescribeVpcAttribute",
 "iam:ListRolePolicies",
 "iam:GetRole",
 "iam:GetRolePolicy",
 "cloudwatch:PutMetricData"
],
 "Resource": [
 "*"
]
 },
 {
 "Effect": "Allow",
 "Action": [
 "s3:CreateBucket",
 "s3:PutBucketPublicAccessBlock"
],
 "Resource": [
```

```
 "arn:aws:s3:::aws-glue-*"
]
},
{
 "Effect": "Allow",
 "Action": [
 "s3:GetObject",
 "s3:PutObject",
 "s3:DeleteObject"
],
 "Resource": [
 "arn:aws:s3:::aws-glue-*/**",
 "arn:aws:s3:::*/*aws-glue-*/**"
]
},
{
 "Effect": "Allow",
 "Action": [
 "s3:GetObject"
],
 "Resource": [
 "arn:aws:s3:::crawler-public**",
 "arn:aws:s3:::aws-glue-*"
]
},
{
 "Effect": "Allow",
 "Action": [
 "logs:CreateLogGroup",
 "logs:CreateLogStream",
 "logs:PutLogEvents",
 "logs:AssociateKmsKey"
],
 "Resource": [
 "arn:aws:logs:*:*:log-group:/aws-glue/*"
]
},
{
 "Effect": "Allow",
 "Action": [
 "ec2:CreateTags",
 "ec2:DeleteTags"
],
 "Condition": {
```

```

 "ForAllValues:StringEquals": {
 "aws:TagKeys": [
 "aws-glue-service-resource"
]
 }
 },
 "Resource": [
 "arn:aws:ec2:*:*:network-interface/*",
 "arn:aws:ec2:*:*:security-group/*",
 "arn:aws:ec2:*:*:instance/*"
]
}
]
}

```

次の表は、このポリシーによって付与されたアクセス権限を示しています。

[アクション]	[リソース]	説明
"glue:*"	"*"	すべての AWS Glue API オペレーションを実行する権限を付与します。
"s3:GetBucketLocation", "s3:ListBucket", "s3:ListAllMyBuckets", "s3:GetBucketAcl",	"*"	クローラ、ジョブ、開発エンドポイント、ノートブックサーバーからの Amazon S3 バケットの一覧表示を許可します。

[アクション]	[リソース]	説明
<code>"ec2:DescribeVpcEndpoints", "ec2:DescribeRouteTables", "ec2:CreateNetworkInterface", "ec2&gt;DeleteNetworkInterface", "ec2:DescribeNetworkInterfaces", "ec2:DescribeSecurityGroups", "ec2:DescribeSubnets", "ec2:DescribeVpcAttribute",</code>	<code>"*"</code>	<p>Virtual Private Cloud (VPC) などの Amazon EC2 ネットワーク項目を、ジョブ、クローラ、開発エンドポイントの実行時に設定することを許可します。</p>
<code>"iam:ListRolePolicies", "iam:GetRole", "iam:GetRolePolicy"</code>	<code>"*"</code>	<p>クローラ、ジョブ、開発エンドポイント、ノートブックサーバーからの IAM ロールの一覧表示を許可します。</p>
<code>"cloudwatch:PutMetricData"</code>	<code>"*"</code>	<p>ジョブの CloudWatch メトリクスの書き込みを許可します。</p>
<code>"s3:CreateBucket", "s3:PutBucketPublicAccessBlock"</code>	<code>"arn:aws:s3:::aws-glue-*"</code>	<p>ジョブおよびノートブックサーバーからアカウントに Amazon S3 バケットを作成できます。</p> <p>命名規則: [aws-glue-] という Amazon S3 フォルダを使用します。</p> <p>AWS Glue による、パブリックアクセスをブロックするバケットの作成を有効化します。</p>

[アクション]	[リソース]	説明
"s3:GetObject", "s3:PutObject", "s3:DeleteObject"	"arn:aws:s3:::aws-glue-*/*", "arn:aws:s3:::*/*aws-glue-*/*"	<p>ETL スクリプトやノートブックサーバーのロケーションなどのオブジェクトを格納するために、アカウントに対する Amazon S3 オブジェクトの取得、配置、および削除を許可します。</p> <p>命名規則: 名前に aws-glue- のプレフィックスが付いている Amazon S3 バケットまたはフォルダにアクセス許可を付与します。</p>
"s3:GetObject"	"arn:aws:s3:::crawler-public*", "arn:aws:s3:::aws-glue-*"	<p>クローラやジョブに関するサンプル、もしくはチュートリアルで使用されている、Amazon S3 オブジェクトの取得を許可します。</p> <p>命名規則: Amazon S3 バケット名は「crawler-public」および「aws-glue-」で始まります。</p>
"logs:CreateLogGroup", "logs:CreateLogStream", "logs:PutLogEvents"	"arn:aws:logs:*:*:log-group:/aws-glue/*"	<p>ログを CloudWatch Logs に書き込むことを許可します。</p> <p>命名規則: AWS Glue は名前が [aws-glue] で始まるロググループにログを書き込みます。</p>

[アクション]	[リソース]	説明
"ec2:CreateTags", "ec2>DeleteTags"	"arn:aws:ec2:*:*:network-interface/*", "arn:aws:ec2:*:*:security-group/*", "arn:aws:ec2:*:*:instance/*"	開発エンドポイント用に作成された Amazon EC2 リソースのタグ付けを許可します。  命名規則: AWS Glue は [aws-glue-service-resource] を使用して、Amazon EC2 ネットワークインターフェイス、セキュリティグループ、およびインスタンスにタグ付けします。

5. [Review Policy (ポリシーの確認)] 画面で、[Policy Name (ポリシー名)] ( [GlueServiceRolePolicy] など ) を入力します。オプションの説明を入力し、ポリシーが適切であることを確認したら、[Create policy] を選択します。

## ステップ 2: AWS Glue 用の IAM ロールを作成する

ユーザーに代わり他のサービスを呼び出す際に AWS Glue が引き受けることができる、IAM ロールのアクセス許可を付与する必要があります。これには、AWS Glue で使用するすべてのソース、ターゲット、スクリプト、および一時ディレクトリでの Amazon S3 へのアクセスが含まれます。クローラ、ジョブ、および開発エンドポイントによって許可が必要です。

AWS Identity and Access Management (IAM) を使用してアクセス権限を提供できます。AWS Glue に渡す IAM ロールにポリシーを追加します。

用に IAM ロールを作成するには AWS Glue

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. 左のナビゲーションペインで、[ロール] を選択します。
3. [ロールの作成] を選択します。
4. ロールタイプについては、[AWS Service] (サービス) を選択してから、[Glue] を見つけて選択し、[Next: Permissions] (次へ: アクセス許可) をクリックします。

5. [アクセス権限ポリシーをアタッチする] のページで、必要なアクセス権限を含むポリシーを選択します。例えば、一般的な AWS 許可の `AWSGlueServiceRole` マネージドポリシー `AWS Glue` や、Amazon S3 リソースへアクセスするための AWS マネージドポリシー `AmazonS3FullAccess` などがあります。続いて、[Next: Review] をクリックします。

#### Note

Amazon S3 のソースとターゲットに対するアクセス許可を、このロールのポリシーの 1 つにより付与してください。独自のポリシーを、特定の Amazon S3 リソースにアクセスするために指定します。データソースには、`s3:ListBucket` および `s3:GetObject` アクセス許可が必要です。データターゲットには、`s3:ListBucket`、`s3:PutObject`、および `s3>DeleteObject` アクセス許可が必要です。リソースの Amazon S3 ポリシーの作成については、「[Specifying Resources in a Policy](#)」を参照してください。Amazon S3 ポリシーの例については、「[Writing IAM Policies: How to Grant Access to an Amazon S3 Bucket](#)」を参照してください。SSE-KMS で暗号化された Amazon S3 のソースとターゲットにアクセスする予定がある場合は、AWS Glue のクローラ、ジョブ、開発エンドポイントに、データを復号化するためのポリシーをアタッチしてください。詳細については、「[Protecting Data Using Server-Side Encryption with AWS KMS-Managed Keys \(SSE-KMS\)](#)」を参照してください。

次に例を示します。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "kms:Decrypt"
],
 "Resource": [
 "arn:aws:kms:*:account-id-without-hyphens:key/key-id"
]
 }
]
}
```

6. [Role name] (ロール名) に、`AWSGlueServiceRoleDefault` などのロール名を入力します。コンソールユーザーからサービスにロールを渡すには、名前に文字列 `AWSGlueServiceRole`

のプレフィックスが付けられたロールを作成します。AWSGlueServiceRole が提供するポリシーでは、IAM サービスロールが AWS Glue で始まることを想定しています。それ以外の場合は、ユーザーにポリシーを追加して、IAM ロールの iam:PassRole アクセス許可をユーザーの命名規則に一致させる必要があります。[ロールの作成] を選択します。

#### Note

ロールを使用してノートブックを作成すると、そのロールがインタラクティブセッションに渡され、その両方で同じロールを使用できるようになります。このように、iam:PassRole のアクセス許可がロールのポリシーの一部として必要です。次の例を使用して、ロール用の新しいポリシーを作成します。アカウント番号とロール名を自分のものに置き換えます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": "iam:PassRole",
 "Resource": "arn:aws:iam::090000000210:role/<role_name>"
 }
]
}
```

### ステップ 3: AWS Glue にアクセスするユーザーまたはグループにポリシーをアタッチする

管理者は、AWS Glue コンソールまたは AWS Command Line Interface (AWS CLI) を使用して、すべてのユーザー、グループ、ロールにアクセス許可を割り当てる必要があります。これらのアクセス許可は、AWS Identity and Access Management (IAM) を使用しながらポリシーを介して提供します。このステップでは、ユーザーまたはグループにアクセス許可を割り当てる方法について説明します。

このステップを完了すると、ユーザーまたはグループに次のポリシーがアタッチされます。

- AWS 管理ポリシー [AWSGlueConsoleFullAccess]、またはカスタムポリシー [GlueConsoleAccessPolicy]。
- **AWSGlueConsoleSageMakerNotebookFullAccess**
- **CloudWatchLogsReadOnlyAccess**
- **AWSCloudFormationReadOnlyAccess**
- **AmazonAthenaFullAccess**

インラインポリシーをアタッチしてこれをユーザーまたはグループに埋め込むには

AWS Glue コンソールにアクセスするユーザーまたはグループに、AWS マネージドポリシーまたはインラインポリシーをアタッチすることができます。このポリシーで指定されているリソースの中には、AWS Glue で (Amazon S3 バケット、Amazon S3 ETL スクリプト、CloudWatch Logs、AWS CloudFormation、および Amazon EC2 リソース用として) 使用されるデフォルトの名前があります。簡素化のため AWS Glue のデフォルトでは、aws-glue-\* のプレフィックスが付いた Amazon S3 オブジェクトをアカウント内のバケットに書き込みます。

#### Note

AWS 管理ポリシー **AWSGlueConsoleFullAccess** を使用する場合は、このステップをスキップできます。

#### Important

AWS Glue には、ユーザーの代理操作を実行するために使用されるロールを引き受けるアクセス権限が必要です。そのためには、**iam:PassRole** アクセス権限を AWS Glue ユーザーまたはグループに追加します。このポリシーは、AWS Glue サービスロールの **AWSGlueServiceRole** で始まるロール、およびノートブックサーバーの作成に必要なロール **AWSGlueServiceNotebookRole** にアクセス権限を与えます。また、命名規則に従った **iam:PassRole** アクセス権限の独自のポリシーを作成することもできます。

セキュリティのベストプラクティスに従って、ポリシーを強化して、Amazon S3 バケットおよび Amazon CloudWatch ロググループへのアクセスをさらに制限することが推奨されます。Amazon S3 ポリシーの例については、「[Writing IAM Policies: How to Grant Access to an Amazon S3 Bucket](#)」を参照してください。

このステップでは、AWSGlueConsoleFullAccess に似たポリシーを作成します。AWSGlueConsoleFullAccess の最新バージョンは IAM コンソールで取得できます。

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. ナビゲーションペインで、[Users] (ユーザー) または [User groups] (ユーザーグループ) を選択します。
3. 一覧から、ポリシーを埋め込むユーザーまたはグループの名前を選択します。
4. [Permissions (アクセス許可)] タブを選択して、必要であれば [Permissions policies (アクセス許可ポリシー)] セクションを展開します。
5. [Add Inline policy] (インラインポリシーの追加) リンクを選択します。
6. [ポリシーの作成] 画面で、JSON 編集のためのタブに移動します。次の JSON ステートメントを使用してポリシードキュメントを作成して、[ポリシーの確認] を選択します。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "glue:*",
 "redshift:DescribeClusters",
 "redshift:DescribeClusterSubnetGroups",
 "iam:ListRoles",
 "iam:ListUsers",
 "iam:ListGroups",
 "iam:ListRolePolicies",
 "iam:GetRole",
 "iam:GetRolePolicy",
 "iam:ListAttachedRolePolicies",
 "ec2:DescribeSecurityGroups",
 "ec2:DescribeSubnets",
 "ec2:DescribeVpcs",
 "ec2:DescribeVpcEndpoints",
 "ec2:DescribeRouteTables",
 "ec2:DescribeVpcAttribute",
 "ec2:DescribeKeyPairs",
 "ec2:DescribeInstances",
 "rds:DescribeDBInstances",
 "rds:DescribeDBClusters",
```

```

 "rds:DescribeDBSubnetGroups",
 "s3:ListAllMyBuckets",
 "s3:ListBucket",
 "s3:GetBucketAcl",
 "s3:GetBucketLocation",
 "cloudformation:DescribeStacks",
 "cloudformation:GetTemplateSummary",
 "dynamodb:ListTables",
 "kms:ListAliases",
 "kms:DescribeKey",
 "cloudwatch:GetMetricData",
 "cloudwatch:ListDashboards"
],
 "Resource": [
 "*"
]
},
{
 "Effect": "Allow",
 "Action": [
 "s3:GetObject",
 "s3:PutObject"
],
 "Resource": [
 "arn:aws:s3::*/*aws-glue-*/*",
 "arn:aws:s3:::aws-glue-*"
]
},
{
 "Effect": "Allow",
 "Action": [
 "tag:GetResources"
],
 "Resource": [
 "*"
]
},
{
 "Effect": "Allow",
 "Action": [
 "s3:CreateBucket",
 "s3:PutBucketPublicAccessBlock"
],
 "Resource": [
 "arn:aws:s3:::aws-glue-*"
]
},

```

```
]
},
{
 "Effect": "Allow",
 "Action": [
 "logs:GetLogEvents"
],
 "Resource": [
 "arn:aws:logs:*:*:/aws-glue/*"
]
},
{
 "Effect": "Allow",
 "Action": [
 "cloudformation:CreateStack",
 "cloudformation>DeleteStack"
],
 "Resource": "arn:aws:cloudformation:*:*:stack/aws-glue*/*"
},
{
 "Effect": "Allow",
 "Action": [
 "ec2:RunInstances"
],
 "Resource": [
 "arn:aws:ec2:*:*:instance/*",
 "arn:aws:ec2:*:*:key-pair/*",
 "arn:aws:ec2:*:*:image/*",
 "arn:aws:ec2:*:*:security-group/*",
 "arn:aws:ec2:*:*:network-interface/*",
 "arn:aws:ec2:*:*:subnet/*",
 "arn:aws:ec2:*:*:volume/*"
]
},
{
 "Action": [
 "iam:PassRole"
],
 "Effect": "Allow",
 "Resource": "arn:aws:iam:*:*:role/AWSGlueServiceRole*",
 "Condition": {
 "StringLike": {
 "iam:PassedToService": [
 "glue.amazonaws.com"
]
 }
 }
}
```

```
]
 }
},
{
 "Action": [
 "iam:PassRole"
],
 "Effect": "Allow",
 "Resource": "arn:aws:iam::*:role/AWSGlueServiceNotebookRole*",
 "Condition": {
 "StringLike": {
 "iam:PassedToService": [
 "ec2.amazonaws.com"
]
 }
 }
},
{
 "Action": [
 "iam:PassRole"
],
 "Effect": "Allow",
 "Resource": [
 "arn:aws:iam::*:role/service-role/AWSGlueServiceRole*"
],
 "Condition": {
 "StringLike": {
 "iam:PassedToService": [
 "glue.amazonaws.com"
]
 }
 }
}
]
```

次の表は、このポリシーによって付与されたアクセス権限を示しています。

[アクション]	[リソース]	説明
"glue:*"	"*"	<p>すべての AWS Glue API オペレーションを実行する権限を付与します。</p> <p>"glue:*" アクションを使用せずに以前にポリシーを作成したことがある場合は、ポリシーに次の個別のアクセス許可を追加する必要があります。</p> <ul style="list-style-type: none"><li>"glue:ListCrawlers"</li><li>"glue:BatchGetCrawlers"</li><li>"glue:ListTriggers"</li><li>"glue:BatchGetTriggers"</li><li>"glue:ListDevEndpoints"</li><li>"glue:BatchGetDevEndpoints"</li><li>"glue:ListJobs"</li><li>"glue:BatchGetJobs"</li></ul>
"redshift:DescribeClusters", "redshift:DescribeClusterSubnetGroups"	"*"	Amazon Redshift に対する接続の作成を許可します。

[アクション]	[リソース]	説明
"iam:ListRoles", "iam:ListRolePolicies", "iam:GetRole", "iam:GetRolePolicy", "iam:ListAttachedRolePolicies"	"*"	クローラ、ジョブ、開発エンドポイント、ノートブックサーバーを使用する際の IAM ロールの一覧表示を許可します。
"ec2:DescribeSecurityGroups", "ec2:DescribeSubnets", "ec2:DescribeVpcs", "ec2:DescribeVpcEndpoints", "ec2:DescribeRouteTables", "ec2:DescribeVpcAttribute", "ec2:DescribeKeyPairs", "ec2:DescribeInstances"	"*"	ジョブ、クローラ、開発エンドポイントを実行する場合に、VPC などの Amazon EC2 ネットワーク項目を設定することを許可します。
"rds:DescribeDBInstances"	"*"	Amazon RDS に対する接続の作成を許可します。
"s3:ListAllMyBuckets", "s3:ListBucket", "s3:GetBucketAcl", "s3:GetBucketLocation"	"*"	クローラ、ジョブ、開発エンドポイント、ノートブックサーバーを使用する際に、Amazon S3 バケットを一覧表示することを許可します。
"dynamodb:ListTables"	"*"	DynamoDB テーブルの一覧表示を許可します。
"kms:ListAliases", "kms:DescribeKey"	"*"	KMS キーを使用できます。

[アクション]	[リソース]	説明
"cloudwatch:GetMetricData", "cloudwatch:ListDashboards"	"*"	CloudWatch メトリクスを使用できます。
"s3:GetObject", "s3:PutObject"	"arn:aws:s3::: aws-glue-*/*", "arn:aws:s3::: */*aws-glue-*/*", "arn:aws:s3::: aws-glue-*"	<p>ETL スクリプトやノートブックサーバーのロケーションなどのオブジェクトを格納する際に、アカウントで Amazon S3 オブジェクトを取得および配置することを許可します。</p> <p>命名規則: 名前に aws-glue- のプレフィックスが付いている Amazon S3 バケットまたはフォルダにアクセス許可を付与します。</p>
"tag:GetResources"	"*"	AWS タグの取得を許可します。

[アクション]	[リソース]	説明
"s3:CreateBucket", "s3:PutBucketPublicAccessBlock"	"arn:aws:s3::: aws-glue-*"	<p>ETL スクリプトやノートブックサーバーのロケーションなどのオブジェクトを格納する際に、アカウントで Amazon S3 バケットを作成することを許可します。</p> <p>命名規則: 名前に aws-glue- のプレフィックスが付いている Amazon S3 バケットにアクセス許可を付与します。</p> <p>AWS Glue による、パブリックアクセスをブロックするバケットの作成を有効化します。</p>
"logs:GetLogEvents"	"arn:aws:logs:*:*: /aws-glue/*"	<p>CloudWatch Logs の取得を許可します。</p> <p>命名規則: AWS Glue は名前が [aws-glue-] で始まるロググループにログを書き込みます。</p>

[アクション]	[リソース]	説明
"cloudformation:CreateStack", "cloudformation>DeleteStack"	"arn:aws:cloudformation:*:*:stack/aws-glue*/**"	<p>ノートブックサーバーで作業するときに AWS CloudFormation スタックの管理を許可します。</p> <p>命名規則: AWS Glue は名前が [aws-glue] で始まるスタックを作成します。</p>
"ec2:RunInstances"	"arn:aws:ec2:*:*:instance/**", "arn:aws:ec2:*:*:key-pair/**", "arn:aws:ec2:*:*:image/**", "arn:aws:ec2:*:*:security-group/**", "arn:aws:ec2:*:*:network-interface/**", "arn:aws:ec2:*:*:subnet/**", "arn:aws:ec2:*:*:volume/**"	開発エンドポイントとノートブックサーバーの実行を許可します。
"iam:PassRole"	"arn:aws:iam:*:*:role/AWSGlueServiceRole*"	AWSGlueServiceRole で始まるロールに対する PassRole の許可を、AWS Glue が引き受けられるようにします。

[アクション]	[リソース]	説明
"iam:PassRole"	"arn:aws:iam::*:role/ AWSGlueServiceNotebookRole*"	AWSGlueServiceNotebookRole で始まるロールに関するアクセス許可 PassRole を引き受けることを、Amazon EC2 に 許可します。
"iam:PassRole"	"arn:aws:iam::*:role/service-role/ AWSGlueServiceRole*"	service-role/ AWSGlueServiceRole で始まるロールに対する PassRole の許可を、AWS Glue が引き受けられるようにします。

7. [ポリシーの確認] セクションで、ポリシー名を入力します ( 例: GlueConsoleAccessPolicy )。ポリシーが完成したら、[Create policy (ポリシーの作成)] を選択します。画面上部の赤いボックスにエラーが表示されていないことを確認します。報告されたエラーがあれば、修正します。

 Note

[Use autoformatting (自動フォーマットを使用する)] を選択した場合は、ポリシーを開いたときおよび [Validate Policy (ポリシーの検証)] を選択したときに毎回、ポリシーが再フォーマットされます。

AWSGlueConsoleFullAccess 管理ポリシーをアタッチするには

AWSGlueConsoleFullAccess ポリシーをアタッチすると、AWS Glue コンソールユーザーが必要とするアクセス許可を提供できます。

**Note**

AWS Glue コンソールアクセス用に独自のポリシーを作成した場合は、このステップをスキップできます。

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. ナビゲーションペインで、[ポリシー] を選択します。
3. ポリシーのリストで、[AWSGlueConsoleFullAccess] ポリシーの横にあるチェックボックスを選択します。[Filter (フィルター)] メニューと検索ボックスを使用して、ポリシーのリストをフィルタリングできます。
4. [Policy actions] を選択して、[Attach] を選択します。
5. ポリシーをアタッチするユーザーを選択します。[Filter] メニューと検索ボックスを使用して、プリンシパルエンティティのリストをフィルタリングできます。ポリシーをアタッチするユーザーを選択し、[ポリシーのアタッチ] を選択します。

**AWSGlueConsoleSageMakerNotebookFullAccess** 管理ポリシーをアタッチするには

AWSGlueConsoleSageMakerNotebookFullAccess ポリシーをユーザーにアタッチして、AWS Glue コンソールで作成した SageMaker ノートブックを管理できます。その他の必要な AWS Glue コンソールアクセス許可に加えて、このポリシーは SageMaker ノートブックの管理に必要なリソースへのアクセス許可を付与します。

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. ナビゲーションペインで、[ポリシー] を選択します。
3. ポリシーのリストで、[AWSGlueConsoleSageMakerNotebookFullAccess] の横にあるチェックボックスを選択します。[Filter (フィルター)] メニューと検索ボックスを使用して、ポリシーのリストをフィルタリングできます。
4. [Policy actions] を選択して、[Attach] を選択します。
5. ポリシーをアタッチするユーザーを選択します。[Filter] メニューと検索ボックスを使用して、プリンシパルエンティティのリストをフィルタリングできます。ポリシーをアタッチするユーザーを選択し、[ポリシーのアタッチ] を選択します。

## CloudWatchLogsReadOnlyAccess 管理ポリシーをアタッチするには

[CloudWatchLogsReadOnlyAccess] ポリシーをユーザーにアタッチして、CloudWatch Logs コンソールで AWS Glue によって作成されたログを表示できます。

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. ナビゲーションペインで、[ポリシー] を選択します。
3. ポリシーのリストで、[CloudWatchLogsReadOnlyAccess] ポリシーの横にあるチェックボックスを選択します。[Filter (フィルター)] メニューと検索ボックスを使用して、ポリシーのリストをフィルタリングできます。
4. [Policy actions] を選択して、[Attach] を選択します。
5. ポリシーをアタッチするユーザーを選択します。[Filter] メニューと検索ボックスを使用して、プリンシパルエンティティのリストをフィルタリングできます。ポリシーをアタッチするユーザーを選択し、[ポリシーのアタッチ] を選択します。

## AWSCloudFormationReadOnlyAccess 管理ポリシーをアタッチするには

[AWSCloudFormationReadOnlyAccess] ポリシーをユーザーにアタッチして、AWS CloudFormation コンソールで AWS Glue が使用する AWS CloudFormation スタックを表示できます。

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. ナビゲーションペインで、[ポリシー] を選択します。
3. ポリシーのリストで、[AWSCloudFormationReadOnlyAccess] ポリシーの横にあるチェックボックスをオンにします。[Filter (フィルター)] メニューと検索ボックスを使用して、ポリシーのリストをフィルタリングできます。
4. [Policy actions] を選択して、[Attach] を選択します。
5. ポリシーをアタッチするユーザーを選択します。[Filter] メニューと検索ボックスを使用して、プリンシパルエンティティのリストをフィルタリングできます。ポリシーをアタッチするユーザーを選択し、[ポリシーのアタッチ] を選択します。

## AmazonAthenaFullAccess 管理ポリシーをアタッチするには

[AmazonAthenaFullAccess] ポリシーをユーザーにアタッチすると、Amazon S3 データを Athena コンソールで表示できます。

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. ナビゲーションペインで、[ポリシー] を選択します。
3. ポリシーのリストで、[AmazonAthenaFullAccess] ポリシーの横にあるチェックボックスを選択します。[Filter (フィルター)] メニューと検索ボックスを使用して、ポリシーのリストをフィルタリングできます。
4. [Policy actions] を選択して、[Attach] を選択します。
5. ポリシーをアタッチするユーザーを選択します。[Filter] メニューと検索ボックスを使用して、プリンシパルエンティティのリストをフィルタリングできます。ポリシーをアタッチするユーザーを選択し、[ポリシーのアタッチ] を選択します。

#### ステップ 4: ノートブックサービス用に IAM ポリシーを作成する

開発エンドポイントでノートブックを使用する予定の場合は、ノートブックサーバーの作成時にアクセス許可を指定する必要があります。AWS Identity and Access Management (IAM) を使用してアクセス権限を提供できます。

このポリシーは、AWS Glue がこのポリシーを使用してロールを引き受ける際に必要となる、アカウント内のリソースを管理する Amazon S3 アクションの一部を許可します。このポリシーで指定されているリソースの一部は、AWS Glue で (Amazon S3 バケット、Amazon S3 ETL スクリプト、および Amazon EC2 リソース用として) 使用されるデフォルトの名前を参照しています。分かりやすさのために、AWS Glue のデフォルトでは一部の Amazon S3 オブジェクトを、アカウント内の `aws-glue-*` のプレフィックスが付いたバケットに書き込んでいます。

#### Note

AWS 管理ポリシー **AWSGlueServiceNotebookRole** を使用する場合は、このステップをスキップできます。

このステップでは、**AWSGlueServiceNotebookRole** に似たポリシーを作成します。**AWSGlueServiceNotebookRole** の最新バージョンは IAM コンソールで取得できます。

ノートブック用に IAM ポリシーを作成するには

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。

2. 左のナビゲーションペインの [ポリシー] を選択します。
3. [ポリシーの作成] を選択します。
4. [ポリシーの作成] 画面で、JSON 編集のためのタブに移動します。次の JSON ステートメントを使用してポリシードキュメントを作成して、[ポリシーの確認] を選択します。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "glue:CreateDatabase",
 "glue:CreatePartition",
 "glue:CreateTable",
 "glue>DeleteDatabase",
 "glue>DeletePartition",
 "glue>DeleteTable",
 "glue:GetDatabase",
 "glue:GetDatabases",
 "glue:GetPartition",
 "glue:GetPartitions",
 "glue:GetTable",
 "glue:GetTableVersions",
 "glue:GetTables",
 "glue:UpdateDatabase",
 "glue:UpdatePartition",
 "glue:UpdateTable",
 "glue:GetJobBookmark",
 "glue:ResetJobBookmark",
 "glue:CreateConnection",
 "glue:CreateJob",
 "glue>DeleteConnection",
 "glue>DeleteJob",
 "glue:GetConnection",
 "glue:GetConnections",
 "glue:GetDevEndpoint",
 "glue:GetDevEndpoints",
 "glue:GetJob",
 "glue:GetJobs",
 "glue:UpdateJob",
 "glue:BatchDeleteConnection",
 "glue:UpdateConnection",

```

```
 "glue:GetUserDefinedFunction",
 "glue:UpdateUserDefinedFunction",
 "glue:GetUserDefinedFunctions",
 "glue>DeleteUserDefinedFunction",
 "glue:CreateUserDefinedFunction",
 "glue:BatchGetPartition",
 "glue:BatchDeletePartition",
 "glue:BatchCreatePartition",
 "glue:BatchDeleteTable",
 "glue:UpdateDevEndpoint",
 "s3:GetBucketLocation",
 "s3:ListBucket",
 "s3:ListAllMyBuckets",
 "s3:GetBucketAcl"
],
 "Resource": [
 "*"
]
},
{
 "Effect": "Allow",
 "Action": [
 "s3:GetObject"
],
 "Resource": [
 "arn:aws:s3:::crawler-public*",
 "arn:aws:s3:::aws-glue*"
]
},
{
 "Effect": "Allow",
 "Action": [
 "s3:PutObject",
 "s3>DeleteObject"
],
 "Resource": [
 "arn:aws:s3:::aws-glue*"
]
},
{
 "Effect": "Allow",
 "Action": [
 "ec2:CreateTags",
 "ec2>DeleteTags"
]
}
```

```

],
 "Condition":{
 "ForAllValues:StringEquals":{
 "aws:TagKeys":[
 "aws-glue-service-resource"
]
 }
 },
 "Resource":[
 "arn:aws:ec2:*:*:network-interface/*",
 "arn:aws:ec2:*:*:security-group/*",
 "arn:aws:ec2:*:*:instance/*"
]
 }
]
}

```

次の表は、このポリシーによって付与されたアクセス権限を示しています。

[アクション]	[リソース]	説明
"glue:*"	"*"	すべての AWS Glue API オペレーションを実行する権限を付与します。
"s3:GetBucketLocation", "s3:ListBucket", "s3:ListAllMyBuckets", "s3:GetBucketAcl"	"*"	ノートブックサーバーからの Amazon S3 バケットの一覧表示を許可します。

[アクション]	[リソース]	説明
"s3:GetObject"	"arn:aws:s3:::crawler-public*", "arn:aws:s3:::aws-glue-*"	<p>ノートブックのサンプルやチュートリアルで使用されている Amazon S3 オブジェクトを取得できます。</p> <p>命名規則: Amazon S3 バケット名は「crawler-public」および「aws-glue-」で始まります。</p>
"s3:PutObject", "s3:DeleteObject"	"arn:aws:s3:::aws-glue*"	<p>ノートブックからアカウントに対する、Amazon S3 オブジェクトの配置と削除を許可します。</p> <p>命名規則: [aws-glue] という Amazon S3 フォルダを使用します。</p>
"ec2:CreateTags", "ec2:DeleteTags"	"arn:aws:ec2:*:*:network-interface/*", "arn:aws:ec2:*:*:security-group/*", "arn:aws:ec2:*:*:instance/*"	<p>ノートブックサーバー用に作成された Amazon EC2 リソースのタグ付けを許可します。</p> <p>命名規則: AWS Glue は Amazon EC2 インスタンスを、[aws-glue-service-resource] によりタグ付けします。</p>

5. [Review Policy (ポリシーの確認)] 画面で、[Policy Name (ポリシー名)] ( [GlueServiceNotebookPolicyDefault] など ) を入力します。オプションの説明を入力し、ポリシーが適切であることを確認したら、[Create policy] を選択します。

## ステップ 5: ノートブックサービス用に IAM ロールを作成する

開発エンドポイントでノートブックを使用する予定がある場合は、IAM ロールのアクセス許可を付与する必要があります。AWS Identity and Access Management IAM を使用しながら IAM ロールを介して、これらのアクセス許可を付与できます。

### Note

IAM コンソールを使用して IAM ロールを作成すると、コンソールによりインスタンスプロファイルが自動的に作成され、対応するロールと同じ名前が付けられます。

ノートブック用に IAM ロールを作成するには

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. 左のナビゲーションペインで、[ロール] を選択します。
3. [ロールの作成] を選択します。
4. ロールタイプについては、[AWS Service] (サービス) を選択した後に [EC2] を見つけて選択し、ユースケースで [EC2] を選択した後、[Next: Permissions] (次へ: アクセス許可) をクリックします。
5. [アクセス権限ポリシーを添付する] のページで、必要なアクセス権を含むポリシーを選択します。例えば、一般的な AWS Glue 許可には AWSGlueServiceNotebookRole、Amazon S3 リソースへのアクセスには AWS マネージドポリシー AmazonS3FullAccess を選択します。続いて、[Next: Review] をクリックします。

### Note

Amazon S3 のソースとターゲットに対するアクセス許可を、このロールのポリシーの 1 つにより付与してください。また、ノートブックサーバーの作成時にノートブックを保管する場所へのフルアクセスがポリシーで許可されていることを確認してください。独自のポリシーを、特定の Amazon S3 リソースにアクセスするために指定します。リソースの Amazon S3 ポリシーの作成については、「[Specifying Resources in a Policy](#)」を参照してください。

SSE-KMS で暗号化された Amazon S3 のソースとターゲットにアクセスする予定がある場合は、ノートブックがデータを復号化できるようにポリシーをアタッチします。詳細

については、「[Protecting Data Using Server-Side Encryption with AWS KMS-Managed Keys \(SSE-KMS\)](#)」を参照してください。

次に例を示します。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "kms:Decrypt"
],
 "Resource": [
 "arn:aws:kms:*:account-id-without-hyphens:key/key-id"
]
 }
]
}
```

6. ロール名 に、ロールの名前を入力します。コンソールユーザーからノートブックサーバーにロールを渡すには、名前に文字列 `AWSGlueServiceNotebookRole` のプレフィックスが付けられたロールを作成します。AWS Glue が提供するポリシーでは、IAM サービスロールが `AWSGlueServiceNotebookRole` で始まることを想定しています。それ以外の場合は、ユーザーにポリシーを追加して、IAM ロールに対する `iam:PassRole` の許可を命名規則に一致させる必要があります。たとえば、`AWSGlueServiceNotebookRoleDefault` と入力します。次に、`[Create role ( ロールの作成 )]` を選択します。

## ステップ 6: SageMaker ノートブック用に IAM ポリシーを作成する

開発エンドポイントで SageMaker ノートブックを使用する予定の場合は、ノートブックの作成時にアクセス許可を指定する必要があります。AWS Identity and Access Management (IAM) を使用してアクセス権限を提供できます。

SageMaker ノートブック用に IAM ポリシーを作成するには

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. 左のナビゲーションペインの `[ポリシー]` を選択します。
3. `[ポリシーの作成]` を選択します。

4. [Create Policy] ページで、JSON を編集するタブに移動します。次の JSON ステートメントを使用して、ポリシードキュメントを作成します。環境の *bucket-name*、*region-code*、*account-id* を編集します。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Action": [
 "s3:ListBucket"
],
 "Effect": "Allow",
 "Resource": [
 "arn:aws:s3:::bucket-name"
]
 },
 {
 "Action": [
 "s3:GetObject"
],
 "Effect": "Allow",
 "Resource": [
 "arn:aws:s3:::bucket-name*"
]
 },
 {
 "Action": [
 "logs:CreateLogStream",
 "logs:DescribeLogStreams",
 "logs:PutLogEvents",
 "logs:CreateLogGroup"
],
 "Effect": "Allow",
 "Resource": [
 "arn:aws:logs:region-code:account-id:log-group:/aws/sagemaker/*",
 "arn:aws:logs:region-code:account-id:log-group:/aws/sagemaker/
:log-stream:aws-glue-"
]
 },
 {
 "Action": [
 "glue:UpdateDevEndpoint",
 "glue:GetDevEndpoint",

```

```

 "glue:GetDevEndpoints"
],
 "Effect": "Allow",
 "Resource": [
 "arn:aws:glue:region-code:account-id:devEndpoint/*"
]
},
{
 "Action": [
 "sagemaker:ListTags"
],
 "Effect": "Allow",
 "Resource": [
 "arn:aws:sagemaker:region-code:account-id:notebook-instance/*"
]
}
]
}

```

次に [ポリシーの確認] を選択します。

次の表は、このポリシーによって付与されたアクセス権限を示しています。

[アクション]	[リソース]	説明
"s3:ListBucket*"	"arn:aws:s3::: <i>bucket-name</i> "	Amazon S3 バケットを一覧表示するアクセス許可を付与します
"s3:GetObject"	"arn:aws:s3::: <i>bucket-name</i> *"	SageMaker ノートブックによって使用される Amazon S3 オブジェクトを取得するための、アクセス許可を付与します。

[アクション]	[リソース]	説明
"logs:CreateLogStream", "logs:DescribeLogStreams", "logs:PutLogEvents", "logs:CreateLogGroup"	"arn:aws:logs: <i>region-code</i> : <i>account-id</i> :log-group:/aws/sagemaker/*", "arn:aws:logs: <i>region-code</i> : <i>account-id</i> :log-group:/aws/sagemaker/*:log-stream:aws-glue-*	ノートブックから Amazon CloudWatch Logs にログを書き込むためのアクセス許可を付与します。  命名規則: 名前が [aws-glue] で始まるロググループに書き込みます。
"glue:UpdateDevEndpoint", "glue:GetDevEndpoint", "glue:GetDevEndpoints"	"arn:aws:glue: <i>region-code</i> : <i>account-id</i> :devEndpoint/*"	SageMaker ノートブックから開発エンドポイントを使用するためのアクセス許可を付与します。
"sagemaker:ListTags"	"arn:aws:sagemaker : <i>region-code</i> : <i>account-id</i> :notebook-instance/*"	SageMaker リソースのタグを返すためのアクセス許可を付与します。この aws-glue-dev-endpoint タグは、ノートブックを開発エンドポイントに接続するために SageMaker ノートブック上で必要です。

5. [Review Policy] (ポリシーの確認) 画面で、[Policy Name] (ポリシー名) に AWSGlueSageMakerNotebook (など) を入力します。オプションの説明を入力し、ポリシーが適切であることを確認したら、[Create policy] を選択します。

## ステップ 7: SageMaker ノートブック用に IAM ロールを作成する

開発エンドポイントで SageMaker ノートブックを使用する予定がある場合は、IAM ロールのアクセス許可を与える必要があります。AWS Identity and Access Management (IAM) を使用しながら IAM ロールを介して、アクセス許可を付与できます。

SageMaker ノートブック用に IAM ロールを作成するには

1. AWS Management Console にサインインして、IAM コンソール (<https://console.aws.amazon.com/iam/>) を開きます。
2. 左のナビゲーションペインで、[ロール] を選択します。
3. [ロールの作成] を選択します。
4. ロールタイプについては、[AWS Service] を選択し、[SageMaker] を見つけて選択してから、[SageMaker - Execution] ユースケースを選択します。その後、[Next] (次へ) を選択します。
5. [Attach permissions policy] ページで、必要なアクセス許可を含むポリシー ([AmazonSageMakerFullAccess] など) を選択します。[次へ: レビュー] を選択します。

SSE-KMS で暗号化された Amazon S3 のソースとターゲットにアクセスする予定がある場合は、次の例に示すように、ノートブックによるデータの復号化を許可するポリシーをアタッチしてください。詳細については、「[Protecting Data Using Server-Side Encryption with AWS KMS-Managed Keys \(SSE-KMS\)](#)」を参照してください。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "kms:Decrypt"
],
 "Resource": [
 "arn:aws:kms:*:account-id-without-hyphens:key/key-id"
]
 }
]
}
```

6. ロール名に、ロールの名前を入力します。コンソールユーザーから SageMaker にロールを渡すには、文字列 `AWSGlueServiceSageMakerNotebookRole` のプレフィックスが付けられた名前を使用します。AWS Glue が提供するポリシーでは、IAM ロールが `AWSGlueServiceSageMakerNotebookRole` で始まることを想定しています。それ以外の場合は、ユーザーにポリシーを追加して、IAM ロールに対する `iam:PassRole` の許可を命名規則に一致させる必要があります。

例えば、AWSGlueServiceSageMakerNotebookRole-Default を入力した上で、[Create role] (ロールを作成) をクリックします。

7. ロールを作成した後に、AWS Glue から SageMaker ノートブックを作成するために必要な、追加のアクセス許可を付与するポリシーをアタッチします。

作成したロール AWSGlueServiceSageMakerNotebookRole-Default を開き、[Attach policies] (ポリシーのアタッチ) をクリックします。作成した AWSGlueSageMakerNotebook ポリシーを、ロールにアタッチします。

## AWS Glue のアクセスコントロールポリシーの例

このセクションでは、アイデンティティベースの (IAM) アクセスコントロールポリシーと、AWS Glue リソースポリシーの両方の例を示します。

### 目次

- [AWS Glue のアイデンティティベースポリシーの例](#)
  - [ポリシーのベストプラクティス](#)
  - [特定の AWS Glue オブジェクトにのみ適用されるリソースレベルのアクセス許可](#)
  - [AWS Glue コンソールの使用](#)
  - [自分の権限の表示をユーザーに許可する](#)
  - [テーブルへの読み取り専用のアクセスを許可する](#)
  - [GetTables アクセス許可によるテーブルのフィルタリング](#)
  - [テーブルとすべてのパーティションへのフルアクセスを付与する](#)
  - [名前のプレフィックスと明示的な拒否によりアクセスを制御する](#)
  - [タグを使用してアクセスを許可する](#)
  - [タグを使用してアクセスを拒否する](#)
  - [リストとバッチ API オペレーションでタグを使用する](#)
  - [条件キーまたはコンテキストキーを使用して設定を制御する](#)
    - [条件キーを使って設定を制御するポリシーを制御する](#)
    - [条件キーを使って設定を制御するポリシーを制御する](#)
  - [ID によるデータプレビューセッションの作成を拒否する](#)
- [AWS Glue のリソースベースのポリシーの例](#)

- [AWS Glue でリソースベースのポリシーを使用する際の考慮事項](#)
- [同じアカウント内のアクセスをコントロールするリソースポリシーの使用](#)

## AWS Glue のアイデンティティベースポリシーの例

デフォルトでは、ユーザーおよびロールには、AWS Glue リソースを作成または変更するアクセス許可は付与されていません。また、AWS Management Console、AWS Command Line Interface (AWS CLI)、または AWS API を使用してタスクを実行することもできません。IAM 管理者は、リソースで必要なアクションを実行するための権限をユーザーに付与する IAM ポリシーを作成できます。その後、管理者はロールに IAM ポリシーを追加し、ユーザーはロールを引き受けることができます。

これらサンプルの JSON ポリシードキュメントを使用して、IAM アイデンティティベースのポリシーを作成する方法については、『IAM ユーザーガイド』の「[IAM ポリシーの作成](#)」を参照してください。

AWS Glue で定義されるアクションとリソースタイプ (リソースタイプごとの ARN の形式を含む) の詳細については、「サービス認証リファレンス」の「[AWS Glue のアクション、リソース、および条件キー](#)」を参照してください。

### Note

このセクションで取り上げる例は、すべて us-west-2 リージョンを使用しています。このリージョンは、ご自身が使用する AWS リージョンに置き換えることが可能です。

### トピック

- [ポリシーのベストプラクティス](#)
- [特定の AWS Glue オブジェクトにのみ適用されるリソースレベルのアクセス許可](#)
- [AWS Glue コンソールの使用](#)
- [自分の権限の表示をユーザーに許可する](#)
- [テーブルへの読み取り専用のアクセスを許可する](#)
- [GetTables アクセス許可によるテーブルのフィルタリング](#)
- [テーブルとすべてのパーティションへのフルアクセスを付与する](#)
- [名前のプレフィックスと明示的な拒否によりアクセスを制御する](#)
- [タグを使用してアクセスを許可する](#)

- [タグを使用してアクセスを拒否する](#)
- [リストとバッチ API オペレーションでタグを使用する](#)
- [条件キーまたはコンテキストキーを使用して設定を制御する](#)
- [ID によるデータプレビューセッションの作成を拒否する](#)

## ポリシーのベストプラクティス

アイデンティティベースのポリシーは、ユーザーのアカウントで誰が AWS Glue リソースを作成し、これにアクセスし、これを削除できるかを決定します。これらのアクションを実行すると、AWS アカウント に料金が発生する可能性があります。アイデンティティベースポリシーを作成したり編集したりする際には、以下のガイドラインと推奨事項に従ってください:

- AWS マネージドポリシーを使用して開始し、最小特権の権限に移行する – ユーザーとワークロードへの権限の付与を開始するには、多くの一般的なユースケースのために権限を付与する AWS マネージドポリシーを使用します。これらは AWS アカウントで使用できます。ユースケースに応じた AWS カスタマーマネージドポリシーを定義することで、権限をさらに減らすことをお勧めします。詳細については、『IAM ユーザーガイド』の「[AWS マネージドポリシー](#)」または「[AWS ジョブ機能の管理ポリシー](#)」を参照してください。
- 最小特権を適用する – IAM ポリシーで権限を設定するときは、タスクの実行に必要な権限のみを付与します。これを行うには、特定の条件下で特定のリソースに対して実行できるアクションを定義します。これは、最小特権権限とも呼ばれています。IAM を使用して権限を適用する方法の詳細については、『IAM ユーザーガイド』の「[IAM でのポリシーと権限](#)」を参照してください。
- IAM ポリシーで条件を使用してアクセスをさらに制限する - ポリシーに条件を追加して、アクションやリソースへのアクセスを制限できます。例えば、ポリシー条件を記述して、すべてのリクエストを SSL を使用して送信するように指定できます。また、AWS CloudFormation などの特定の AWS のサービスを介して使用する場合、条件を使ってサービスアクションへのアクセス権を付与することもできます。詳細については、『IAM ユーザーガイド』の「[IAM JSON policy elements: Condition](#)」(IAM JSON ポリシー要素 : 条件) を参照してください。
- IAM Access Analyzer を使用して IAM ポリシーを検証し、安全で機能的な権限を確保する - IAM Access Analyzer は、新規および既存のポリシーを検証して、ポリシーが IAM ポリシー言語 (JSON) および IAM のベストプラクティスに準拠するようにします。IAM アクセスアナライザーは 100 を超えるポリシーチェックと実用的な推奨事項を提供し、安全で機能的なポリシーの作成をサポートします。詳細については、『IAM ユーザーガイド』の「[IAM Access Analyzer ポリシーの検証](#)」を参照してください。
- 多要素認証 (MFA) を要求する – AWS アカウント で IAM ユーザーまたはルートユーザーを要求するシナリオがある場合は、セキュリティを強化するために MFA をオンにします。API オペレー

シオンが呼び出されるときに MFA を必須にするには、ポリシーに MFA 条件を追加します。詳細については、『IAM ユーザーガイド』の「[MFA 保護 API アクセスの設定](#)」を参照してください。

IAM でのベストプラクティスの詳細については、『IAM ユーザーガイド』の「[IAM でのセキュリティのベストプラクティス](#)」を参照してください。

#### 特定の AWS Glue オブジェクトにのみ適用されるリソースレベルのアクセス許可

AWS Glue の特定のオブジェクトに限り、詳細なコントロールを定義できます。したがって、Resource ステートメントの Amazon リソースネーム (ARN) を許可する API オペレーションと ARN を許可しない API オペレーションが混在しないようにクライアントの IAM ポリシーを記述する必要があります。

たとえば、次の IAM ポリシーは GetClassifier と GetJobRun の API オペレーションを許可します。AWS Glue は分類子とジョブ実行の ARN を許可しないため、このポリシーでは Resource を \* として定義します。ARN は GetDatabase や GetTable などの特定の API オペレーションに対して許可されるため、ARN はポリシーの後半で指定できます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "glue:GetClassifier*",
 "glue:GetJobRun*"
],
 "Resource": "*"
 },
 {
 "Effect": "Allow",
 "Action": [
 "glue:Get*"
],
 "Resource": [
 "arn:aws:glue:us-east-1:123456789012:catalog",
 "arn:aws:glue:us-east-1:123456789012:database/default",
 "arn:aws:glue:us-east-1:123456789012:table/default/e*1*",
 "arn:aws:glue:us-east-1:123456789012:connection/connection2"
]
 }
]
}
```

```
]
}
```

ARN を許可する AWS Glue オブジェクトのリストについては、「[リソース ARN](#)」を参照してください。

## AWS Glue コンソールの使用

AWS Glue コンソールにアクセスするには、最小限のアクセス許可が必要になります。これらのアクセス許可により、AWS アカウントにある AWS Glue リソースの詳細を一覧で表示できます。最小限必要なアクセス許可よりも制限が厳しいアイデンティティベースのポリシーを作成すると、そのポリシーを持つエンティティ (ユーザーまたはロール) ではコンソールが意図したとおりに機能しません。

AWS CLI または AWS API のみを呼び出すユーザーには、最小限のコンソール権限を付与する必要はありません。代わりに、実行しようとしている API オペレーションに一致するアクションのみへのアクセスを許可します。

ユーザーとロールが引き続き AWS Glue コンソールを使用できるようにするには、これらのエンティティに AWS Glue *ConsoleAccess* または *ReadOnly* AWS マネージドポリシーもアタッチします。詳細については、「IAM ユーザーガイド」の「[ユーザーへのアクセス許可の追加](#)」を参照してください。

ユーザーが AWS Glue コンソールを使用するには、AWS アカウントで AWS Glue リソースの使用を許可する最小限のアクセス許可セットが必要です。これらの AWS Glue アクセス許可に加えて、コンソールでは次のサービスからのアクセス許可が必要になります。

- ログを表示するための Amazon CloudWatch Logs のアクセス許可。
- ロールをリスト化して渡すための AWS Identity and Access Management (IAM) のアクセス許可。
- スタックを操作する AWS CloudFormation のアクセス許可。
- Amazon Elastic Compute Cloud (Amazon EC2) が VPC、サブネット、セキュリティグループ、インスタンス、およびその他のオブジェクトをリストするためのアクセス許可。
- Amazon Simple Storage Service (Amazon S3) が バケットとオブジェクトをリストし、スクリプトを取得および保存するためのアクセス許可。
- クラスターでの作業に必要な Amazon Redshift のアクセス許可。
- インスタンスをリスト化するための Amazon Relational Database Service (Amazon RDS) のアクセス許可。

ユーザーが AWS Glue コンソールを表示して操作するために必要なアクセス許可の詳細については、「[ステップ 3: AWS Glue にアクセスするユーザーまたはグループにポリシーをアタッチする](#)」を参照してください。

これらの最小限必要なアクセス許可よりも制限された IAM ポリシーを作成している場合、その IAM ポリシーを使用するユーザーに対してコンソールは意図したとおりには機能しません。これらのユーザーが引き続き AWS Glue コンソールを使用できるようにするには、[AWS の管理 \(定義済み\) ポリシー AWS Glue](#) で述べたとおり `AWSGlueConsoleFullAccess` マネージドポリシーもアタッチします。

### 自分の権限の表示をユーザーに許可する

この例では、ユーザーアイデンティティにアタッチされたインラインおよびマネージドポリシーの表示を IAM ユーザーに許可するポリシーの作成方法を示します。このポリシーには、コンソールで、または AWS CLI か AWS API を使用してプログラマ的に、このアクションを完了する権限が含まれています。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "ViewOwnUserInfo",
 "Effect": "Allow",
 "Action": [
 "iam:GetUserPolicy",
 "iam:ListGroupsWithUser",
 "iam:ListAttachedUserPolicies",
 "iam:ListUserPolicies",
 "iam:GetUser"
],
 "Resource": ["arn:aws:iam::*:user/${aws:username}"]
 },
 {
 "Sid": "NavigateInConsole",
 "Effect": "Allow",
 "Action": [
 "iam:GetGroupPolicy",
 "iam:GetPolicyVersion",
 "iam:GetPolicy",
 "iam:ListAttachedGroupPolicies",
 "iam:ListGroupPolicies",
 "iam:ListPolicyVersions",

```

```
 "iam:ListPolicies",
 "iam:ListUsers"
],
 "Resource": "*"
}
]
```

## テーブルへの読み取り専用のアクセスを許可する

次のポリシーは、データベース db1 の books テーブルに対する読み取り専用アクセス許可を付与します。Amazon リソースネーム (ARN) の詳細については、「[データカタログの ARN](#)」を参照してください。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "GetTablesActionOnBooks",
 "Effect": "Allow",
 "Action": [
 "glue:GetTables",
 "glue:GetTable"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:catalog",
 "arn:aws:glue:us-west-2:123456789012:database/db1",
 "arn:aws:glue:us-west-2:123456789012:table/db1/books"
]
 }
]
}
```

このポリシーでは db1 という名前のデータベースにある books という名前のテーブルへの読み取り専用許可を付与します。テーブルへの Get アクセスを許可するには、カタログとデータベースリソースへのアクセス許可も必要になります。

次のポリシーは、データベース db1 にテーブル tb1 を作成するための必要最小限のアクセス許可を付与します。

```
{
```

```
"Version": "2012-10-17",
"Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "glue:CreateTable"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:table/db1/tbl1",
 "arn:aws:glue:us-west-2:123456789012:database/db1",
 "arn:aws:glue:us-west-2:123456789012:catalog"
]
 }
]
```

### GetTables アクセス許可によるテーブルのフィルタリング

データベース db1 内に、customers、stores、および store\_sales の 3 つのテーブルがあるとします。次のポリシーは、GetTables アクセス許可を stores および store\_sales に付与しますが、customers には付与しません。このポリシーで GetTables を呼び出すと、結果には 2 つの認可されたテーブルのみが含まれます (customers テーブルは返されません)。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "GetTablesExample",
 "Effect": "Allow",
 "Action": [
 "glue:GetTables"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:catalog",
 "arn:aws:glue:us-west-2:123456789012:database/db1",
 "arn:aws:glue:us-west-2:123456789012:table/db1/store_sales",
 "arn:aws:glue:us-west-2:123456789012:table/db1/stores"
]
 }
]
}
```

store\* を使用して store で始まる任意のテーブル名に一致させることで、前述のポリシーを簡素化できます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "GetTablesExample2",
 "Effect": "Allow",
 "Action": [
 "glue:GetTables"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:catalog",
 "arn:aws:glue:us-west-2:123456789012:database/db1",
 "arn:aws:glue:us-west-2:123456789012:table/db1/store*"
]
 }
]
}
```

同様に、/db1/\* を使用して db1 内のすべてのテーブルと一致させることで、次のポリシーによって db1 内のすべてのテーブルへのアクセスが GetTables に付与されます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "GetTablesReturnAll",
 "Effect": "Allow",
 "Action": [
 "glue:GetTables"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:catalog",
 "arn:aws:glue:us-west-2:123456789012:database/db1",
 "arn:aws:glue:us-west-2:123456789012:table/db1/*"
]
 }
]
}
```

テーブルの ARN を指定しない場合、GetTables の呼び出しは成功しますが、空のリストが返されます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "GetTablesEmptyResults",
 "Effect": "Allow",
 "Action": [
 "glue:GetTables"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:catalog",
 "arn:aws:glue:us-west-2:123456789012:database/db1"
]
 }
]
}
```

ポリシーでデータベースの ARN が欠落している場合は、GetTables の呼び出しが失敗して AccessDeniedException が発生します。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "GetTablesAccessDeny",
 "Effect": "Allow",
 "Action": [
 "glue:GetTables"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:catalog",
 "arn:aws:glue:us-west-2:123456789012:table/db1/*"
]
 }
]
}
```

## テーブルとすべてのパーティションへのフルアクセスを付与する

次のポリシーは、データベース db1 で books という名前のテーブルのすべての許可を付与します。これには、テーブル自体、テーブルのアーカイブされたバージョン、テーブルのすべてのパーティションでの読み込みおよび書き込みアクセス許可が含まれます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "FullAccessOnTable",
 "Effect": "Allow",
 "Action": [
 "glue:CreateTable",
 "glue:GetTable",
 "glue:GetTables",
 "glue:UpdateTable",
 "glue>DeleteTable",
 "glue:BatchDeleteTable",
 "glue:GetTableVersion",
 "glue:GetTableVersions",
 "glue>DeleteTableVersion",
 "glue:BatchDeleteTableVersion",
 "glue:CreatePartition",
 "glue:BatchCreatePartition",
 "glue:GetPartition",
 "glue:GetPartitions",
 "glue:BatchGetPartition",
 "glue:UpdatePartition",
 "glue>DeletePartition",
 "glue:BatchDeletePartition"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:catalog",
 "arn:aws:glue:us-west-2:123456789012:database/db1",
 "arn:aws:glue:us-west-2:123456789012:table/db1/books"
]
 }
]
}
```

前述のポリシーは、実際に次のように簡素化できます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "FullAccessOnTable",
 "Effect": "Allow",
 "Action": [
 "glue:*Table*",
 "glue:*Partition*"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:catalog",
 "arn:aws:glue:us-west-2:123456789012:database/db1",
 "arn:aws:glue:us-west-2:123456789012:table/db1/books"
]
 }
]
}
```

きめ細かなアクセスコントロールの最小粒度はテーブルレベルであることに注意してください。つまり、テーブル内の一部のパーティションへのアクセス許可をユーザーに付与できませんが、他のパーティションには付与できます。また、一部のテーブルの列には付与できませんが、他のテーブルの列には付与できます。ユーザーはすべてのテーブルへのアクセス許可を持っているか、どのアクセス許可も持っていません。

#### 名前のプレフィックスと明示的な拒否によりアクセスを制御する

ここでは、AWS Glue Data Catalog 内にあるデータベースとテーブルが、名前のプレフィックスを使用して整理されていると想定します。開発ステージのデータベースには、名前のプレフィックス `dev-` があります。実稼働環境のデータベースには、名前のプレフィックス `prod-` があります。次のポリシーを使用して、開発者に `dev-` プレフィックスがある、すべてのデータベース、テーブル、UDF などへのフルアクセスを付与できます。ただし、`prod-` プレフィックスを使用して、すべての読み取り専用アクセス許可を付与します。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "DevAndProdFullAccess",
 "Effect": "Allow",
 "Action": [
```

```

 "glue:*Database*",
 "glue:*Table*",
 "glue:*Partition*",
 "glue:*UserDefinedFunction*",
 "glue:*Connection*"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:catalog",
 "arn:aws:glue:us-west-2:123456789012:database/dev-*",
 "arn:aws:glue:us-west-2:123456789012:database/prod-*",
 "arn:aws:glue:us-west-2:123456789012:table/dev-*/*",
 "arn:aws:glue:us-west-2:123456789012:table/*/dev-*",
 "arn:aws:glue:us-west-2:123456789012:table/prod-*/*",
 "arn:aws:glue:us-west-2:123456789012:table/*/prod-*",
 "arn:aws:glue:us-west-2:123456789012:userDefinedFunction/dev-*/*",
 "arn:aws:glue:us-west-2:123456789012:userDefinedFunction/*/dev-*",
 "arn:aws:glue:us-west-2:123456789012:userDefinedFunction/prod-*/*",
 "arn:aws:glue:us-west-2:123456789012:userDefinedFunction/*/prod-*",
 "arn:aws:glue:us-west-2:123456789012:connection/dev-*",
 "arn:aws:glue:us-west-2:123456789012:connection/prod-*"
]
},
{
 "Sid": "ProdWriteDeny",
 "Effect": "Deny",
 "Action": [
 "glue:*Create*",
 "glue:*Update*",
 "glue:*Delete*"
],
 "Resource": [
 "arn:aws:glue:us-west-2:123456789012:database/prod-*",
 "arn:aws:glue:us-west-2:123456789012:table/prod-*/*",
 "arn:aws:glue:us-west-2:123456789012:table/*/prod-*",
 "arn:aws:glue:us-west-2:123456789012:userDefinedFunction/prod-*/*",
 "arn:aws:glue:us-west-2:123456789012:userDefinedFunction/*/prod-*",
 "arn:aws:glue:us-west-2:123456789012:connection/prod-*"
]
}
]
}
}

```

前述のポリシーの 2 番目のステートメントは、明示的な deny を使用します。明示的な deny を使用して、プリンシパルに付与された任意の allow アクセス許可を上書きできます。これにより、重要なリソースへのアクセスをロックダウンして、別のポリシーによって重要なリソースへのアクセスが誤って付与されることを防ぐことができます。

前述の例では、最初のステートメントによって prod- リソースへのフルアクセスが付与されていますが、2 番目のステートメントによって明示的にリソースへの書き込みアクセスが取り消され、prod- リソースへの読み取りアクセスのみが残されています。

### タグを使用してアクセスを許可する

たとえば、お使いのアカウントの Tom という名前の特定のユーザーに対して、トリガー t2 へのアクセスを制限するとします。その他すべてのユーザー (Sam など) は、トリガー t1 にアクセスできます。トリガー t1 と t2 には、以下のプロパティがあります。

```
aws glue get-triggers
{
 "Triggers": [
 {
 "State": "CREATED",
 "Type": "SCHEDULED",
 "Name": "t1",
 "Actions": [
 {
 "JobName": "j1"
 }
],
 "Schedule": "cron(0 0/1 * * ? *)"
 },
 {
 "State": "CREATED",
 "Type": "SCHEDULED",
 "Name": "t2",
 "Actions": [
 {
 "JobName": "j1"
 }
],
 "Schedule": "cron(0 0/1 * * ? *)"
 }
]
}
```

AWS Glue 管理者は、トリガー t2 にタグ値 Tom (`aws:ResourceTag/Name": "Tom"`) を添付しました。AWS Glue 管理者は、タグに基づく条件ステートメントを使用して Tom に IAM ポリシーも提供しました。その結果、Tom はタグ値 Tom を持つリソースを対象とした AWS Glue オペレーションのみ使用できます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": "glue:*",
 "Resource": "*",
 "Condition": {
 "StringEquals": {
 "aws:ResourceTag/Name": "Tom"
 }
 }
 }
]
}
```

Tom がトリガー t1 にアクセスしようとする、アクセス拒否メッセージが返されます。一方、トリガー t2 は正常に取得できます。

```
aws glue get-trigger --name t1
```

```
An error occurred (AccessDeniedException) when calling the GetTrigger operation:
User: Tom is not authorized to perform: glue:GetTrigger on resource: arn:aws:glue:us-east-1:123456789012:trigger/t1
```

```
aws glue get-trigger --name t2
```

```
{
 "Trigger": {
 "State": "CREATED",
 "Type": "SCHEDULED",
 "Name": "t2",
 "Actions": [
 {
 "JobName": "j1"
 }
],
 "Schedule": "cron(0 0/1 * * ? *)"
 }
}
```

```
 }
 }
}
```

この API オペレーションではタグのフィルタリングがサポートされていないため、Tom は、複数の GetTriggers API オペレーションを使用してトリガーを一覧表示することができません。

Tom に GetTriggers へのアクセス権を与えるには、AWS Glue 管理者はアクセス許可を 2 つのセクションに分割するポリシーを作成します。1 つのセクションでは、GetTriggers API オペレーションを使用してすべてのトリガーにアクセスすることを Tom に許可します。2 番目のセクションでは、Tom という値でタグ付けされている API オペレーションにアクセスすることを Tom に許可します。このポリシーでは、Tom はトリガー t2 に対する GetTriggers と GetTrigger の両方のアクセスが許可されます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": "glue:GetTriggers",
 "Resource": "*"
 },
 {
 "Effect": "Allow",
 "Action": "glue:*",
 "Resource": "*",
 "Condition": {
 "StringEquals": {
 "aws:ResourceTag/Name": "Tom"
 }
 }
 }
]
}
```

## タグを使用してアクセスを拒否する

リソースポリシーのもう 1 つの方法は、リソースへのアクセスを明示的に拒否することです。

**⚠ Important**

明示的な拒否ポリシーは、GetTriggers など複数の API オペレーションに使用することはできません。

次のポリシー例では、すべての AWS Glue ジョブ操作が許可されています。ただし、2 番目の Effect ステートメントでは、Team キーと Special 値がタグ付けされたジョブへのアクセスが明示的に拒否されています。

管理者が次のポリシーをアイデンティティにアタッチすると、このアイデンティティは、Team キーと Special 値がタグ付けされたジョブを除くすべてのジョブにアクセスできます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": "glue:*",
 "Resource": "arn:aws:glue:us-east-1:123456789012:job/*"
 },
 {
 "Effect": "Deny",
 "Action": "glue:*",
 "Resource": "arn:aws:glue:us-east-1:123456789012:job/*",
 "Condition": {
 "StringEquals": {
 "aws:ResourceTag/Team": "Special"
 }
 }
 }
]
}
```

### リストとバッチ API オペレーションでタグを使用する

リソースポリシーを書き込む 3 つ目の方法としては、List API オペレーションを使用してリソースにアクセスし、タグ値のリソースを一覧表示することを許可します。次に、対応する Batch API オペレーションを使用して、特定のリソースの詳細にアクセスすることを許可します。この方法では、管理者は複数形の GetCrawlers、GetDevEndpoints、GetJobs、または GetTriggers API オ

ペレーションへのアクセスを許可する必要がありません。代わりに、次の API オペレーションを使用してリソースを一覧表示することを許可できます。

- ListCrawlers
- ListDevEndpoints
- ListJobs
- ListTriggers

また、以下の API オペレーションを使用して、個々のリソースに関する詳細を取得することを許可できます。

- BatchGetCrawlers
- BatchGetDevEndpoints
- BatchGetJobs
- BatchGetTriggers

管理者は、この方法を使用して、次の操作を実行できます。

1. クローラ、開発エンドポイント、ジョブ、およびトリガーにタグを追加します。
2. GetCrawlers、GetDevEndpoints、GetJobs、および GetTriggers などの Get API オペレーションへのユーザーアクセスを拒否します。
3. アクセスできるタグ付きリソースをユーザーが見つげられるようにするには、ListCrawlers、ListDevEndpoints、ListJobs、および ListTriggers などの List API オペレーションへのユーザーアクセスを許可します。
4. TagResource や UntagResource などの AWS Glue タグ付け API へのユーザーアクセスを拒否します。
5. BatchGetCrawlers、BatchGetDevEndpoints、BatchGetJobs、および BatchGetTriggers などの BatchGet API オペレーションを使用して、リソースの詳細に対するアクセスをユーザーに許可します。

たとえば、ListCrawlers オペレーションを呼び出すときに、ユーザー名と一致するタグ値を指定します。結果として、指定したタグ値と一致するクローラのリストが返されます。指定されたタグを持つ各クローラの詳細を取得するには、BatchGetCrawlers に名前リストを提供します。

例えば、Tom に、Tom でタグ付けされたトリガーの詳細の取得のみを許可する場合、管理者は、Tom のトリガーにタグを追加し、すべてのユーザーに対して GetTriggers API オペレーションへのアクセスを拒否し、ListTriggers および BatchGetTriggers へのすべてのユーザーのアクセスを許可することができます。

次に示すのは、AWS Glue 管理者が Tom に付与するリソースポリシーです。ポリシーの最初のセクションでは、AWS Glue API オペレーションは GetTriggers に対して拒否されています。ポリシーの 2 番目のセクションでは、ListTriggers がすべてのリソースに対して許可されています。ただし、3 番目のセクションでは、Tom でタグ付けされたリソースは BatchGetTriggers アクセスでアクセスが許可されています。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Deny",
 "Action": "glue:GetTriggers",
 "Resource": "*"
 },
 {
 "Effect": "Allow",
 "Action": [
 "glue:ListTriggers"
],
 "Resource": [
 "*"
]
 },
 {
 "Effect": "Allow",
 "Action": [
 "glue:BatchGetTriggers"
],
 "Resource": [
 "*"
],
 "Condition": {
 "StringEquals": {
 "aws:ResourceTag/Name": "Tom"
 }
 }
 }
]
}
```

```
]
}
```

前の例と同じトリガーを使用して、Tom はトリガー t2 にはアクセスできますが、トリガー t1 にはアクセスできません。次の例は、Tom が BatchGetTriggers で t1 と t2 にアクセスしようとしたときの結果を示しています。

```
aws glue batch-get-triggers --trigger-names t2
{
 "Triggers": {
 "State": "CREATED",
 "Type": "SCHEDULED",
 "Name": "t2",
 "Actions": [
 {
 "JobName": "j2"
 }
],
 "Schedule": "cron(0 0/1 * * ? *)"
 }
}
```

```
aws glue batch-get-triggers --trigger-names t1
```

```
An error occurred (AccessDeniedException) when calling the BatchGetTriggers operation:
No access to any requested resource.
```

次の例は、Tom が同じ BatchGetTriggers 呼び出しでトリガー t2 とトリガー t3 (存在しない) の両方にアクセスしようとしたときの結果を示しています。Tom は t2 をトリガーするアクセス権を持っており、それが存在するため、t2 のみが返されることに注意してください。Tom はトリガー t3 へのアクセスを許可されていますが、トリガー t3 は存在しないため、t3 はレスポンスの "TriggersNotFound": [] のリストで返されます。

```
aws glue batch-get-triggers --trigger-names t2 t3
{
 "Triggers": {
 "State": "CREATED",
 "Type": "SCHEDULED",
 "Name": "t2",
 "Actions": [
 {
 "JobName": "j2"
 }
]
 }
}
```

```
 }
],
 "TriggersNotFound": ["t3"],
 "Schedule": "cron(0 0/1 * * ? *)"
}
}
```

条件キーまたはコンテキストキーを使用して設定を制御する

ジョブを作成および更新する権限を付与する場合は、条件キーまたはコンテキストキーを使用できません。ここでは、キーについて説明します。

- [条件キーを使って設定を制御するポリシーを制御する](#)
- [条件キーを使って設定を制御するポリシーを制御する](#)

条件キーを使って設定を制御するポリシーを制御する

AWS Glue は glue:VpcIds、glue:SubnetIds、glue:SecurityGroupIds の3つの IAM 条件キーを提供します。ジョブを作成および更新する権限を付与する場合は、IAM ポリシーで条件キーを使用できます。この設定を使用して、目的の VPC 環境外で実行するジョブまたはセッションが作成されない (または更新されない) ようにすることができます。VPC 設定の情報は、CreateJob リクエストから直接入力されるのではなく、AWS Glue 接続を指すジョブの「接続」フィールドから推測されます。

使用例

希望する VpcId 「vpc-id1234」、SubnetIds、SecurityGroupIds を持つ 「traffic-monitored-connection」という名前の AWS Glue ネットワーク型接続を作成します。

IAM ポリシーのアクション CreateJob および UpdateJob の条件キーを指定します。

```
{
 "Effect": "Allow",
 "Action": [
 "glue:CreateJob",
 "glue:UpdateJob"
],
 "Resource": [
 "*"
],
 "Condition": {
 "ForAnyValue:StringLike": {
```

```
 "glue:VpcIds": [
 "vpc-id1234"
]
 }
}
}
```

同様の IAM ポリシーを作成して、接続情報を指定しない AWS Glue ジョブの作成を禁止することができます。

### VPC でのセッションの制限

指定した VPC 内で実行するように作成済みセッションを強制するには、`glue:vpc-id` が `vpc-<123>` と等しくないという条件付きで `glue:CreateSession` アクションに `Deny effect` を追加することで、ロールのアクセス許可を制限します。例:

```
"Effect": "Deny",
"Action": [
 "glue:CreateSession"
],
"Condition": {
 "StringNotEquals" : {"glue:VpcIds" : ["vpc-123"]}
}
```

また、`glue:vpc-id` が `NULL` であるという条件付きで `glue:CreateSession` アクションに `Deny` 効果を追加しても、VPC 内で実行するように作成済みセッションを強制できます。例:

```
{
 "Effect": "Deny",
 "Action": [
 "glue:CreateSession"
],
 "Condition": {
 "Null": {"glue:VpcIds": true}
 }
},
{
 "Effect": "Allow",
 "Action": [
 "glue:CreateSession"
],
```

```
"Resource": ["*"]
}
```

### 条件キーを使って設定を制御するポリシーを制御する

AWS Glue は、`glue:CredentialIssuingService= glue.amazonaws.com` でジョブや開発者のエンドポイントから利用できるように、各ロールセッションにコンテキストキー(AWS Glue)を提供します。これにより、AWS Glue スクリプトが実行するアクションに対し、セキュリティ制御を実装できます。AWS Glue は、別のコンテキストキー(`glue:RoleAssumedBy=glue.amazonaws.com`) を各ロールセッションに提供し、そこで AWS Glue が顧客に代わって別のサービス AWS を (ジョブ/デバイスのエンドポイントではなく AWS Glue サービスにより直接) 呼び出します。

### 使用例

IAM ポリシーで条件付きアクセス許可を指定し、これを AWS Glue ジョブが使用するロールにアタッチします。これで、特定のアクションが、ロールセッションが AWS Glue ジョブランタイム環境で使用されているかどうかに基づいて、許可/拒否されます。

```
{
 "Effect": "Allow",
 "Action": "s3:GetObject",
 "Resource": "arn:aws:s3:::confidential-bucket/*",
 "Condition": {
 "StringEquals": {
 "glue:CredentialIssuingService": "glue.amazonaws.com"
 }
 }
}
```

### ID によるデータプレビューセッションの作成を拒否する

このセクションには、ID によるデータプレビューセッションの作成を拒否するために使用される IAM ポリシーの例が含まれています。このポリシーを ID にアタッチします。ID は、データプレビューセッションの実行中に使用するロールとは別のものです。

```
{
 "Sid": "DatapreviewDeny",
 "Effect": "Deny",
 "Action": [
 "glue:CreateSession"
]
}
```

```
],
 "Resource": [
 "arn:aws:glue:*:*:session/glue-studio-datapreview*"
]
}
```

## AWS Glue のリソースベースのポリシーの例

このセクションでは、リソースベースのポリシーの例を紹介します。これには、クロスアカウントアクセスを許可するポリシーが含まれます。

次の例では、AWS Command Line Interface (AWS CLI) を使用して AWS Glue サービス API オペレーションを操作します。AWS Glue コンソール、または AWS SDK のいずれかを使用して、これと同じオペレーションを実行できます。

### Important

AWS Glue リソースポリシーを変更することで、アカウント内の既存の AWS Glue ユーザーのアクセス許可を誤って取り消したり、予期しない障害が発生したりする場合があります。これらの例は開発やテストアカウントでのみ試して、変更を行う前に既存のワークフローを壊さないことを確認します。

## トピック

- [AWS Glue でリソースベースのポリシーを使用する際の考慮事項](#)
- [同じアカウント内のアクセスをコントロールするリソースポリシーの使用](#)

## AWS Glue でリソースベースのポリシーを使用する際の考慮事項

### Note

IAM ポリシーと AWS Glue リソースポリシーが、それぞれ反映されるまでには数秒かかります。新しいポリシーをアタッチした後、新しいポリシーがシステム全体に反映されるまで古いポリシーが有効なままの場合があります。

JSON 形式で記述されたポリシードキュメントを使用して、リソースポリシーを作成または変更します。ポリシー構文は、以下の例外を除き、アイデンティティベースの IAM ポリシーと同じです ([IAM JSON ポリシーリファレンス](#)を参照のこと)。

- ポリシーステートメントごとに "Principal" または "NotPrincipal" ブロックが必要です。
- "Principal" または "NotPrincipal" では、有効な既存のプリンシパルを特定する必要があります。ワイルドカードパターン (arn:aws:iam::*account-id*:user/\* など) は使用できません。
- ポリシー内の "Resource" ブロックでは、すべてのリソース ARN が、次の正規表現構文と一致している必要があります (最初の %s は *region*、2 番目の %s は *account-id*)。

```
arn:aws:glue:%s:%s:(|[a-zA-Z*]+\/*?.*)
```

たとえば、arn:aws:glue:us-west-2:*account-id*:\* と arn:aws:glue:us-west-2:*account-id*:database/default はいずれも許可されますが、\* は許可されません。

- アイデンティティベースのポリシーとは異なり、AWS Glue リソースポリシーには、このポリシーがアタッチされているカタログに属するリソースの、Amazon リソースネーム (ARN) のみを含めることができます。このような ARN は常に arn:aws:glue: で始まります。
- ポリシーは、それを作成しているアイデンティティを、それ以降のポリシーの作成または変更から排除することはできません。
- リソースポリシーの JSON ドキュメントのサイズは 10 KB を超えることはできません。

## 同じアカウント内のアクセスをコントロールするリソースポリシーの使用

この例では、Account A の管理者ユーザーが、Account A 内の IAM ユーザー Alice にカタログへの完全なアクセス許可を付与するリソースポリシーを作成します。Alice に添付されている IAM ポリシーはありません。

これを行うには、管理者ユーザーは次の AWS CLI コマンドを実行します。

```
Run as admin of Account A
$ aws glue put-resource-policy --profile administrator-name --region us-west-2 --
policy-in-json '{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Principal": {
 "AWS": [
 "arn:aws:iam::account-A-id:user/Alice"
]
 },
 "Effect": "Allow",
```

```

 "Action": [
 "glue:*"
],
 "Resource": [
 "arn:aws:glue:us-west-2:account-A-id:*"
]
 }
]
}'

```

AWS CLI コマンドの一部として JSON ポリシードキュメントを入力する代わりに、ポリシードキュメントをファイルに保存して、file:// でプレフィックスが付けられた AWS CLI コマンドのファイルパスをリファレンスできます。その方法の例は次のとおりです。

```

$ echo '{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Principal": {
 "AWS": [
 "arn:aws:iam::account-A-id:user/Alice"
]
 },
 "Effect": "Allow",
 "Action": [
 "glue:*"
],
 "Resource": [
 "arn:aws:glue:us-west-2:account-A-id:*"
]
 }
]
}' > /temp/policy.json

$ aws glue put-resource-policy --profile admin1 \
 --region us-west-2 --policy-in-json file:///temp/policy.json

```

このリソースポリシーが反映されると、Alice は Account A のすべての AWS Glue リソースにアクセスできます。

```

Run as user Alice
$ aws glue create-database --profile alice --region us-west-2 --database-input '{

```

```
"Name": "new_database",
"Description": "A new database created by Alice",
"LocationUri": "s3://my-bucket"
}'

$ aws glue get-table --profile alice --region us-west-2 --database-name "default" --
table-name "tbl1"}
```

Alice の `get-table` 呼び出しに応じて、AWS Glue サービスは次を返します。

```
{
 "Table": {
 "Name": "tbl1",
 "PartitionKeys": [],
 "StorageDescriptor": {

 },

 }
}
```

## AWSAWS Glue の管理ポリシー

AWS 管理ポリシーは、によって作成および管理されるスタンドアロンのポリシーです。AWS AWS 管理ポリシーは、ユーザー、グループ、ロールにアクセス権限を割り当てることができるように、多くの一般的な使用事例にアクセス許可を与えるように設計されています。

AWS 管理ポリシーでは、AWS すべての顧客が使用できるようになっているため、特定のユースケースでは最小権限のアクセス権限が付与されない場合があることに注意してください。ユースケース別に[カスタマーマネージドポリシー](#)を定義して、マネージドポリシーを絞り込むことをお勧めします。

AWS 管理ポリシーで定義されている権限は変更できません。AWS 管理ポリシーで定義されている権限を更新すると AWS、その更新はポリシーがアタッチされているすべてのプリンシパル ID (ユーザー、グループ、ロール) に影響します。AWS AWS 管理ポリシーが更新される可能性が最も高いのは、新しい API 操作が既存のサービスで開始されたときや、新しい API AWS のサービス操作が使用可能になったときです。

詳細については、「IAM ユーザーガイド」の「[AWS マネージドポリシー](#)」を参照してください。

## AWS の管理 (定義済み) ポリシー AWS Glue

AWS によって作成および管理されるスタンドアロンの IAM ポリシーを提供することで、多くの一般的なユースケースに対応します。AWS これらの管理ポリシーは、一般的なユースケースに必要なアクセス権限を付与するため、どのような権限が必要かを調査する必要がありません。詳細については、「IAM ユーザーガイド」の「[AWS マネージドポリシー](#)」を参照してください。

アカウント内の ID AWS にアタッチできる以下の管理ポリシーは、AWS Glue ユースケースシナリオに固有のもので、ユースケースシナリオごとにグループ化されています。

- [AWSGlueConsoleFullAccess](#)— ポリシーがアタッチされている ID がを使用している場合、AWS Glue リソースへのフルアクセスを許可します。AWS Management Console このポリシーで指定されたリソースの命名規則に従った場合、ユーザーは完全なコンソール機能を使用できます。このポリシーは、通常 AWS Glue コンソールのユーザーにアタッチされています。
- [AWSGlueServiceRole](#)— AWS Glue ユーザーに代わってさまざまなプロセスを実行する必要があるリソースへのアクセスを許可します。これらのリソースには AWS Glue、Amazon S3、IAM、CloudWatch ログ、Amazon EC2 が含まれます。このポリシーで指定されたリソースの命名規則に従った場合、AWS Glue プロセスは必要なアクセス権限を使用できます。このポリシーは、通常、クローラ、ジョブ、開発エンドポイントを定義するときに指定されたロールにアタッチされます。
- [AwsGlueSessionUserRestrictedServiceRole](#)— AWS Glue セッションを除くすべてのリソースへのフルアクセスを提供します。ユーザーに、自分に関連付けられたインタラクティブセッションのみ、作成と使用を許可します。このポリシーには、AWS Glue AWS Glue AWS 他のサービスのリソースを管理するために必要なその他の権限が含まれています。このポリシーでは、AWS Glue AWS 他のサービスのリソースにタグを追加することもできます。

### Note

セキュリティ上のメリットを最大限に引き出すには、[AWSGlueServiceRole](#)、[AWSGlueConsoleFullAccess](#)、または [AWSGlueConsoleSageMakerNotebookFullAccess](#) のポリシーが割り当てられたユーザーにはこのポリシーを付与しないでください。

- [AwsGlueSessionUserRestrictedPolicy](#)— タグキー「owner」AWS と担当者のユーザー ID と一致する値が提供されている場合のみ、CreateSession API AWS Glue オペレーションを使用してインタラクティブセッションを作成するためのアクセスを提供します。この ID ポリシーは、CreateSession API オペレーションを呼び出す IAM ユーザーにアタッチされます。このポリシーでは、担当者がユーザー ID と一致する「owner」AWS Glue タグと値で作成されたインタラクティブセッションリソースを操作することも許可されます。AWS このポリシーは、セッショ

ンが作成された後に、AWS Glue セッションリソースから「所有者」タグを変更または削除する許可を拒否します。

 Note

セキュリティ上のメリットを最大限に引き出すには、AWSGlueServiceRole、AWSGlueConsoleFullAccess、またはAWSGlueConsoleSageMakerNotebookFullAccess のポリシーが割り当てられたユーザーにはこのポリシーを付与しないでください。

- [AwsGlueSessionUserRestrictedNotebookServiceRole](#)— AWS Glue Studio ノートブックセッションへの十分なアクセスを提供し、AWS Glue特定のインタラクティブセッションリソースを操作できるようにします。これらは、ノートブックを作成したプリンシパル (IAM AWS ユーザーまたはロール) のユーザー ID と一致する「owner」タグ値を使用して作成されたリソースです。これらのタグの詳細については、IAM ユーザーガイドにある[プリンシパルキーの値](#)の表を参照してください。

このサービスロールポリシーは、ノートブック内でマジックコマンドにより指定されたロールと、CreateSession API オペレーションにロールとして渡されたロールにアタッチされます。また、このポリシーでは、タグキー「owner」AWS と値がプリンシパルのユーザー ID と一致する場合のみ、AWS GlueAWS Glue Studioプリンシパルがノートブックインターフェースからインタラクティブセッションを作成することを許可します。このポリシーは、セッションが作成された後に、AWS Glue セッションリソースから「所有者」タグを変更または削除する許可を拒否します。このポリシーには、Amazon S3 バケットへの書き込みと読み取り、CloudWatch ログの書き込み、および使用する Amazon EC2 リソースのタグの作成と削除を行う権限も含まれています。AWS Glue

 Note

セキュリティ上のメリットを最大限に引き出すには、AWSGlueServiceRole、AWSGlueConsoleFullAccess、またはAWSGlueConsoleSageMakerNotebookFullAccess のポリシーが割り当てられたユーザーにはこのポリシーを付与しないでください。

- [AwsGlueSessionUserRestrictedNotebookPolicy](#)— AWS Glue Studio ノートブックを作成したプリンシパル (IAM ユーザーまたはロール) AWS のユーザー ID と一致するタグキー「owner」と値がある場合のみ、AWS Glueノートブックインターフェイスからインタラクティブセッションを作

成するためのアクセスを提供します。これらのタグの詳細については、IAM ユーザーガイドにある[プリンシパルキーの値](#)の表を参照してください。

このポリシーは、AWS Glue Studio ノートブックインターフェイスからセッションを作成するプリンシパル (IAM ユーザーまたはロール) にアタッチされます。このポリシーでは、特定の AWS Glue インタラクティブセッションリソースのやり取りを行うための、AWS Glue Studio ノートブックへの十分なアクセス権限も付与されます。これらは、AWS プリンシパルのユーザー ID と一致する「owner」タグ値で作成されるリソースです。このポリシーは、セッションが作成された後に、AWS Glue セッションリソースから「所有者」タグを変更または削除する許可を拒否します。

- [AWSGlueServiceNotebookRole](#)— AWS Glue AWS Glue Studio ノートブックで開始されたセッションへのアクセスを許可します。このポリシーでは、すべてのセッションのセッション情報を一覧表示して取得できますが、AWS ユーザーは自分のユーザー ID でタグ付けされたセッションの作成と使用のみを許可します。このポリシーは、ID AWS Glue がタグ付けされたセッションリソースの「所有者」タグを変更または削除する権限を拒否します。AWS

このポリシーは、AWS のノートブックインターフェイスを使用してジョブを作成するユーザーに割り当てます。AWS Glue Studio

- [AWSGlueConsoleSageMakerNotebookFullAccess](#)— ポリシーがアタッチされている ID がを使用している場合、AWS Glue SageMaker とリソースへのフルアクセスを許可します AWS Management Console。このポリシーで指定されたリソースの命名規則に従った場合、ユーザーは完全なコンソール機能を使用できます。このポリシーは通常、AWS Glue SageMaker ノートブックを管理するコンソールのユーザーに添付されます。
- [AWSGlueSchemaRegistryFullAccess](#)— ポリシーがアタッチされている ID AWS Management Console AWS CLIがまたはを使用している場合、AWS Glueスキーマレジストリリソースへのフルアクセスを許可します。このポリシーで指定されたリソースの命名規則に従った場合、ユーザーは完全なコンソール機能を使用できます。このポリシーは通常、AWS Glue AWS CLI AWS Glueコンソールのユーザーまたはスキーマレジストリを管理するユーザーに添付されます。
- [AWSGlueSchemaRegistryReadOnlyAccess](#)— ポリシーがアタッチされている ID AWS Management Console AWS CLIがまたはを使用している場合、AWS Glueスキーマレジストリリソースへの読み取り専用アクセスを許可します。このポリシーで指定されたリソースの命名規則に従った場合、ユーザーは完全なコンソール機能を使用できます。このポリシーは通常、AWS Glue AWS CLI AWS Glueコンソールのユーザーまたはスキーマレジストリを使用するユーザーに添付されます。

**Note**

IAM コンソールにサインインし、特定のポリシーを検索することで、これらの許可ポリシーを確認できます。

これに加え、独自のカスタム IAM ポリシーを作成して、AWS Glue アクションとリソースにアクセス許可を付与することもできます。これらのカスタムポリシーは、それらのアクセス許可が必要な IAM ユーザーまたはグループにアタッチできます。

## AWS Glue AWS を管理ポリシーに付加

このサービスが変更の追跡を開始して以降の AWS Glue AWS の管理ポリシーの更新に関する詳細を表示します。このページへの変更に関する自動通知を受け取るには、AWS Glue ドキュメント履歴ページの RSS フィードを購読してください。

変更	説明	日付
AwsGlueSessionUserRestrictedPolicy — 既存のポリシーのマイナーアップデート。	glue:StartCompletion と glue:GetCompletion をポリシーに追加します。AWS Glue で Amazon Q データを統合する場合に必要です。	2024 年 4 月 30 日
AwsGlueSessionUserRestrictedNotebookServiceRole — 既存のポリシーのマイナーアップデート。	glue:StartCompletion と glue:GetCompletion をポリシーに追加します。AWS Glue で Amazon Q データを統合する場合に必要です。	2024 年 4 月 30 日
AwsGlueSessionUserRestrictedServiceRole — 既存のポリシーのマイナーアップデート。	glue:StartCompletion と glue:GetCompletion をポリシーに追加します。AWS Glue で Amazon Q データを統合する場合に必要です。	2024 年 4 月 30 日

変更	説明	日付
	タを統合する場合に必要です。	
AWSGlueServiceNotebookRole — 既存のポリシーのマイナーアップデート。	glue:StartCompletion と glue:GetCompletion をポリシーに追加します。AWS Glue で Amazon Q データを統合する場合に必要です。	2024 年 1 月 30 日
AwsGlueSessionUserRestrictedNotebookPolicy — 既存のポリシーのマイナーアップデート。	glue:StartCompletion と glue:GetCompletion をポリシーに追加します。AWS Glue で Amazon Q データを統合する場合に必要です。	2023 年 11 月 29 日
AWSGlueServiceNotebookRole — 既存のポリシーのマイナーアップデート。	codewhisperer:GenerateRecommendations をポリシーに追加します。AWS Glue CodeWhisperer がレコメンデーションを生成する新機能が必要です。	2023 年 10 月 9 日
AWSGlueServiceRole — 既存のポリシーのマイナーアップデート。	AWS Glue CloudWatch のロギングをより適切に反映するように権限の範囲を厳しくしてください。	2023 年 8 月 4 日
AWSGlueConsoleFullAccess — 既存のポリシーのマイナーアップデート。	databrew レシピのリストと説明のアクセス許可をポリシーに追加します。AWS Glue がレシピにアクセスできる新機能の完全な管理アクセスを提供するために必要です。	2023 年 5 月 9 日

変更	説明	日付
<p>AWSGlueConsoleFullAccess — 既存のポリシーのマイナーアップデート。</p>	<p>cloudformation:ListStacks をポリシーに追加します。AWS CloudFormation 承認要件が変更された後も、既存の機能を維持します。</p>	<p>2023 年 3 月 28 日</p>
<p>インタラクティブセッション機能で新たに追加されたマネージドポリシー:</p> <ul style="list-style-type: none"> <li>• AwsGlueSessionUserRestrictedServiceRole</li> <li>• AwsGlueSessionUserRestrictedPolicy</li> <li>• AwsGlueSessionUserRestrictedNotebookServiceRole</li> <li>• AwsGlueSessionUserRestrictedNotebookPolicy</li> </ul>	<p>これらのポリシーは、AWS Glue Studio のインタラクティブセッションとノートブックのセキュリティを強化するように設計されています。これらのポリシーは、所有者だけがアクセスできるように、CreateSession API オペレーションへのアクセスを制限します。</p>	<p>2021 年 11 月 30 日</p>

変更	説明	日付
<p>AWSGlueConsoleSageMakerNotebookFullAccess — 既存のポリシーを更新します。</p>	<p>arn:aws:s3:::aws-glue-*/* がスクリプトと一時ファイルを保存する際に使用する、Amazon S3 バケットに対する読み込み/書き込みのアクセス許可を付与するアクションで、冗長していたリソース ARN (AWS Glue) を削除しました。</p> <p>"StringEquals" を "ForAnyValue:StringLike" に変更し、順序に誤りがあった (各所の "Effect": "Allow" 行を "Action": 行の前に移動することにより、構文の問題を修正しました。</p>	<p>2021 年 7 月 15 日</p>
<p>AWSGlueConsoleFullAccess — 既存のポリシーへの更新。</p>	<p>arn:aws:s3:::aws-glue-*/* がスクリプトと一時ファイルを保存する際に使用する、Amazon S3 バケットに対する読み込み/書き込みのアクセス許可を付与するアクションで、冗長していたリソース ARN (AWS Glue) を削除しました。</p>	<p>2021 年 7 月 15 日</p>
<p>AWS Glue は変更の追跡を開始しました</p>	<p>AWS Glue AWS 管理ポリシーの変更の追跡を開始しました。</p>	<p>2021 年 6 月 10 日</p>

## AWS Glue リソース ARN の指定

AWS Glue では、AWS Identity and Access Management (IAM) ポリシーを使用してリソースへのアクセスを制御できます。ポリシーで Amazon リソースネーム (ARN) を使用して、ポリシーを適用するリソースを識別します。AWS Glue のすべてのリソースで ARN がサポートされているわけではありません。

### トピック

- [データカタログの ARN](#)
- [AWS Glue にあるカタログ以外のオブジェクトの ARN](#)
- [AWS Glue のカタログ以外の単一の API オペレーション](#)
- [複数の項目を取得する AWS Glue のカタログ以外の API オペレーションのアクセスコントロール](#)
- [AWS Glue のカタログ以外のバッチ取得 API オペレーションのアクセス制御](#)

### データカタログの ARN

Data Catalog リソースには、catalog をルートとする階層構造があります。

```
arn:aws:glue:region:account-id:catalog
```

各 AWS アカウントは、カタログ ID として 12 桁のアカウント ID を使用する単一の Data Catalog を、AWS リージョン内に持っています。リソースには、次の表に示すように、一意の ARN が関連付けられています。

リソースタイプ:	ARN 形式
カタログ	arn:aws:glue: <i>region</i> : <i>account-id</i> :catalog 例: arn:aws:glue:us-east-1:123456789012:catalog
データベース	arn:aws:glue: <i>region</i> : <i>account-id</i> :database/ <i>database name</i> 例: arn:aws:glue:us-east-1:123456789012:database/db1
テーブル	arn:aws:glue: <i>region</i> : <i>account-id</i> :table/ <i>database name</i> / <i>table name</i>

リソースタイプ:	ARN 形式 例: arn:aws:glue:us-east-1:123456789012:table/db1/tb11
ユーザー定義関数	arn:aws:glue: <i>region:account-id</i> :userDefinedFunction/ <i> database name/user-defined function name</i> 例: arn:aws:glue:us-east-1:123456789012:userDefinedFunction/db1/func1
Connection	arn:aws:glue: <i>region:account-id</i> :connection/ <i> connection name</i> 例: arn:aws:glue:us-east-1:123456789012:connection/connection1
インタラクティブセッション	arn:aws:glue: <i>region:account-id</i> :session/ <i> interactive session id</i> 例: arn:aws:glue:us-east-1:123456789012:session/a1b2c3d4-5678-90ab-cdef-EXAMPLE11111

きめ細かなアクセスコントロールを有効にするには、IAM ポリシーとリソースポリシーでこれらの ARN を使用して、特定のリソースへのアクセス許可を付与または拒否します。ポリシーではワイルドカードを使用できます。たとえば、次の ARN はデータベース default のすべてのテーブルと一致します。

```
arn:aws:glue:us-east-1:123456789012:table/default/*
```

### Important

Data Catalog リソースで実行されるすべてのオペレーションには、リソースとその継承元すべてに対するアクセス許可が必要です。たとえば、テーブルのパーティションを作成するには、テーブル、データベース、およびテーブルが存在するカタログに対するアクセス許可が必要です。以下に、Data Catalog のデータベース PrivateDatabase 内のテーブル PrivateTable で、パーティションを作成するのに必要なアクセス許可の例を示します。

```
{
 "Sid": "GrantCreatePartitions",
 "Effect": "Allow",
 "Action": [
 "glue:BatchCreatePartitions"
],
 "Resource": [
 "arn:aws:glue:us-east-1:123456789012:table/PrivateDatabase/PrivateTable",
 "arn:aws:glue:us-east-1:123456789012:database/PrivateDatabase",
 "arn:aws:glue:us-east-1:123456789012:catalog"
]
}
```

リソースとそのすべての祖先のアクセス許可に加えて、すべての削除オペレーションにはそのリソースのすべての子のアクセス許可が必要です。たとえば、データベースを削除するには、データベースおよびデータベースが存在するカタログに加えて、データベース内のすべてのテーブルおよびユーザー定義関数に対するアクセス許可が必要です。以下は、Data Catalog でデータベース PrivateDatabase を削除するのに必要なアクセス許可の例です。

```
{
 "Sid": "GrantDeleteDatabase",
 "Effect": "Allow",
 "Action": [
 "glue>DeleteDatabase"
],
 "Resource": [
 "arn:aws:glue:us-east-1:123456789012:table/PrivateDatabase/*",
 "arn:aws:glue:us-east-1:123456789012:userDefinedFunction/PrivateDatabase/*",
 "arn:aws:glue:us-east-1:123456789012:database/PrivateDatabase",
 "arn:aws:glue:us-east-1:123456789012:catalog"
]
}
```

要約すると、Data Catalog リソースに対するアクションは、以下のアクセス許可のルールに従います。

- カタログに対するアクションには、カタログのみに対するアクセス許可が必要です。
- データベースに対するアクションには、データベースおよびカタログに対するアクセス許可が必要です。

- データベースに対する削除アクションには、データベースとカタログに対するアクセス許可に加えて、データベース内のすべてのテーブルおよびユーザー定義関数に対するアクセス許可も必要です。
- テーブル、パーティション、またはテーブルバージョンに対するアクションには、テーブル、データベース、およびカタログに対するアクセス許可が必要です。
- ユーザー定義関数に対するアクションには、ユーザー定義関数、データベース、およびカタログに対するアクセス許可が必要です。
- 接続に対するアクションには、接続およびカタログに対するアクセス許可が必要です。

## AWS Glue にあるカタログ以外のオブジェクトの ARN

一部の AWS Glue リソースによって、ARN を使用してアクセスをコントロールするリソースレベルの許可が付与されます。IAM ポリシーでこれらの ARN を使用して、きめ細かなアクセスコントロールを有効にします。次の表は、リソース ARN を含むことができるリソースの一覧です。

リソースタイプ:	ARN 形式
Crawler	arn:aws:glue: <i>region:account-id</i> :crawler/ <i>crawler-name</i>  例: arn:aws:glue:us-east-1:123456789012:crawler/mycrawler
ジョブ	arn:aws:glue: <i>region:account-id</i> :job/ <i>job-name</i>  例: arn:aws:glue:us-east-1:123456789012:job/testjob
Trigger トリガー)	arn:aws:glue: <i>region:account-id</i> :trigger/ <i>trigger-name</i>  例: arn:aws:glue:us-east-1:123456789012:trigger/sampletrigger
開発エンドポイント	arn:aws:glue: <i>region:account-id</i> :devEndpoint/ <i>development-endpoint-name</i>

リソースタイプ:	ARN 形式
	例: arn:aws:glue:us-east-1:123456789012:devEndpoint/temporarydevendpoint
機械学習変換	arn:aws:glue: <i>region:account-id</i> :mlTransform/ <i>transform-id</i>  例: arn:aws:glue:us-east-1:123456789012:mlTransform/tfm-1234567890

## AWS Glue のカタログ以外の単一の API オペレーション

AWS Glue のカタログ以外の単一の API オペレーションは、単一の項目 (開発エンドポイント) で動作します。例は、GetDevEndpoint, CreateUpdateDevEndpoint および UpdateDevEndpoint です。これらのオペレーションでは、ポリシーは API 名を "action" ブロックに、リソース ARN を "resource" ブロックに置く必要があります。

GetDevEndpoint オペレーションの呼び出しをユーザーに許可するとします。次のポリシーは、myDevEndpoint-1 という名前のエンドポイントに必要な最小限のアクセス許可を付与します。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "MinimumPermissions",
 "Effect": "Allow",
 "Action": "glue:GetDevEndpoint",
 "Resource": "arn:aws:glue:us-east-1:123456789012:devEndpoint/myDevEndpoint-1"
 }
]
}
```

次のポリシーは、ワイルドカード (\*) を使用して myDevEndpoint- と一致するリソースへのアクセスを UpdateDevEndpoint に許可します。

```
{
 "Version": "2012-10-17",
 "Statement": [
```

```

 {
 "Sid": "PermissionWithWildcard",
 "Effect": "Allow",
 "Action": "glue:UpdateDevEndpoint",
 "Resource": "arn:aws:glue:us-east-1:123456789012:devEndpoint/myDevEndpoint-
*"
 }
]
}

```

次の例のように 2 つのポリシーを組み合わせることができます。名前が A で始まるすべての開発エンドポイントに対して `EntityNotFoundException` が表示される場合があります。一方、その他の開発エンドポイントにアクセスしようとすると、アクセス拒否エラーが表示されます。

```

{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Sid": "CombinedPermissions",
 "Effect": "Allow",
 "Action": [
 "glue:UpdateDevEndpoint",
 "glue:GetDevEndpoint"
],
 "Resource": "arn:aws:glue:us-east-1:123456789012:devEndpoint/A*"
 }
]
}

```

## 複数の項目を取得する AWS Glue のカタログ以外の API オペレーションのアクセスコントロール

一部の AWS Glue API オペレーション (`GetDevEndpoints` など) は、複数の項目 (複数の開発エンドポイントなど) を取得します。このオペレーションでは、特定の ARN ではなく、ワイルドカード (\*) リソースのみを指定できます。

たとえば、`GetDevEndpoints` をポリシーに含めるには、リソースの範囲をワイルドカード (\*) に限定する必要があります。単一のオペレーション (`GetDevEndpoint`、`CreateDevEndpoint`、および `DeleteDevendpoint`) の範囲も、例のすべての (\*) リソースに限定されます。

```

{

```

```
 "Sid": "PluralAPIIncluded",
 "Effect": "Allow",
 "Action": [
 "glue:GetDevEndpoints",
 "glue:GetDevEndpoint",
 "glue:CreateDevEndpoint",
 "glue:UpdateDevEndpoint"
],
 "Resource": [
 "*"
]
}
```

## AWS Glue のカタログ以外のバッチ取得 API オペレーションのアクセス制御

一部の AWS Glue API オペレーション (BatchGetDevEndpoints など) は、複数の項目 (複数の開発エンドポイントなど) を取得します。このオペレーションでは、アクセス可能なリソースの範囲を制限するために ARN を指定できます。

たとえば、特定の開発エンドポイントへのアクセスを許可するには、そのリソース ARN を指定して BatchGetDevEndpoints をポリシーに含めます。

```
{
 "Sid": "BatchGetAPIIncluded",
 "Effect": "Allow",
 "Action": [
 "glue:BatchGetDevEndpoints"
],
 "Resource": [
 "arn:aws:glue:us-east-1:123456789012:devEndpoint/de1"
]
}
```

このポリシーでは、de1 という名前の開発エンドポイントに正常にアクセスできます。ただし、de2 という名前の開発エンドポイントにアクセスしようとすると、エラーが返されます。

An error occurred (AccessDeniedException) when calling the BatchGetDevEndpoints operation: No access to any requested resource.

**⚠ Important**

List や BatchGet API オペレーションを使用するなど、IAM ポリシーを設定するための別の方法については、「[AWS Glue のアイデンティティベースポリシーの例](#)」を参照してください。

## クロスアカウントアクセス許可の付与

アカウント間で Data Catalog のリソースへのアクセスを許可することで、抽出、変換、ロード (ETL) ジョブで、異なるアカウントからデータをクエリし、そのデータを結合できるようになります。

### トピック

- [クロスアカウントのアクセスを許可する AWS Glue のメソッド](#)
- [Data Catalog リソースポリシーの追加または更新](#)
- [クロスアカウント API コールの実行](#)
- [クロスアカウント ETL コールの実行](#)
- [CloudTrail のクロスアカウントロギング](#)
- [クロスアカウントリソースの所有権および請求](#)
- [クロスアカウントアクセスの制約事項](#)

## クロスアカウントのアクセスを許可する AWS Glue のメソッド

自分のデータに対するアクセスを外部の AWS アカウントに許可するためには AWS Glue メソッドまたは AWS Lake Formation クロスアカウント付与を使用する必要があります。AWS Glue メソッドは AWS Identity and Access Management (IAM) ポリシーを使用して、きめ細かなアクセスコントロールを実現します。Lake Formation では、リレーショナルデータベースシステムでの GRANT/REVOKE コマンドに類似した、よりシンプルなアクセス許可 (GRANT/REVOKE) のモデルを使用します。

このセクションでは、AWS Glue メソッドの使用について説明します。Lake Formation のクロスアカウント付与の使用については、AWS Lake Formation デベロッパーガイドの「[Granting Lake Formation Permissions](#)」を参照してください。

リソースにクロスアカウントアクセスを許可するためには、2 つの AWS Glue メソッドが用意されています。

- Data Catalog のリソースポリシーを使用
- IAM ロール を使用

### リソースポリシーを使用してのクロスアカウントアクセスの許可

Data Catalog リソースポリシーを使用してクロスアカウントアクセスを許可するための、一般的な手順を以下に示します。

1. アカウント A の管理者 (または権限付与されたアイデンティティ) が、リソースポリシーをアカウント A 内の Data Catalog にアタッチします。このポリシーは、アカウント A のカタログにあるリソースに対してオペレーションを実行するための、特定のクロスアカウントアクセス許可をアカウント B に付与します。
2. アカウント B の管理者は、アカウント A から受け取ったアクセス許可を委任する、アカウント B の IAM アイデンティティに IAM ポリシーをアタッチします。

これで、アカウント B のアイデンティティは、アカウント A にある指定されたリソースにアクセスできるようになります。

このアイデンティティがこのリソースにアクセスするには、リソースの所有者 (アカウント A) およびその親アカウント (アカウント B) の両方からアクセスを許可される必要があります。

### IAM ロールを使用したクロスアカウントアクセス許可の付与

以下に、IAM ロールを使用してクロスアカウントアクセスを許可する一般的な手順を示します。

1. リソースを所有するアカウント (アカウント A) 内の管理者 (または権限許可された他のアイデンティティ) が IAM ロールを作成します。
2. アカウント A 内の管理者は、対象のリソースにアクセスするためのクロスアカウントアクセス許可を付与するロールにポリシーをアタッチします。
3. アカウント A 内の管理者により、ロールへの信頼ポリシーのアタッチが行われます。このロールは、別のアカウント (アカウント B) 内の IAM アイデンティティを、そのロールを引き受けることのできるプリンシパルとして識別します。

信頼ポリシーのプリンシパルに、ロールを引き受けるための AWS サービスへのアクセス許可を付与する場合は、このプリンシパルを AWS サービスのプリンシパルとすることができます。

- これで、アカウント B の管理者はアカウント B の 1 つ以上の IAM アイデンティティにアクセス許可を委任して、ロールを引き受けることを許可します。これにより、アカウント B のこれらのアイデンティティに対してアカウント A のリソースに対するアクセス権が付与されます。

IAM を使用したアクセス許可の委任の詳細については、「IAM ユーザーガイド」の「[アクセス管理](#)」を参照してください。ユーザー、グループ、ロール、アクセス権限の詳細については、「IAM ユーザーガイド」の「[ID \(ユーザー、グループ、ロール\)](#)」を参照してください。

この 2 つの方法の比較については、IAM ユーザーガイドの「[IAM ロールとリソースベースのポリシーとの相違点](#)」を参照してください。AWS Glue は両方の方法をサポートしていますが、リソースポリシーが許可できるのはデータカタログリソースへのアクセスのみ、という制限があります。

例えば、アカウント B の Dev ロールにアカウント A のデータベース db1 へのアクセスを許可するときは、次のリソースポリシーをアカウント A のカタログにアタッチします。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "glue:GetDatabase"
],
 "Principal": {"AWS": [
 "arn:aws:iam::account-B-id:role/Dev"
]},
 "Resource": [
 "arn:aws:glue:us-east-1:account-A-id:catalog",
 "arn:aws:glue:us-east-1:account-A-id:database/db1"
]
 }
]
}
```

さらに、アカウント B がアカウント A の db1 に実際にアクセスする前に、アカウント B で次の IAM ポリシーを Dev ロールにアタッチする必要があります。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
```

```
"Effect": "Allow",
"Action": [
 "glue:GetDatabase"
],
"Resource": [
 "arn:aws:glue:us-east-1:account-A-id:catalog",
 "arn:aws:glue:us-east-1:account-A-id:database/db1"
]
}
]
```

## Data Catalog リソースポリシーの追加または更新

コンソール、API、または AWS Glue (AWS Command Line Interface) により、AWS CLI Data Catalog のリソースポリシーを追加したり、更新したりできます。

### Important

自分のアカウントの AWS Lake Formation で、既にクロスアカウントのアクセス許可付与を作成済みの場合は、Data Catalog のリソースポリシーを追加または更新するには追加の手順が必要となります。詳細については、AWS Lake Formation デベロッパーガイドの「[AWS Glue Glue と Lake Formation の両方を使用したクロスアカウント許可の管理](#)」を参照してください。

Lake Formation クロスアカウント許可が存在するかどうかを確認するときには、`glue:GetResourcePolicies` API オペレーションまたは AWS CLI を使用します。この `glue:GetResourcePolicies` が既存のデータカタログポリシー以外のポリシーを返した場合、Lake Formation の許可は存在します。詳細については、AWS Lake Formation デベロッパーガイドの「[Viewing All Cross-Account Grants Using the GetResourcePolicies API](#)」(`GetResourcePolicies` API を使用したすべてのクロスアカウント許可の表示) を参照してください。

Data Catalog のリソースポリシーを追加または更新するには (コンソール)

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。

`glue:PutResourcePolicy` アクセス許可を持つ AWS Identity and Access Management (IAM) 管理ユーザーとしてサインインします。

2. ナビゲーションペインで [設定] を選択します。

3. [Data catalog settings] (データカタログ設定) ページの [Permissions] (アクセス許可) の下にあるテキスト欄に、リソースポリシーを貼り付けます。次に、[Save] (保存) を選択します。

ポリシー内のアクセス許可が、Lake Formation を使用して付与されたアクセス許可に追加されることを示す警告がコンソールに表示される場合は、[Proceed] (続行する) をクリックします。

Data Catalog のリソースポリシーを追加または更新するには (AWS CLI)

- `aws glue put-resource-policy` コマンドを送信します。Lake Formation・ グラントが既に存在する場合には、`--enable-hybrid` オプションに値 'TRUE' を設定する必要があります。

このコマンドの使用例については、「[AWS Glue のリソースベースのポリシーの例](#)」を参照してください。

## クロスアカウント API コールの実行

すべての AWS Glue Data Catalog オペレーションに `CatalogId` フィールドがあります。クロスアカウントアクセスを有効にするために必要なアクセス許可が付与済みであれば、複数のアカウントをまたいで Data Catalog API を呼び出すことが可能です。呼び出し元は、ターゲット AWS アカウント ID を `CatalogId` で渡すことで、ターゲットアカウントのリソースに対するアクセスを実行します。

`CatalogId` 値を指定しない場合、AWS Glue はデフォルトで発信者自身のアカウント ID を使用するため、コールはクロスアカウントにはなりません。

## クロスアカウント ETL コールの実行

一部の AWS Glue PySpark および Scala API には、カタログ ID フィールドがあります。クロスアカウントアクセスを有効にする、すべての必要なアクセス許可が付与されていれば、ETL ジョブはアカウントをまたがって、PySpark および Scala による API オペレーション呼び出しを行えます。これには、ターゲットアカウント内の Data Catalog リソースにアクセスするように、ターゲットの AWS アカウント ID をカタログ ID フィールドで指定し受け渡す必要があります。

カタログ ID 値を指定しない場合、AWS Glue はデフォルトで発信者自身のアカウント ID を使用するため、コールはクロスアカウントにはなりません。

catalog\_id をサポートする PySpark API については、「[GlueContext クラス](#)」を参照してください。catalogId をサポートする Scala API については、「[AWS Glue スカラ API GlueContext](#)」を参照してください。

次の例は、ETL ジョブを実行するために被付与者に必要となるアクセス許可を示しています。この例で、*grantee-account-id* はジョブを実行するクライアントの catalog-id、*grantor-account-id* はリソースの所有者です。この例では、権限付与者のアカウントにあるすべてのカタログリソースに対するアクセス許可を付与します。付与されたリソースの範囲を制限するには、カタログ、データベース、テーブル、および接続の特定の ARN を指定できます。

```
{
 "Version": "2012-10-17",
 "Statement": [
 {
 "Effect": "Allow",
 "Action": [
 "glue:GetConnection",
 "glue:GetDatabase",
 "glue:GetTable",
 "glue:GetPartition"
],
 "Principal": {"AWS": ["arn:aws:iam::grantee-account-id:root"]},
 "Resource": [
 "arn:aws:glue:us-east-1:grantor-account-id:*"
]
 }
]
}
```

#### Note

権限付与者のアカウントのテーブルで Amazon S3 の場所を指しており、この場所が権限付与者のアカウントにも含まれる場合、被付与者のアカウントで ETL ジョブを実行するために使用する IAM ロールには、権限付与者のアカウントからオブジェクトをリスト化し取得するための、アクセス許可が必要となります。

アカウント A のクライアントが ETL ジョブを作成して実行するアクセス許可をすでに付与されている場合、クロスアカウントアクセスのために ETL ジョブをセットアップする基本的な手順は次のとおりです。

1. クロスアカウントのデータアクセスを許可します (Amazon S3 クロスアカウントアクセスがセットアップ済みである場合は、このステップをスキップします)。
  - a. アカウント B の Amazon S3 バケットポリシーを更新して、アカウント A からのクロスアカウントアクセスを許可します。
  - b. アカウント A の IAM ポリシーを更新して、アカウント B のバケットへのアクセスを許可します。
2. Data Catalog へのクロスアカウントアクセスを許可します
  - a. アカウント B の Data Catalog にアタッチされたリソースポリシーを作成または更新して、アカウント A からのアクセスを許可します。
  - b. アカウント A の IAM ポリシーを更新して、アカウント B の Data Catalog へのアクセスを許可します。

## CloudTrail のクロスアカウントロギング

AWS Glue の抽出、変換、ロード (ETL) ジョブが、AWS Lake Formation クロスアカウント付与を介して共有された Data Catalog テーブルの基板となるデータにアクセスする際には、AWS CloudTrail ログ記録に追加的な動作が発生します。

この点を考えるために、テーブルを共有する AWS アカウントは所有者アカウントであり、テーブルの共有先アカウントは受け取りアカウントとします。受取人アカウントの ETL ジョブが、所有者アカウントのテーブルにあるデータにアクセスすると、受取人アカウントのログに追加された CloudTrail のデータアクセスイベントが所有者アカウントの CloudTrail ログにコピーされます。これは、所有者アカウントが、さまざまな受け取りアカウントによるデータアクセスを追跡できるようにするためです。デフォルトでは、CloudTrail イベントには、人間が読める形式のプリンシパル ID (プリンシパル ARN) は含まれていません。受け取りアカウントの管理者は、ログにプリンシパル ARN を含めるようにオプトインできます。

詳細については、AWS Lake Formation デベロッパーガイドの「[CloudTrail のクロスアカウントロギング](#)」を参照してください。

### その他の参照資料

- [the section called “ログ記録とモニタリング”](#)

## クロスアカウントリソースの所有権および請求

例えば、AWS アカウント (アカウント A) のユーザーが、別のアカウント (アカウント B) にデータベースなどの新しいリソースを作成すると、このリソースは、リソースの作成先のアカウントであるアカウント B に所有されます。アカウント B の管理者は、新しいリソースにアクセスする完全なアクセス許可を自動的に取得します。これには、読み取り、書き込み、および第 3 のアカウントに対するアクセス許可の付与が含まれます。アカウント A のユーザーは、アカウント B によって付与された適切なアクセス許可を持っている場合のみ、前の手順で作成したリソースにアクセスできます。

ストレージコストおよび新しいリソースに直接関連するその他のコストは、リソース所有者であるアカウント B に請求されます。リソースを作成したユーザーからのリクエストのコストは、リクエストのアカウントであるアカウント A に請求されます。

AWS Glue での請求と料金の詳細については、「[How AWS Pricing Works](#)」を参照してください。

## クロスアカウントアクセスの制約事項

AWS Glue クロスアカウントアクセスには、次の制約事項があります。

- リージョンが AWS Glue をサポートする前に、Amazon Athena または Amazon Redshift Spectrum を使用してデータベースとテーブルを作成し、リソースオーナーアカウントが Amazon Athena データカタログを AWS Glue に移行していない場合、AWS Glue へのクロスアカウントアクセスは許可されません。現在の移行ステータスは [GetCatalogImportStatus \(get\\_catalog\\_import\\_status\)](#) を使用して確認できます。Athena カタログを AWS Glue に移行する方法の詳細については、「Amazon Athena ユーザーガイド」の「[ステップバイステップの AWS Glue Data Catalog へのアップグレード](#)」を参照してください。
- クロスアカウントアクセスは、データベース、テーブル、ユーザー定義関数、接続などの Data Catalog リソースでのみサポートされています。
- Athena DataCatalog リソースからデータカタログへのクロスアカウントアクセスには、カタログを Athena として登録する必要があります。手順については、Amazon Athena ユーザーガイドの「[別のアカウントからの AWS Glue Data Catalog の登録](#)」を参照してください。

## AWS Glue のアイデンティティとアクセスのトラブルシューティング

次の情報は、AWS Glue と IAM の使用時に発生する可能性がある一般的な問題の診断と修正に役立ちます。

### トピック

- [AWS Glue でアクションを実行する権限がない](#)
- [iam を実行する権限がありません。PassRole](#)
- [自分の 以外のユーザーに AWS Glue リソース AWS アカウント へのアクセスを許可したい](#)

## AWS Glue でアクションを実行する権限がない

「I am not authorized to perform an action in Amazon Bedrock」というエラーが表示された場合、そのアクションを実行できるようにポリシーを更新する必要があります。

次のエラー例は、mateojackson IAM ユーザーがコンソールを使用して、ある *my-example-widget* リソースに関する詳細情報を表示しようとしたことを想定して、その際に必要な `glue:GetWidget` アクセス許可を持っていない場合に発生するものです。

```
User: arn:aws:iam::123456789012:user/mateojackson is not authorized to perform:
glue:GetWidget on resource: my-example-widget
```

この場合、`glue:GetWidget` アクションを使用して *my-example-widget* リソースへのアクセスを許可するように、mateojackson ユーザーのポリシーを更新する必要があります。

サポートが必要な場合は、AWS 管理者にお問い合わせください。サインイン認証情報を提供した担当者が管理者です。

## iam を実行する権限がありません。PassRole

`iam:PassRole` アクションを実行する権限がないというエラーが表示された場合は、ポリシーを更新して AWS Glue にロールを渡すことができるようにする必要があります。

一部の AWS のサービスでは、新しいサービスロールまたはサービスにリンクされたロールを作成する代わりに、そのサービスに既存のロールを渡すことができます。そのためには、サービスにロールを渡す権限が必要です。

次のエラー例は、marymajor という名前の IAM ユーザーが、コンソールから AWS Glue でアクションを実行しようとするると発生します。ただし、このアクションをサービスが実行するには、サービスロールから付与された権限が必要です。メアリーには、ロールをサービスに渡す許可がありません。

```
User: arn:aws:iam::123456789012:user/marymajor is not authorized to perform:
iam:PassRole
```

この場合、Mary のポリシーを更新してメアリーに iam:PassRole アクションの実行を許可する必要があります。

サポートが必要な場合は、AWS 管理者にお問い合わせください。サインイン認証情報を提供した担当者が管理者です。

## 自分の 以外のユーザーに AWS Glue リソース AWS アカウント へのアクセスを許可したい

他のアカウントのユーザーや組織外の人が、リソースにアクセスするために使用できるロールを作成できます。ロールの引き受けを委託するユーザーを指定できます。リソースベースのポリシーまたはアクセスコントロールリスト (ACL) をサポートするサービスの場合、それらのポリシーを使用して、リソースへのアクセスを付与できます。

詳細については、以下を参照してください:

- AWS Glue がこれらの機能をサポートしているかどうかを確認するには、「」を参照してください [AWS Glue と IAM の連携方法](#)。
- 所有 AWS アカウント している のリソースへのアクセスを提供する方法については、[IAM ユーザーガイドの「所有 AWS アカウント している別の の IAM ユーザーへのアクセスを提供する」](#)を参照してください。
- リソースへのアクセスをサードパーティー に提供する方法については AWS アカウント、IAM ユーザーガイドの [「サードパーティー AWS アカウント が所有する へのアクセスを提供する」](#)を参照してください。
- ID フェデレーションを介してアクセスを提供する方法については、IAM ユーザーガイドの [外部で認証されたユーザー \(ID フェデレーション\) へのアクセスの許可](#)を参照してください。
- クロスアカウントアクセスでのロールとリソースベースのポリシーの使用の違いについては、IAM ユーザーガイドの [「IAM でのクロスアカウントリソースアクセス」](#)を参照してください。

## AWS Glue でのログ記録とモニタリング

ETL (抽出、変換、ロード) ジョブの実行を自動化できます。AWS Glue は、モニタリングできるクローラとジョブのメトリクスを提供します。必要なメタデータを使用して AWS Glue Data Catalog を設定すると、AWS Glue は環境のヘルスに関する統計を提供します。クローラとジョブの呼び出しを、cron に基づく時間ベースのスケジュールで自動化することができます。イベントベースのトリガーが発生したときにジョブをトリガーすることもできます。

AWS Glue は、AWS Glue で ユーザー、ロール、または AWS のサービスによって実行されたアクションを記録するサービスである AWS CloudTrail と統合されています。証跡を作成する場合は、Amazon Simple Storage Service (Amazon S3) バケット、Amazon CloudWatch Logs、および Amazon CloudWatch Events への CloudTrail イベントの継続的な配信を有効にします。各イベントまたはログエントリには、リクエストの生成者に関する情報が含まれます。

Amazon CloudWatch Events を使用して、AWS のサービスを自動化し、アプリケーションの可用性の問題やリソースの変更などのシステムイベントに自動的に応答します。AWS のサービスからのイベントは、ほぼリアルタイムで CloudWatch Events に配信されます。簡単なルールを記述して、注目するイベントと、イベントがルールに一致した場合に自動的に実行するアクションを指定できます。

 以下も参照してください。

- [CloudWatch Events を使用した AWS Glue の自動化](#)
- [CloudTrail のクロスアカウントロギング](#)

クラウドのセキュリティの重要な側面に、ログ記録があります。クラウドインフラストラクチャのデバッグと保護に必要な情報をキャプチャしながら、シークレットや機密情報をキャプチャしない方法でログを設定する必要があります。ログに記録されているものを十分に理解してください。

## のコンプライアンス検証 AWS Glue

AWS のサービス が特定のコンプライアンスプログラムの範囲内にあるかどうかを確認するには、コンプライアンスプログラム [AWS のサービス による対象範囲内のコンプライアンスプログラム](#) を参照し、関心のあるコンプライアンスプログラムを選択します。一般的な情報については、[AWS 「コンプライアンスプログラム」](#) を参照してください。

を使用して、サードパーティーの監査レポートをダウンロードできます AWS Artifact。詳細については、[「でのレポートのダウンロード AWS Artifact」](#) の」を参照してください。

を使用する際のお客様のコンプライアンス責任 AWS のサービス は、お客様のデータの機密性、貴社のコンプライアンス目的、適用される法律および規制によって決まります。 は、コンプライアンスに役立つ以下のリソース AWS を提供しています。

- [セキュリティとコンプライアンスのクイックスタートガイド](#) – これらのデプロイガイドでは、アーキテクチャ上の考慮事項について説明し、セキュリティとコンプライアンスに重点を置いたベースライン環境 AWS を にデプロイする手順について説明します。
- [アマゾン ウェブ サービスにおける HIPAA セキュリティとコンプライアンスのアーキテクチャ](#) – このホワイトペーパーでは、企業が AWS を使用して HIPAA 対象アプリケーションを作成する方法について説明します。

 Note

すべて AWS のサービス HIPAA の対象となるわけではありません。詳細については、「[HIPAA 対応サービスのリファレンス](#)」を参照してください。

- [AWS コンプライアンスリソース](#) – このワークブックとガイドのコレクションは、お客様の業界や地域に適用される場合があります。
- [AWS カスタマーコンプライアンスガイド](#) – コンプライアンスの観点から責任共有モデルを理解します。このガイドでは、ガイダンスを保護し AWS のサービス、複数のフレームワーク (米国国立標準技術研究所 (NIST)、Payment Card Industry Security Standards Council (PCI)、国際標準化機構 (ISO) を含む) のセキュリティコントロールにマッピングするためのベストプラクティスをまとめています。
- 「[デベロッパーガイド](#)」の「[ルールによるリソースの評価](#)」 – この AWS Config サービスは、リソース設定が社内プラクティス、業界ガイドライン、および規制にどの程度準拠しているかを評価します。AWS Config
- [AWS Security Hub](#) – これにより AWS のサービス、内のセキュリティ状態を包括的に確認できます AWS。Security Hub では、セキュリティコントロールを使用して AWS リソースを評価し、セキュリティ業界標準とベストプラクティスに対するコンプライアンスをチェックします。サポートされているサービスとコントロールのリストについては、「[Security Hub のコントロールリファレンス](#)」を参照してください。
- [Amazon GuardDuty](#) – これにより AWS アカウント、疑わしいアクティビティや悪意のあるアクティビティがないか環境を監視することで、、、ワークロード、コンテナ、データに対する潜在的な脅威 AWS のサービス を検出します。GuardDuty は、特定のコンプライアンスフレームワークで義務付けられている侵入検知要件を満たすことで、PCI DSS などのさまざまなコンプライアンス要件への対応に役立ちます。
- [AWS Audit Manager](#) – これにより AWS のサービス、AWS 使用状況を継続的に監査し、リスクの管理方法と規制や業界標準への準拠を簡素化できます。

## のレジリエンス AWS Glue

AWS AWS グローバルインフラストラクチャはリージョンとアベイラビリティゾーンを中心に構築されています。AWS リージョンには、物理的に分離され隔離された複数のアベイラビリティゾーンがあり、低レイテンシー、高スループット、冗長性の高いネットワークで接続されています。アベイラビリティゾーンでは、アベイラビリティゾーン間で中断せずに、自動的にフェイルオーバーするアプリケーションとデータベースを設計および運用することができます。アベイラビリティゾーンは、従来の単一または複数のデータセンターインフラストラクチャよりも可用性、耐障害性、およびスケーラビリティが優れています。

[AWS リージョンとアベイラビリティゾーンの詳細については、「グローバルインフラストラクチャ」を参照してください。](#) AWS

AWS Glue ジョブの耐障害性について詳しくは、「[エラー:内の VPC 間のフェイルオーバー動作](#)」を参照してください。 AWS Glue

## AWS Glue 内のインフラストラクチャセキュリティ

マネージドサービスである AWS Glue は、[Amazon Web Services のセキュリティプロセスの概要](#) ホワイトペーパーに記載されている AWS グローバルネットワークセキュリティの手順で保護されています。

AWS 公開版 API コールを使用して、ネットワーク経由で AWS Glue にアクセスします。クライアントで Transport Layer Security (TLS) 1.0 以降がサポートされている必要があります。TLS 1.2 以降が推奨されています。また、Ephemeral Diffie-Hellman (DHE) や Elliptic Curve Ephemeral Diffie-Hellman (ECDHE) などの Perfect Forward Secrecy (PFS) を使用した暗号スイートもクライアントでサポートされている必要があります。これらのモードは、Java 7 以降など、最近のほとんどのシステムでサポートされています。

また、リクエストには、アクセスキー ID と、IAM プリンシパルに関連付けられているシークレットアクセスキーを使用して署名する必要があります。または、[AWS Security Token Service](#) AWS STS を使用して、一時的なセキュリティ認証情報を生成し、リクエストに署名することもできます。

### トピック

- [AWS Glue とインターフェース VPC エンドポイント \(AWS PrivateLink\)](#)
- [共有 Amazon VPC](#)

## AWS Glue とインターフェイス VPC エンドポイント (AWS PrivateLink)

VPC と AWS Glue とのプライベート接続を確立するには、インターフェイス VPC エンドポイントを作成します。インターフェイスエンドポイントは、インターネットゲートウェイ、NAT デバイス、VPN 接続、AWS Direct Connect 接続のいずれも必要とせずに AWS Glue API にプライベートにアクセスできるテクノロジーである [AWS PrivateLink](#) を利用しています。VPC のインスタンスは、パブリック IP アドレスがなくても AWS Glue API と通信できます。VPC と AWS Glue 間のトラフィックは、Amazon ネットワークを離れません。

各インターフェイスエンドポイントは、サブネット内の 1 つ以上の [Elastic Network Interface](#) によって表されます。

詳細については、「Amazon VPC ユーザーガイド」の「[インターフェイス VPC エンドポイント \(AWS PrivateLink\)](#)」を参照してください。

### AWS Glue VPC エンドポイントに関する考慮事項

AWS Glue のインターフェイス VPC エンドポイントを設定する前に、Amazon VPC ユーザーガイドの [インターフェイスエンドポイントのプロパティと制限](#) を確認してください。

AWS Glue は、VPC からのすべての API アクションの呼び出しをサポートしています。

### AWS Glue 用のインターフェイス VPC エンドポイントの作成

AWS Glue サービス用の VPC エンドポイントは、Amazon VPC コンソールまたは AWS Command Line Interface (AWS CLI) を使用して作成できます。詳細については、「Amazon VPC ユーザーガイド」の「[インターフェイスエンドポイントの作成](#)」を参照してください。

AWS Glue 用の VPC エンドポイントは、以下のサービス名を使用して作成します。

- `com.amazonaws.regionglue`

エンドポイントのプライベート DNS を有効にすると、リージョンのデフォルト DNS 名 (AWS Glue など) を使用して、`glue.us-east-1.amazonaws.com` への API リクエストを実行できます。

詳細については、「Amazon VPC ユーザーガイド」の「[インターフェイスエンドポイントを介したサービスへのアクセス](#)」を参照してください。

## AWS Glue 用の VPC エンドポイントポリシーの作成

VPC エンドポイントには、AWS Glue へのアクセスを制御するエンドポイントポリシーをアタッチできます。このポリシーでは、以下の情報を指定します。

- アクションを実行できるプリンシパル。
- 実行可能なアクション。
- このアクションを実行できるリソース。

詳細については、「Amazon VPC ユーザーガイド」の「[VPC エンドポイントでサービスへのアクセスを制御する](#)」を参照してください。

例: AWS Glue の VPC エンドポイントポリシージョブの作成と更新を許可する

以下は、AWS Glue 用のエンドポイントポリシーの例です。エンドポイントにアタッチされると、このポリシーは、すべてのリソースですべてのプリンシパルに、リストされている AWS Glue アクションへのアクセス権を付与します。

```
{
 "Statement": [
 {
 "Principal": "*",
 "Effect": "Allow",
 "Action": [
 "glue:CreateJob",
 "glue:UpdateJob",
 "iam:PassRole"
],
 "Resource": "*"
 }
]
}
```

例: VPC エンドポイントポリシーでデータカタログアクセスを読み取り専用で許可する場合

以下は、AWS Glue 用のエンドポイントポリシーの例です。エンドポイントにアタッチされると、このポリシーは、すべてのリソースですべてのプリンシパルに、リストされている AWS Glue アクションへのアクセス権を付与します。

```
{
```

```
"Statement": [
 {
 "Principal": "*",
 "Effect": "Allow",
 "Action": [
 "glue:GetDatabase",
 "glue:GetDatabases",
 "glue:GetTable",
 "glue:GetTables",
 "glue:GetTableVersion",
 "glue:GetTableVersions",
 "glue:GetPartition",
 "glue:GetPartitions",
 "glue:BatchGetPartition",
 "glue:SearchTables"
],
 "Resource": "*"
 }
]
```

## 共有 Amazon VPC

AWS Glue は、Amazon Virtual Private Cloud において、共有Virtual Private Cloud (VPC) をサポートしています。Amazon VPC の共有では、複数の AWS アカウントが、Amazon EC2 インスタンスや Amazon Relational Database Service (Amazon RDS) データベースなどのアプリケーションリソースを、共有の一元管理された Amazon VPC に作成できます。このモデルでは、VPC を所有するアカウント (所有者) は、同じ組織に属する他のアカウント (参加者) と 1 つまたは複数のサブネットを共有します。AWS Organizations サブネットが共有されると、参加者は共有しているサブネット内にある自分のアプリケーションリソースを表示、作成、変更、および削除できます。

AWS Glue で、共有サブネットとの接続を作成するには、アカウント内にセキュリティグループを作成し、そのセキュリティグループを共有サブネットにアタッチする必要があります。

詳細については、以下のトピックを参照してください。

- 詳細については、Amazon VPC ユーザーガイドの「[Working with Shared VPCs](#)」を参照してください。
- AWS Organizations ユーザーガイドの「[What Is AWS Organizations?](#)」

# AWS Glue のトラブルシューティング

## トピック

- [AWS Glue トラブルシューティング情報の収集](#)
- [AWS Glue for Spark のエラーのトラブルシューティング](#)
- [ログから AWS Glue for Ray エラーをトラブルシューティングする](#)
- [AWS Glue 機械学習の例外](#)
- [AWS Glue のクォータ](#)

## AWS Glue トラブルシューティング情報の収集

AWS Glue でエラーまたは予期しない動作が発生して AWS Support に問い合わせる必要がある場合、名前、ID、および失敗したアクションに関連したログに関する情報をまず収集する必要があります。この情報が用意してあると、AWS Support は発生している問題を解決しやすくなります。

アカウント ID に加えて、次の失敗のタイプそれぞれについて以下の情報を収集します。

クローラが失敗した場合、以下の情報を収集します。

- クローラの名前

クローラ実行からのログは `/aws-glue/crawlers` の下の CloudWatch Logs にあります。

テスト接続が失敗した場合、以下の情報を収集します。

- 接続名
- 接続 ID
- `jdbc:protocol://host:port/database-name` フォームの JDBC 接続文字列。

テスト接続からのログは `/aws-glue/testconnection` の下の CloudWatch Logs にあります。

ジョブが失敗した場合、以下の情報を収集します。

- ジョブ名
- `jr_xxxxx` フォームのジョブ実行 ID。

ジョブ実行からのログは `/aws-glue/jobs` の下の CloudWatch Logs にあります。

# AWS Glue for Spark のエラーのトラブルシューティング

でエラーが発生した場合は AWS Glue、次の解決策を使用して問題の原因を見つけて修正してください。

## Note

AWS Glue GitHub リポジトリには、[AWS Glue よくある質問](#) のトラブルシューティングに関する追加のガイダンスが含まれています。

## トピック

- [エラー: リソースを利用できません。](#)
- [エラー: VPC の subnetId に S3 エンドポイントまたは NAT ゲートウェイが見つかりませんでした](#)
- [エラー: 必要なセキュリティグループのインバウンドルール](#)
- [エラー: 必要なセキュリティグループのアウトバウンドルール](#)
- [エラー: 渡されたロールに AWS Glue サービスのロールを引き受けるアクセス許可を付与する必要があるため、ジョブの実行に失敗しました](#)
- [エラー: DescribeVpcEndpoints アクションが不正です。VPC ID vpc-id を検証できません](#)
- [エラー: DescribeRouteTables アクションが不正です。サブネット ID: VPC id: vpc-id のサブネット ID を検証できません](#)
- [エラー: ec2 の呼び出しに失敗しました: DescribeSubnets](#)
- [エラー: ec2 の呼び出しに失敗しました。DescribeSecurityGroups](#)
- [エラー: AZ のサブネットが見つかりませんでした](#)
- [エラー: JDBC ターゲットへの書き込み時のジョブ実行の例外](#)
- [Error: Amazon S3: The operation is not valid for the object's storage class](#)
- [エラー: Amazon S3 タイムアウト](#)
- [エラー: Amazon S3 へのアクセスが拒否されました](#)
- [エラー: Amazon S3 のアクセスキー ID が存在しません](#)
- [エラー: URI に s3a:// を使用しながら Amazon S3 にアクセスするとジョブ実行が失敗します](#)
- [エラー: Amazon S3 のサービストークンが有効期限切れです](#)
- [エラー: ネットワークインターフェイスのプライベート DNS が見つかりません](#)
- [エラー: 開発エンドポイントのプロビジョニングに失敗しました](#)

- [エラー: ノートブックサーバー CREATE\\_FAILED](#)
- [エラー: ローカルノートブックの起動に失敗する](#)
- [エラー: クローラの実行に失敗しました](#)
- [エラー: パーティションが更新されませんでした](#)
- [Error: Job bookmark update failed due to version mismatch \(エラー: バージョンの不一致のため、ジョブのブックマークの更新に失敗しました\)](#)
- [エラー: ジョブのブックマークが有効なときにジョブがデータを再処理しています](#)
- [エラー: の VPCs間のフェイルオーバー動作 AWS Glue](#)
- [クローラーが Lake Formation の認証情報を使用している場合のクローラーエラーのトラブルシューティング](#)

## エラー: リソースを利用できません。

がリソース使用不可メッセージを AWS Glue 返す場合は、エラーメッセージまたはログを表示して、問題の詳細を知ることができます。ここでは、トラブルシューティングするための一般的な方法について説明します。

- 使用する接続および開発エンドポイントについては、クラスターに Elastic Network Interface が不足していないかどうかを確認してください。

## エラー: VPC の subnetId に S3 エンドポイントまたは NAT ゲートウェイが見つかりませんでした

問題の診断に役立つメッセージのサブネット ID と VPC ID を確認します。

- AWS Glueに必要な Amazon S3 VPC エンドポイントがセットアップ済みであることを確認してください。さらに、設定の一部である場合には、NAT ゲートウェイを確認します。詳細については、[Amazon S3 用の VPC エンドポイント](#) を参照してください。

## エラー: 必要なセキュリティグループのインバウンドルール

少なくとも1つのセキュリティグループのすべての受信ポートを開く必要があります。トラフィックを制限するために、インバウンドルールのソースセキュリティグループを同じセキュリティグループに制限できます。

- 使用するすべての接続について、セキュリティグループで自己参照のインバウンドルールを確認します。詳細については、[データストアへのネットワークアクセスを設定する](#) を参照してください。
- 開発エンドポイントを使用している場合は、自己参照のインバウンドルールのセキュリティグループを確認します。詳細については、[データストアへのネットワークアクセスを設定する](#) を参照してください。

## エラー: 必要なセキュリティグループのアウトバウンドルール

少なくとも1つのセキュリティグループのすべての発信ポートを開く必要があります。トラフィックを制限するために、アウトバウンドルールのソースセキュリティグループを同じセキュリティグループに制限できます。

- 使用するすべての接続について、セキュリティグループで自己参照のアウトバウンドルールを確認します。詳細については、[データストアへのネットワークアクセスを設定する](#) を参照してください。
- 開発エンドポイントを使用している場合は、自己参照のアウトバウンドルールのセキュリティグループを確認します。詳細については、「[データストアへのネットワークアクセスを設定する](#)」を参照してください。

## エラー: 渡されたロールに AWS Glue サービスのロールを引き受けるアクセス許可を付与する必要があるため、ジョブの実行に失敗しました

ジョブを定義するユーザーは AWS Glue の `iam:PassRole` のアクセス権限を持っている必要があります。

- ユーザーが AWS Glue ジョブを作成するときに、ユーザーのロールに `iam:PassRole` のを含むポリシーが含まれていることを確認します AWS Glue。詳細については、「[ステップ 3: AWS Glue にアクセスするユーザーまたはグループにポリシーをアタッチする](#)」を参照してください。

## エラー: DescribeVpcEndpoints アクションが不正です。VPC ID vpc-id を検証できません

- に渡されたポリシーで `ec2:DescribeVpcEndpoints` 許可 AWS Glue を確認します。

エラー : DescribeRouteTables アクションが不正です。サブネット ID: VPC id: vpc-id のサブネット ID を検証できません

- に渡されたポリシーで アクセス ec2:DescribeRouteTables 許可 AWS Glue を確認します。

エラー: ec2 の呼び出しに失敗しました : DescribeSubnets

- に渡されたポリシーで アクセス ec2:DescribeSubnets 許可 AWS Glue を確認します。

エラー: ec2 の呼び出しに失敗しました。 DescribeSecurityGroups

- に渡されたポリシーで アクセス ec2:DescribeSecurityGroups 許可 AWS Glue を確認します。

エラー: AZ のサブネットが見つかりませんでした

- アベイラビリティゾーンは で使用できない場合があります AWS Glue。メッセージで指定されているものとは異なるアベイラビリティゾーンに新しいサブネットを作成して使用します。

エラー: JDBC ターゲットへの書き込み時のジョブ実行の例外

JDBC ターゲットに書き込むジョブを実行すると、以下のシナリオのようにジョブでエラーが発生する可能性があります。

- ジョブが Microsoft SQL Server テーブルに書き込む場合、テーブルに Boolean 型として定義された列がある場合には、SQL Server データベースでテーブルを事前に定義する必要があります。SQL Server ターゲット を使用して AWS Glue コンソールでジョブを定義する場合、データターゲット にテーブルを作成する オプションを使用して、ソース列をデータ型 のターゲット列にマッピングしないでください Boolean。ジョブの実行時にエラーが発生する可能性があります。

次のようにして、エラーを回避できます。

- [Boolean (ブール)] 列を使用して既存のテーブルを選択します。
- ApplyMapping 変換を編集し、ソース内の [Boolean] (ブール) 列をターゲット内の数値または文字列にマップします。
- ApplyMapping 変換を編集して、[Boolean] (ブール) 列をソースから削除します。

- ジョブが Oracle テーブルに書き込む場合は、Oracle オブジェクトの名前の長さを調整する必要があります。Oracle の一部のバージョンでは、識別子の最大長は 128 バイトまたは 30 バイトに制限されています。この制限は、Oracle ターゲットデータストアのテーブル名および列名に影響を与えます。

次のようにして、エラーを回避できます。

- ご使用のバージョンの制限内で Oracle ターゲットテーブルに名前を付けます。
- デフォルトの列名は、データのフィールド名から生成されます。列名が制限よりも長い場合は処理のために、ApplyMapping または RenameField 変換を使用して列の名前を制限内に変更します。

## Error: Amazon S3: The operation is not valid for the object's storage class

がこのエラーを AWS Glue 返す場合、AWS Glue ジョブは Amazon S3 ストレージクラス階層にまたがるパーティションを持つテーブルからデータを読み取っている可能性があります。

- ストレージクラスの除外を使用すると、ジョブ AWS Glue がこれらのストレージクラス階層にまたがるパーティションを持つテーブルで動作するようにできます。除外しないと、これらの階層からデータを読み取るジョブは AmazonS3Exception: The operation is not valid for the object's storage class のエラーで失敗します。

詳細については、「[Amazon S3 ストレージクラスの除外](#)」を参照してください。

## エラー:Amazon S3 タイムアウト

が接続タイムアウトエラーを AWS Glue 返した場合、別の AWS リージョンの Amazon S3 バケットにアクセスしようとしている可能性があります。

- Amazon S3 VPC エンドポイントは、AWS リージョン内のバケットにのみトラフィックをルーティングできます。他のリージョンのバケットに接続する必要がある場合、考えられる回避策は NAT ゲートウェイを使用することです。詳細については、「[NAT ゲートウェイ](#)」を参照してください。

## エラー: Amazon S3 へのアクセスが拒否されました

が Amazon S3 バケットまたはオブジェクトへのアクセス拒否エラーを AWS Glue 返す場合、提供された IAM ロールにデータストアへのアクセス許可を持つポリシーがないためである可能性があります。

- ETL ジョブは、ソースまたはターゲットとして使用される Amazon S3 データストアにアクセスできる必要があります。クローラは、クロールする Amazon S3 データストアにアクセスできる必要があります。詳細については、[ステップ 2: AWS Glue 用の IAM ロールを作成する](#) を参照してください。

## エラー: Amazon S3 のアクセスキー ID が存在しません

がジョブの実行時にアクセスキー ID が存在しないというエラーを AWS Glue 返す場合、次のいずれかの理由が原因である可能性があります。

- ETL ジョブは IAM ロールを使用してデータストアにアクセスし、IAM ロールがそのジョブの開始前に削除されていないことを確認します。
- IAM ロールには、データストアにアクセスするためのアクセス許可が含まれており、アタッチ済みの `s3:ListBucket` を含む Amazon S3 ポリシーが正しいことを確認します。

## エラー: URI に `s3a://` を使用しながら Amazon S3 にアクセスするとジョブ実行が失敗します

ジョブ実行から「XML ドキュメントをハンドラクラスで解析できませんでした」などのエラーが返された場合は、`s3a://` URI を使用して数百のファイルをリストアップしようとして失敗した可能性があります。代わりに `s3://` URI を使用してデータストアにアクセスしてください。次の例外トレースを調べてエラーを確認できます。

```
1. com.amazonaws.SdkClientException: Failed to parse XML document with handler class
 com.amazonaws.services.s3.model.transform.XmlResponsesSaxParser$ListBucketHandler
2. at
 com.amazonaws.services.s3.model.transform.XmlResponsesSaxParser.parseXmlInputStream(XmlResponses
3. at
 com.amazonaws.services.s3.model.transform.XmlResponsesSaxParser.parseListBucketObjectsResponse
4. at com.amazonaws.services.s3.model.transform.Unmarshallers
 $ListObjectsUnmarshaller.unmarshall(Unmarshallers.java:70)
```

```
5. at com.amazonaws.services.s3.model.transform.Unmarshallers
 $ListObjectsUnmarshaller.unmarshall(Unmarshallers.java:59)
6. at
 com.amazonaws.services.s3.internal.S3XmlResponseHandler.handle(S3XmlResponseHandler.java:62)
7. at
 com.amazonaws.services.s3.internal.S3XmlResponseHandler.handle(S3XmlResponseHandler.java:31)
8. at
 com.amazonaws.http.response.AwsResponseHandlerAdapter.handle(AwsResponseHandlerAdapter.java:70)
9. at com.amazonaws.http.AmazonHttpClient
 $RequestExecutor.handleResponse(AmazonHttpClient.java:1554)
10. at com.amazonaws.http.AmazonHttpClient
 $RequestExecutor.executeOneRequest(AmazonHttpClient.java:1272)
11. at com.amazonaws.http.AmazonHttpClient
 $RequestExecutor.executeHelper(AmazonHttpClient.java:1056)
12. at com.amazonaws.http.AmazonHttpClient
 $RequestExecutor.doExecute(AmazonHttpClient.java:743)
13. at com.amazonaws.http.AmazonHttpClient
 $RequestExecutor.executeWithTimer(AmazonHttpClient.java:717)
14. at com.amazonaws.http.AmazonHttpClient
 $RequestExecutor.execute(AmazonHttpClient.java:699)
15. at com.amazonaws.http.AmazonHttpClient$RequestExecutor.access
 $500(AmazonHttpClient.java:667)
16. at com.amazonaws.http.AmazonHttpClient
 $RequestExecutionBuilderImpl.execute(AmazonHttpClient.java:649)
17. at com.amazonaws.http.AmazonHttpClient.execute(AmazonHttpClient.java:513)
18. at com.amazonaws.services.s3.AmazonS3Client.invoke(AmazonS3Client.java:4325)
19. at com.amazonaws.services.s3.AmazonS3Client.invoke(AmazonS3Client.java:4272)
20. at com.amazonaws.services.s3.AmazonS3Client.invoke(AmazonS3Client.java:4266)
21. at com.amazonaws.services.s3.AmazonS3Client.listObjects(AmazonS3Client.java:834)
22. at org.apache.hadoop.fs.s3a.S3AFileSystem.getFileStatus(S3AFileSystem.java:971)
23. at
 org.apache.hadoop.fs.s3a.S3AFileSystem.deleteUnnecessaryFakeDirectories(S3AFileSystem.java:115)
24. at org.apache.hadoop.fs.s3a.S3AFileSystem.finishedWrite(S3AFileSystem.java:1144)
25. at org.apache.hadoop.fs.s3a.S3AOutputStream.close(S3AOutputStream.java:142)
26. at org.apache.hadoop.fs.FSDataOutputStream
 $PositionCache.close(FSDataOutputStream.java:74)
27. at org.apache.hadoop.fs.FSDataOutputStream.close(FSDataOutputStream.java:108)
28. at org.apache.parquet.hadoop.ParquetFileWriter.end(ParquetFileWriter.java:467)
29. at
 org.apache.parquet.hadoop.InternalParquetRecordWriter.close(InternalParquetRecordWriter.java:1)
30. at
 org.apache.parquet.hadoop.ParquetRecordWriter.close(ParquetRecordWriter.java:112)
31. at
 org.apache.spark.sql.execution.datasources.parquet.ParquetOutputWriter.close(ParquetOutputWriter.java:1)
```

```
32. at org.apache.spark.sql.execution.datasources.FileFormatWriter
$SingleDirectoryWriteTask.releaseResources(FileFormatWriter.scala:252)
33. at org.apache.spark.sql.execution.datasources.FileFormatWriter$$anonfun
orgapache$spark$sql$execution$databases$FileFormatWriter$$executeTask
$3.apply(FileFormatWriter.scala:191)
34. at org.apache.spark.sql.execution.datasources.FileFormatWriter$$anonfun
orgapache$spark$sql$execution$databases$FileFormatWriter$$executeTask
$3.apply(FileFormatWriter.scala:188)
35. at org.apache.spark.util.Utils
$.tryWithSafeFinallyAndFailureCallbacks(Utils.scala:1341)
36. at org.apache.spark.sql.execution.datasources.FileFormatWriter$.org$apache$spark
sqlexecution$databases$FileFormatWriter$$executeTask(FileFormatWriter.scala:193)
37. at org.apache.spark.sql.execution.datasources.FileFormatWriter$$anonfun$write$1$
$anonfun$3.apply(FileFormatWriter.scala:129)
38. at org.apache.spark.sql.execution.datasources.FileFormatWriter$$anonfun$write$1$
$anonfun$3.apply(FileFormatWriter.scala:128)
39. at org.apache.spark.scheduler.ResultTask.runTask(ResultTask.scala:87)
40. at org.apache.spark.scheduler.Task.run(Task.scala:99)
41. at org.apache.spark.executor.Executor$TaskRunner.run(Executor.scala:282)
42. at java.util.concurrent.ThreadPoolExecutor.runWorker(ThreadPoolExecutor.java:1149)
43. at java.util.concurrent.ThreadPoolExecutor$Worker.run(ThreadPoolExecutor.java:624)
44. at java.lang.Thread.run(Thread.java:748)
```

## エラー: Amazon S3 のサービストークンが有効期限切れです

Amazon Redshift との間でデータを移行する際には、1 時間で有効期限が切れる、Amazon S3 の一時的な認証情報が使用されます。長時間実行するジョブがある場合は、失敗する可能性があります。長時間実行するジョブを設定して Amazon Redshift との間でデータを移動する方法については、「[aws-glue-programming-etl-connect-redshift-home](#)」を参照してください。

## エラー: ネットワークインターフェイスのプライベート DNS が見つかりません

ジョブが失敗したり、開発エンドポイントがプロビジョニングに失敗した場合は、ネットワーク設定が問題の原因である可能性があります。

- Amazon が提供する DNS を使用している場合は、enableDnsHostnames の値を true に設定する必要があります。詳細については、「[DNS](#)」を参照してください。

## エラー: 開発エンドポイントのプロビジョニングに失敗しました

が開発エンドポイントのプロビジョニングに成功 AWS Glue しなかった場合は、ネットワーク設定の問題が原因である可能性があります。

- 開発エンドポイントを定義すると、VPC、サブネット、およびセキュリティグループが検証され、特定の要件を満たしていることが確認されます。
- オプションの SSH パブリックキーを指定した場合、有効な SSH パブリックキーになっていることを確認します。
- VPC コンソールで、VPC が有効な [DHCP option set (DHCP オプションセット)] を使用していることを確認します。詳細については、「[DHCP オプションセット](#)」を参照してください。
- クラスターが PROVISIONING 状態のままになっている場合は、AWS Supportにお問い合わせください。

## エラー: ノートブックサーバー CREATE\_FAILED

が開発エンドポイントのノートブックサーバーの作成に AWS Glue 失敗した場合、次のいずれかの問題が原因である可能性があります。

- AWS Glue は、ノートブックサーバーの設定時に IAM ロールを Amazon EC2 に渡します。IAM ロールでは、Amazon EC2 との信頼関係が必要です。
- IAM ロールには、同じ名前のインスタンスプロファイルが必要です。IAM コンソールで Amazon EC2 のロールを作成すると、同じ名前のインスタンスプロファイルが自動的に作成されます。無効なインスタンスプロファイル名 `iamInstanceProfile.name` に関するログのエラーをチェックします。詳細については、「[インスタンスプロファイルの使用](#)」を参照してください。
- ロールにノートブックサーバーを作成するために渡したポリシーの `aws-glue*` バケットにアクセスする権限があることを確認します。

## エラー: ローカルノートブックの起動に失敗する

ローカルノートブックの起動に失敗し、ディレクトリまたはフォルダが見つからないというエラーが報告された場合は、次のいずれかの問題が原因である可能性があります。

- Microsoft Windows で実行している場合は、`JAVA_HOME` 環境変数が適切な Java ディレクトリを指していることを確認します。この変数を更新せずに Java が更新される可能性があります。存在しなくなったフォルダを変数が指していると、Jupyter Notebook は起動に失敗します。

## エラー: クローラの実行に失敗しました

がクローラを正常に実行してデータをカタログ化 AWS Glue できない場合、次のいずれかの理由が考えられます。まず、AWS Glue コンソールのクローラリストにエラーがあるかどうかを確認します。クローラ名の横に感嘆符アイコンがあるかどうかを確認し、アイコンの上にカーソルを置いて関連するメッセージを確認します。

- のログで実行されているクローラの CloudWatch ログを確認します /aws-glue/crawlers。

## エラー: パーティションが更新されませんでした

ETL ジョブの実行時に Data Catalog でパーティションが更新されなかった場合は、ログの DataSink クラスからのログステートメントが役立つ CloudWatch ことがあります。

- 「Attempting to fast-forward updates to the Catalog - nameSpace:」 — このジョブによって変更しようとしているデータベース、テーブル、および catalogId を示します。このステートメントがここにはない場合は、enableUpdateCatalog が true に設定され、getSink() パラメータとして、または additional\_options で正しく渡されるかどうかを確認します。
- 「Schema change policy behavior:」 — 渡したスキーマ updateBehavior 値を示します。
- 「Schemas qualify (schema compare):」 — true または false になります。
- 「Schemas qualify (case-insensitive compare):」 — true または false になります。
- 両方が false で、updateBehavior が に設定されていない場合、DynamicFrame スキーマは同一であるか UPDATE\_IN\_DATABASE、Data Catalog テーブルスキーマに表示される列のサブセットが含まれている必要があります。

パーティションの更新の詳細については、「[AWS Glue ETL ジョブを使用してデータカタログのスキーマを更新し、新規パーティションを追加する](#)」を参照してください。

## Error: Job bookmark update failed due to version mismatch (エラー: パーティションの不一致のため、ジョブのブックマークの更新に失敗しました)

Amazon S3 の異なるデータセットに同じ変換/ロジックを適用するために、AWS Glue ジョブをパラメータ化しようとしている可能性があります。ユーザーは指定の場所で処理されたファイルを追跡し

たいと考えています。同じソースバケットで同じジョブを実行し、同じ/異なる送信先に同時に書き込んだ場合 (同時実行数 >1)、次のエラーによりジョブは失敗します。

```
py4j.protocol.Py4JJavaError: An error occurred while
callingz:com.amazonaws.services.glue.util.Job.commit.:com.amazonaws.services.gluejobexecutor.m
Continuation update failed due to version mismatch. Expected version 2 but found
version 3
```

解決策として、同時実行数を 1 に設定するか、ジョブを同時に実行しないようにします。

現在、AWS Glue ブックマークは同時ジョブ実行をサポートしておらず、コミットは失敗します。

## エラー：ジョブのブックマークが有効なときにジョブがデータを再処理しています

AWS Glue ジョブのブックマークを有効にしても、ETL ジョブが以前の実行で既に処理されたデータを再処理している場合があります。このエラーの一般的な原因を確認します。

### 最大同時実行数

ジョブの同時実行の最大数をデフォルト値の 1 より大きく設定すると、ジョブのブックマークを妨げる可能性があります。ジョブのブックマークがオブジェクトの最終更新日時を確認し、どのオブジェクトを再処理する必要があるか確認するときに発生します。詳細については、「[での Spark ジョブのジョブプロパティの設定 AWS Glue](#)」の最大同時実行数に関する説明を参照してください。

### Job オブジェクトが見つからない

ジョブの実行スクリプトが次のコメントで終わることを確認します。

```
job.commit()
```

このオブジェクトを含めると、 はジョブ実行のタイムスタンプとパス AWS Glue を記録します。同じパスでジョブを再度実行すると、 は新しいファイルのみ AWS Glue を処理します。このオブジェクトを含まず、ジョブのブックマークが有効になっている場合、ジョブはすでに処理されたファイルと共に新しいファイルを再処理して、ジョブのターゲットデータストアで冗長化されます。

### 変換コンテキストパラメータが見つからない

変換コンテキストは、GlueContext クラスのオプションのパラメータですが、ジョブのブックマークは含めなければ機能しません。このエラーを解決するには、次に示すように、[の作成 DynamicFrame](#)時に変換コンテキストパラメータを追加します。

```
sample_dynF=create_dynamic_frame_from_catalog(database,
table_name,transformation_ctx="sample_dynF")
```

## 入力ソース

入力ソースのリレーショナルデータベース (JDBC 接続) を使用している場合は、テーブルのプライマリキーが順番になっている場合のみジョブのブックマークが機能します。ジョブのブックマークは新しい行に対して機能しますが、更新された行に対しては機能しません。これは、ジョブのブックマークで探すプライマリキーがすでに存在しているからです。これは、入力ソースが Amazon Simple Storage Service (Amazon S3) である場合には適用されません。

## 最終更新日時

Amazon S3 入力ソースの場合、ジョブのブックマークは、ファイル名ではなくオブジェクトの最終更新日時を確認して、どのオブジェクトを再処理する必要があるのかを確認します。入力ソースデータが最後のジョブ実行以降に変更されている場合、ジョブを再度実行すると、ファイルが再度処理されます。

## エラー: の VPCs間のフェイルオーバー動作 AWS Glue

AWS Glue 4.0 以前のバージョンのジョブのフェイルオーバーには、次のプロセスが使用されます。

概要: AWS Glue 接続は、ジョブ実行の送信時に選択されます。ジョブ実行で何らかの問題 (IP アドレスの不足、ソースへの接続、ルーティングの問題) が発生した場合、ジョブ実行は失敗します。再試行が設定されている場合、は同じ接続で再試行 AWS Glue します。

1. 実行の試行ごとに、は、使用できる接続が見つかるまで、ジョブ設定にリストされている順序で接続の状態 AWS Glue をチェックします。アベイラビリティーゾーン (AZ) に障害が発生した場合、その AZ からの接続はチェックに失敗し、スキップされます。
2. AWS Glue は、以下を使用して接続を検証します。
  - 有効な Amazon VPC ID とサブネットを確認する。
  - NAT ゲートウェイまたは Amazon VPC エンドポイントが存在することを確認する。
  - サブネットに 0 個を超える IP アドレスが割り当てられていることを確認する。
  - AZ が正常であることを確認する。

AWS Glue は、ジョブ実行の送信時に接続を検証できません。

3. Amazon VPC を使用するジョブの場合、すべてのドライバーとエグゼキューターは、ジョブ実行の送信時に選択された接続を使用して同じ AZ に作成されます。

- 再試行が設定されている場合、は同じ接続で再試行 AWS Glue します。これは、この接続に関する問題が長時間継続することを保証できないためです。AZ に障害が発生すると、その AZ 内の既存のジョブ実行 (ジョブ実行の段階によって異なる) が失敗する可能性があります。再試行すると AZ の障害が検出され、新しい実行で別の AZ が選択されるはずですが。

## クローラーが Lake Formation の認証情報を使用している場合のクローラーエラーのトラブルシューティング

以下の情報を使用して、Lake Formation の認証情報を使用してクローラーを設定する際のさまざまな問題を診断し、修正します。

エラー: 「The S3 location: s3://examplepath is not registered」 (S3 ロケーション: s3://examplepath は登録されていません)

Lake Formation の認証情報を使用してクローラーを実行する場合は、最初に Lake Formation のアクセス許可を設定する必要があります。このエラーを解決するには、ターゲットの Amazon S3 ロケーションを Lake Formation で登録してください。詳細については、「[Amazon S3 ロケーションの登録](#)」を参照してください。

エラー: 「User/Role is not authorized to perform: lakeformation:GetDataAccess on resource」 (ユーザー/ロールは、リソースに対する lakeformation:GetDataAccess の実行が許可されていません)

IAM コンソールまたは AWS CLI を使用して、クローラーロールに lakeformation:GetDataAccess へのアクセス許可を追加してください。この許可があると、Lake Formation がデータにアクセスするための一時的な認証情報のリクエストを承諾します。以下のポリシーを参照してください。

```
{
 "Version": "2012-10-17",
 "Statement": {
 "Effect": "Allow",
 "Action": [
 "lakeformation:GetDataAccess"
],
 "Resource": "*"
 }
}
```

エラー: 「Insufficient Lake Formation permission(s) on (Database name: exampleDatabase, Table Name: exampleTable)」 (データベース名: exampleDatabase、テーブル名: exampleTable に対する Lake Formation へのアクセス許可が不十分です)

Lake Formation コンソール (<https://console.aws.amazon.com/lakeformation/>) で、出力データベースとして指定されたデータベースに対するクローラーロールへのアクセス許可 (Create、Describe、Alter) を付与します。テーブルに対するアクセス許可を付与することもできます。詳細については、「[名前付きリソース方式を使用した Data Catalog 許可の付与](#)」を参照してください。

エラー: 「Insufficient Lake Formation permission(s) on s3://examplepaths」 (s3://examplepath に対する Lake Formation へのアクセス許可が不十分です)

#### 1. クロスアカウントでのクローリング

- a. Amazon S3 バケットが登録されているアカウント (アカウント B) を使用して、Lake Formation コンソール (<https://console.aws.amazon.com/lakeformation/>) にログインします。クローラーを実行するアカウント (アカウント A) に、データロケーションへのアクセス許可を付与します。これにより、クローラーはターゲットの Amazon S3 ロケーションからデータを読み取ることができます。
- b. クローラーが作成されたアカウント (アカウント A) で、クローラーの実行に使用される IAM ロールに、ターゲットの Amazon S3 ロケーションのデータロケーションへのアクセス許可を付与し、クローラーが Lake Formation の送信先からデータを読み取れるようにします。詳細については、「[データロケーション許可の付与 \(外部アカウント\)](#)」を参照してください。

2. アカウント内 (クローラーと登録されている Amazon S3 ロケーションが同じアカウントにある) クローリング - Amazon S3 ロケーションで実行されるクローラーに使用される IAM ロールにデータロケーションへのアクセス許可を付与し、クローラーが Lake Formation のターゲットからデータを読み取れるようにします。詳細については、「[データロケーション許可の付与 \(同じアカウント\)](#)」を参照してください。

### Lake Formation の認証情報を使用したクローラーの設定に関するよくある質問

1. AWS コンソールを使用して、Lake Formation の認証情報を使用して実行できるようにクローラーを設定するにはどうすればよいですか。

AWS Glue コンソール (<https://console.aws.amazon.com/glue/>) で、クローラーを設定する際に、オプション [Use Lake Formation credentials for crawling Amazon S3 data source] (Amazon S3 データソースのクローリングに Lake Formation の認証情報を使用する) を選択してください。クロスアカウントクローリングの場合は、ターゲットの Amazon S3 ロケーションが Lake Formation で登録されている AWS アカウント ID を指定します。アカウント内クローリングでは、accountId フィールドはオプションです。

2. AWS CLI を使用して、Lake Formation の認証情報を使用できるようにクローラーを設定するにはどうすればよいですか。

CreateCrawler API コール中に LakeFormationConfiguration を追加します。

```
"LakeFormationConfiguration": {
 "UseLakeFormationCredentials": true,
 "AccountId": "111111111111" (AWS account ID where the target Amazon S3 location
 is registered with Lake Formation)
}
```

3. Lake Formation の認証情報を使用するクローラーでサポートされているターゲットは何ですか。

Lake Formation の認証情報を使用するクローラーは、Amazon S3 (アカウント内およびクロスアカウントのクローリング) と、アカウント内のデータカタログターゲット (基盤となるロケーションが Amazon S3 である)、および Apache Iceberg ターゲットでのみサポートされています。

4. Lake Formation の認証情報を使用して、複数の Amazon S3 バケットを単一のクローラーの一部としてクローリングできますか。

いいえ。Lake Formation の認証情報供給を使用するターゲットをクローリングするには、基盤となる Amazon S3 ロケーションが同じバケットに属している必要があります。例えば、複数のお客様が同じバケット (bucket1) に属している場合は、複数のターゲットロケーション (s3://bucket1/folder1, s3://bucket1/folder2) を使用できます。別のバケット (s3://bucket1/folder1, s3://bucket2/folder2) の指定はサポートされていません。

## ログから AWS Glue for Ray エラーをトラブルシューティングする

AWS Glue では、ジョブの実行中に Ray プロセスが生成するログにアクセス可能です。Ray ジョブでエラーや予期しない動作が発生した場合は、まずログから情報を収集して失敗の原因を特定します。また、インタラクティブセッションについても同様のログが提供されています。セッションのログは、先頭に /aws-glue/ray/sessions が付きます。

ジョブの実行時、ログの行はリアルタイムで CloudWatch に送信されます。実行が完了すると、Print ステートメントが CloudWatch ログに追加されます。ログはジョブの実行後 2 週間保持されます。

## Ray ジョブのログを検査する

ジョブが失敗した場合は、ジョブ名とジョブ実行 ID を収集します。これらは AWS Glue コンソールで確認できます。ジョブページに移動し、[Runs] (実行) タブに移動します。Ray ジョブのログは、次の専用の CloudWatch ロググループに保存されます。

- `/aws-glue/ray/jobs/script-log/` — メインの Ray スクリプトが生成するログを保存します。
- `/aws-glue/ray/jobs/ray-monitor-log/` — Ray のオートスケーラープロセスが生成するログを保存します。これらのログはヘッドノードについて生成され、他のワーカーノードでは生成されません。
- `/aws-glue/ray/jobs/ray-gcs-logs/` — GCS (グローバルコントロールストア) プロセスから出力されたログを保存します。これらのログはヘッドノードについて生成され、他のワーカーノードでは生成されません。
- `/aws-glue/ray/jobs/ray-process-logs/` — ヘッドノードで実行されている他の Ray プロセス (主にダッシュボードエージェント) から出力されたログを保存します。これらのログはヘッドノードについて生成され、他のワーカーノードでは生成されません。
- `/aws-glue/ray/jobs/ray-raylet-logs/` — 各 raylet プロセスから出力されたログを保存します。これらのログは、ワーカーノード (ヘッドノードを含む) ごとに 1 つのストリームにまとめられます。
- `/aws-glue/ray/jobs/ray-worker-out-logs/` — クラスター内の各ワーカーの stdout ログを保存します。これらのログは、ワーカーノード (ヘッドノードを含む) ごとに生成されます。
- `/aws-glue/ray/jobs/ray-worker-err-logs/` — クラスター内の各ワーカーの stderr ログを保存します。これらのログは、ワーカーノード (ヘッドノードを含む) ごとに生成されます。
- `/aws-glue/ray/jobs/ray-runtime-env-log/` — Ray のセットアッププロセスに関するログを保存します。これらのログは、ワーカーノード (ヘッドノードを含む) ごとに生成されます。

## Ray ジョブのエラーをトラブルシューティングする

Ray のロググループの構成を理解し、エラーのトラブルシューティングの助けになるロググループを見つけるには、Ray のアーキテクチャに関する背景情報があると役立ちます。

AWS Glue ETL では、ワーカーはインスタンスに対応します。AWS Glue ジョブのワーカーを設定する際は、そのジョブ専用のインスタンスタイプおよびインスタンス数を設定します。Ray では、ワーカーという単語がさまざまな形で使われています。

Ray は、ヘッドノードとワーカーノードを使用して、Ray クラスター内のインスタンスの役割を区別します。Ray のワーカーノードは、分散計算の結果を得るため計算を実行する、複数のアクタープロセスをホストできます。関数のレプリカを実行するアクターはレプリカと呼ばれます。レプリカのアクターはワーカープロセスとも呼ばれます。レプリカはヘッドノードでも実行できます。ヘッドノードは、追加のプロセスを実行してクラスターを調整することから、ヘッドと呼ばれます。

計算に関与する各アクターは、独自のログストリームを生成します。これにより、いくつかのインサイトが得られます。

- ログを生成するプロセスの数は、ジョブに割り当てられるワーカーの数よりも多い場合があります。多くの場合、各インスタンスのコアには 1 つのアクターがあります。
- Ray のヘッドノードは、クラスターの管理および起動ログを生成します。それに対して Ray のワーカーノードは、そのノードで実行される処理のログのみを生成します。

Ray のアーキテクチャの詳細については、Ray のドキュメントの「[Architecture Whitepapers](#)」(アーキテクチャホワイトペーパー) を参照してください。

## 問題の領域: Amazon S3 へのアクセス

ジョブ実行のエラーメッセージを確認します。それでも十分な情報が得られない場合は、`/aws-glue/ray/jobs/script-log/` を確認します。

## 問題領域: PIP の依存性の管理

`/aws-glue/ray/jobs/ray-runtime-env-log/` をチェックします。

## 問題の領域: メインプロセスにおける中間値の検査

メインスクリプトから `stderr` または `stdout` に書き込みを行います。また、`/aws-glue/ray/jobs/script-log/` からログを取得します。

## 問題の領域: 子プロセスにおける中間値の検査

`remote` 関数から `stderr` または `stdout` に書き込みを行います。その後、`/aws-glue/ray/jobs/ray-worker-out-logs/` または `/aws-glue/ray/jobs/ray-worker-err-logs/` から

ログを取得します。関数はどのレプリカでも実行されている可能性があるため、目的の出力を見つけるには複数のログを調べなければならない場合があります。

## 問題の領域: エラーメッセージ内の IP アドレスの解釈

特定のエラー状況では、ジョブは IP アドレスを含むエラーメッセージを出力することがあります。これらの IP アドレスは一時的なもので、クラスターがノードを識別し、ノード間で通信するために使用されます。ノードのログは、IP アドレスに基づく一意のサフィックスが付いたログストリームに発行されます。

CloudWatch では、このサフィックスを識別することでログを絞り込み、この IP アドレス固有のログを検査できます。例えば、*FAILED\_IP* と *JOB\_RUN\_ID* が指定されている場合、サフィックスを以下のように識別できます。

```
filter @logStream like /JOB_RUN_ID/
| filter @message like /IP-/
| parse @message "IP-[*]" as ip
| filter ip like /FAILED_IP/
| fields replace(ip, ":", "_") as uIP
| stats count_distinct by uIP as logStreamSuffix
| display logStreamSuffix
```

## AWS Glue 機械学習の例外

このトピックでは、機械学習に関連する AWS Glue 例外の HTTP エラーコードとエラー文字列について説明します。機械学習アクティビティごとに、オペレーションの実行時に発生する可能性があるエラーコードとエラー文字列が表示されます。また、エラーが発生したオペレーションを再試行できるかどうかを確認できます。

### CancelMLTaskRunActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」
  - 「[taskRunId] の ML タスク実行が見つかりません: 変換 [transformName] のアカウント [accountId] 内。」

再試行してもいいですか: いいえ。

## CreateMLTaskRunActivity

このアクティビティには、次の例外があります。

- InvalidInputException (400)
  - 「予期しない入力による内部サービスエラー。」
  - 「AWS Glue テーブルの入カソースを変換に指定する必要があります。」
  - 「入カソース列 [columnName] で、カタログに無効なデータ型が定義されています。」
  - 「入カレコードテーブルは 1 つだけ指定する必要があります。」
  - 「データベース名を指定する必要があります。」
  - 「テーブル名を指定する必要があります。」
  - 「スキーマが変換に定義されていません。」
  - 「スキーマには、特定の主キー [PrimaryKey] を含める必要があります。」
  - 「データカタログスキーマの取得中に問題が発生しました: [message]。」
  - 「最大容量とワーカー番号/タイプを同時に設定することはできません。」
  - 「WorkerType と NumberOfWorkers の両方を設定する必要があります。」
  - 「MaxCapacity は  $\geq$  [maxCapacity] にする必要があります。」
  - 「NumberOfWorkers は  $\geq$  [maxCapacity] にする必要があります。」
  - 「最大再試行数は負以外にする必要があります。」
  - 「一致検索パラメータが設定されていません。」
  - 「一致検索パラメータに主キーを指定する必要があります。」

再試行してもいいですか: いいえ。

- AlreadyExistsException (400)
  - 「[transformName] という名前の変換はすでに存在します。」

再試行してもいいですか: いいえ。

- IdempotentParameterMismatchException (400)
  - 「変換 [transformName] のべき等作成リクエストにパラメータの不一致がありました。」

再試行してもいいですか: いいえ。

- InternalServiceException (500)

- 「依存関係エラー。」

再試行してもいいですか: はい。

- ResourceNumberLimitExceededException (400)

- 「ML 変換数 ([count]) が [limit] 変換数の制限を超えています。」

再試行してもいいですか: はい。変換を削除して、この新しい変換用のスペースを作成してから再試行できます。

## DeleteMLTransformActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)

- 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません」

再試行してもいいですか: いいえ。

## GetMLTaskRunActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)

- 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」
- 「[taskRunId] の ML タスク実行が見つかりません: 変換 [transformName] のアカウント [accountId] 内。」

再試行してもいいですか: いいえ。

## GetMLTaskRunsActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)

- 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」

- 「[taskRunId] の ML タスク実行が見つかりません: 変換 [transformName] のアカウント [accountId] 内。」

再試行してもいいですか: いいえ。

## GetMLTransformActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」

再試行してもいいですか: いいえ。

## GetMLTransformsActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」

再試行してもいいですか: いいえ。

- InvalidInputException (400)
  - 「アカウント ID を空白にすることはできません。」
  - 「列 [column] での並べ替えはサポートされていません。」
  - 「[column] を空白にすることはできません。」
  - 「予期しない入力による内部サービスエラー。」

再試行してもいいですか: いいえ。

## GetSaveLocationForTransformArtifactActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」

再試行してもいいですか: いいえ。

- InvalidInputException (400)
  - 「サポートされていないアーティファクトタイプ [artifactType]。」
  - 「予期しない入力による内部サービスエラー。」

再試行してもいいですか: いいえ。

## GetTaskRunArtifactActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」
  - 「[taskRunId] の ML タスク実行が見つかりません: 変換 [transformName] のアカウント [accountId] 内。」

再試行してもいいですか: いいえ。

- InvalidInputException (400)
  - 「ファイル名 [fileName] は発行に対して無効です。」
  - 「[taskType] タスクタイプのアーティファクトを取得できません。」
  - 「[artifactType] のアーティファクトを取得できません。」
  - 「予期しない入力による内部サービスエラー。」

再試行してもいいですか: いいえ。

## PublishMLTransformModelActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」
  - 「アカウント ID - [accountId] - および変換 ID - [transformId] で、バージョン - [version] を持つ既存のモデルが見つかりません。」

再試行してもいいですか: いいえ。

- `InvalidInputException` (400)
  - 「ファイル名 [fileName] は発行に対して無効です。」
  - 「符号なし文字列 [string] の先頭のマイナス記号が不正です。」
  - 「[string] の末尾に不正な数字があります。」
  - 「文字列値 [string] が符号なし長整数の範囲を超えています。」
  - 「予期しない入力による内部サービスエラー。」

再試行してもいいですか: いいえ。

## PullLatestMLTransformModelActivity

このアクティビティには、次の例外があります。

- `EntityNotFoundException` (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」

再試行してもいいですか: いいえ。

- `InvalidInputException` (400)
  - 「予期しない入力による内部サービスエラー。」

再試行してもいいですか: いいえ。

- `ConcurrentModificationException` (400)
  - 「競合する挿入のパラメータが一致しないため、トレーニング対象のモデルバージョンを作成できません。」
  - 「変換 ID [transformId] の ML 変換モデルが古いか、別のプロセスによって更新中です。再試行してください。」

再試行してもいいですか: はい。

## PutJobMetadataForMLTransformActivity

このアクティビティには、次の例外があります。

- `EntityNotFoundException` (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」

- 「[taskRunId] の ML タスク実行が見つかりません: 変換 [transformName] のアカウント [accountId] 内。」

再試行してもいいですか: いいえ。

- InvalidInputException (400)
  - 「予期しない入力による内部サービスエラー。」
  - 「不明なジョブメタデータタイプ [jobType]。」
  - 「更新するタスク実行 ID を指定する必要があります。」

再試行してもいいですか: いいえ。

## StartExportLabelsTaskRunActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」
  - 「アカウント ID [accountId] に変換 ID [transformId] のラベルセットがありません。」

再試行してもいいですか: いいえ。

- InvalidInputException (400)
  - 「[message]。」
  - 「指定した S3 パスは、変換と同じリージョン内にありません。」 リージョン - [region] が必要ですが、-[region] が提供されました。」

再試行してもいいですか: いいえ。

## StartImportLabelsTaskRunActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」

再試行してもいいですか: いいえ。

- InvalidInputException (400)

- 「[message]。」
- 「無効なラベルファイルパスです。」
- 「[labelPath] のラベルファイルにアクセスできません。[message]。」
- 「変換で提供された IAM ロールを使用できません。ロール: [role]。」
- 「サイズ 0 の無効なラベルファイル。」
- 「指定した S3 パスは、変換と同じリージョン内にありません。」 リージョン - [region] が必要ですが、-[region] が提供されました。」

再試行してもいいですか: いいえ。

- ResourceNumberLimitExceededException (400)
  - 「ラベルファイルが [limit] MB の制限を超えています。」

再試行してもいいですか: いいえ。ラベルファイルをいくつかの小さなファイルに分割することを検討してください。

## StartMLEvaluationTaskRunActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)
  - 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」

再試行してもいいですか: いいえ。

- InvalidInputException (400)
  - 「入力レコードテーブルは 1 つだけ指定する必要があります。」
  - 「データベース名を指定する必要があります。」
  - 「テーブル名を指定する必要があります。」
  - 「一致検索パラメータが設定されていません。」
  - 「一致検索パラメータに主キーを指定する必要があります。」

再試行してもいいですか: いいえ。

- MLTransformNotReadyException (400)
  - 「このオペレーションは READY 状態の変換にのみ適用できます。」

再試行してもいいですか: いいえ。

- InternalServiceException (500)

- 「依存関係エラー。」

再試行してもいいですか: はい。

- ConcurrentRunsExceededException (400)

- 「ML タスク実行数 [count] が [limit] タスク実行の変換制限を超えています。」
- 「ML タスク実行数 [count] が [limit] タスク実行数の制限を超えています。」

再試行してもいいですか: はい。タスク実行が完了するまで待つてから再試行できます。

## StartMLLabelingSetGenerationTaskRunActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)

- 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」

再試行してもいいですか: いいえ。

- InvalidInputException (400)

- 「入力レコードテーブルは 1 つだけ指定する必要があります。」
- 「データベース名を指定する必要があります。」
- 「テーブル名を指定する必要があります。」
- 「一致検索パラメータが設定されていません。」
- 「一致検索パラメータに主キーを指定する必要があります。」

再試行してもいいですか: いいえ。

- InternalServiceException (500)

- 「依存関係エラー。」

再試行してもいいですか: はい。

- ConcurrentRunsExceededException (400)

- 「ML タスク実行数 [count] が [limit] タスク実行の変換制限を超えています。」

再試行してもいいですか: はい。タスク実行の完了後に再試行できます。

## UpdateMLTransformActivity

このアクティビティには、次の例外があります。

- EntityNotFoundException (400)

- 「ハンドル [transformName] を持つアカウント [accountId] に ML 変換が見つかりません。」

再試行してもいいですか: いいえ。

- InvalidInputException (400)

- 「[TransformName] という名前の別の変換がすでに存在します。」
- 「[message]。」
- 「変換名を空白にすることはできません。」
- 「最大容量とワーカー番号/タイプを同時に設定することはできません。」
- 「WorkerType と NumberOfWorkers の両方を設定する必要があります。」
- 「MaxCapacity は  $\geq$  [minMaxCapacity] にする必要があります。」
- 「NumberOfWorkers は  $\geq$  [minNumWorkers] にする必要があります。」
- 「最大再試行数は負以外にする必要があります。」
- 「予期しない入力による内部サービスエラー。」
- 「一致検索パラメータが設定されていません。」
- 「一致検索パラメータに主キーを指定する必要があります。」

再試行してもいいですか: いいえ。

- AlreadyExistsException (400)

- 「[transformName] という名前の変換はすでに存在します。」

再試行してもいいですか: いいえ。

- IdempotentParameterMismatchException (400)

- 「変換 [transformName] のべき等作成リクエストにパラメータの不一致がありました。」

再試行してもいいですか: いいえ。

## AWS Glue のクォータ

AWS Support に連絡して、「AWS 全般のリファレンス」に記載されているサービスクォータに関し、[クォータの引き上げをリクエスト](#)できます。特に明記されていない限り、クォータは地域固有です。詳細については、「[AWS Glue エンドポイントとクォータ](#)」を参照してください。

# AWS Glue パフォーマンスの向上

## パフォーマンスチューニングのベースライン戦略

AWS Glue パフォーマンスを向上させるために、AWS Glue 特定のパフォーマンス関連パラメータを更新することを検討してください。パラメータのチューニングを準備する際は、次のベストプラクティスを使用します。

- 問題の特定を始める前に、パフォーマンス目標を決定します。
- パラメータのチューニングを変更する前に、メトリクスを使用して問題を特定します。

ジョブのチューニングで一貫した結果を得るには、チューニング作業のベースライン戦略を立てます。

通常、パフォーマンスチューニングは次のワークフローで実施します。

1. パフォーマンス目標を決定します。
2. メトリクスを測定します。
3. ボトルネックを特定します。
4. ボトルネックの影響を軽減します。
5. 想定目標を達成するまで、ステップ 2~4 を繰り返します。

## ジョブタイプのチューニング戦略

Spark ジョブ — 「規範的ガイダンス」の「[Apache Spark AWS Glue ジョブのパフォーマンスチューニングのベストプラクティス](#)」のガイダンスに従ってください。AWS

その他のジョブ — 他のランタイム環境で使用できるストラテジーを適応させることで、Ray や AWS Glue Python AWS Glue のシエルジョブに合わせてチューニングできます。

## AWS Glue for Apache Spark ジョブのためのパフォーマンスの向上

AWS Glue for Spark のパフォーマンスを向上させるには、AWS Glue 関連のパフォーマンスおよび Spark パラメータを更新することを検討してください。

メトリクスを通じてボトルネックを特定し、その影響を軽減するための具体的な戦略に関する詳細は、「AWS 規範的ガイダンス」の「[Best practices for performance tuning AWS Glue for Apache](#)

[Spark jobs](#)」を参照してください。このガイドでは、Spark アーキテクチャや Resilient Distributed Datasets など、すべてのランタイム環境で Apache Spark に適用できる主なトピックを紹介し、これらのトピックを参考に、シャッフルの最適化やタスクの並列化など、特定のパフォーマンスチューニング戦略を実装する方法を案内します。

AWS Glue を設定して Spark UI を表示することにより、ボトルネックを特定できます。詳細については、「[the section called “Spark UI を使用したモニタリング”](#)」を参照してください。

さらに、AWS Glue はジョブが接続する特定タイプのデータストアに適用できるパフォーマンス機能も提供します。データストアのパフォーマンスパラメータに関する参照情報は、[the section called “接続パラメータ”](#) にあります。

## AWS Glue ETL のプッシュダウンによる読み取りの最適化

プッシュダウンは、データのソースにより近いデータの取得に関するロジックをプッシュする最適化手法です。ソースは、データベースでも、Amazon S3 などのファイルシステムでもかまいません。特定の操作をソース上で直接実行する場合、ネットワーク上のすべてのデータを AWS Glue が管理する Spark エンジンに送らないようにすることで、時間と処理能力を節約できます。

つまり、プッシュダウンによってデータスキャンが減ります。この手法が適切であることを特定するプロセスの詳細については、「AWS 規範的ガイダンス」の「Best practices for performance tuning AWS Glue for Apache Spark jobs」ガイドにある「[Reduce the amount of data scan](#)」を参照してください。

### Amazon S3 に保存されているファイルの述語プッシュダウン

プレフィックス別に整理された Amazon S3 上のファイルを操作する場合、プッシュダウン述語を定義することでターゲット Amazon S3 パスをフィルタリングできます。データセット全体を読み込み、DynamicFrame 内のフィルタを適用する代わりに、AWS Glue データカタログに保存されているパーティションメタデータにフィルタを直接適用できます。この方法では、必要なデータだけを選択的に一覧表示して読み込むことができます。パーティションごとのバケットへの書き込みなど、このプロセスの詳細については、「[the section called “パーティションの管理”](#)」を参照してください。

Amazon S3 で述語プッシュダウンを実現するには、push\_down\_predicate パラメータを使用します。Amazon S3 のバケットを年、月、日ごとに分割したとします。2022 年 6 月の顧客データを取得したい場合は、関連する Amazon S3 パスのみを読み込むように AWS Glue に指示できます。この場合の push\_down\_predicate は year='2022' and month='06' です。まとめると、読み込み操作は以下のようになります。

## Python

```
customer_records = glueContext.create_dynamic_frame.from_catalog(
 database = "customer_db",
 table_name = "customer_tbl",
 push_down_predicate = "year='2022' and month='06'"
)
```

## Scala

```
val customer_records = glueContext.getCatalogSource(
 database="customer_db",
 tableName="customer_tbl",
 pushDownPredicate="year='2022' and month='06'"
).getDynamicFrame()
```

前のシナリオでは、`push_down_predicate` は AWS Glue データカタログからすべてのパーティションのリストを取得してフィルタリングしてから、基になる Amazon S3 ファイルを読み込みます。これはほとんどの場合に役立ちますが、数百万のパーティションがあるデータセットを扱う場合、パーティションを一覧表示するプロセスには時間がかかる場合があります。この問題に対処するには、サーバー側でパーティションをプルーニングすることでパフォーマンスを向上させることができます。そのためには、AWS Glue データカタログでデータのパーティションインデックスを作成します。パーティションインデックスについての詳細は、「[the section called “パーティションインデックスの使用”](#)」を参照してください。その後、`catalogPartitionPredicate` オプションを使用してインデックスを参照できます。`catalogPartitionPredicate` を使用してパーティションを取得する例については、「[the section called “カタログパーティション述語”](#)」を参照してください。

## JDBC ソースを操作するときのプッシュダウン

`GlueContext` で使用されている AWS Glue JDBC リーダーは、ソース上で直接実行できるカスタム SQL クエリを提供することで、サポートされているデータベースへのプッシュダウンをサポートします。これは `sampleQuery` パラメータを設定することで実現できます。サンプルクエリでは、選択する列を指定したり、Spark エンジンに転送されるデータを制限するプッシュダウン述語を指定したりできます。

デフォルトでは、サンプルクエリは 1 つのノードで動作するため、大量のデータを処理するとジョブが失敗する可能性があります。この機能を使用してデータを大規模にクエリするに

は、`enablePartitioningForSampleQuery` を `true` に設定してクエリパーティショニングを設定する必要があります。これにより、選択したキーにまたがる複数のノードにクエリが分散されます。クエリパーティショニングには、他にもいくつかの必要な設定パラメータがあります。クエリパーティショニングについての詳細は、「[the section called “JDBC からの並列読み取り”](#)」を参照してください。

`enablePartitioningForSampleQuery` を設定すると、AWS Glue はデータベースをクエリするときに、プッシュダウン述語とパーティショニング述語を組み合わせます。パーティショニング条件を追加するには、AWS Glue の `sampleQuery` が AND で終わる必要があります (プッシュダウン述語を指定しない場合は、`sampleQuery` は WHERE で終わる必要があります)。以下の例を参照してください。ここでは、述語をプッシュダウンして `id` が 1000 を超える行のみを取得しています。この `sampleQuery` により、`id` が指定した値より大きい行の名前と場所の列だけが返されます。

## Python

```
sample_query = "select name, location from customer_tbl WHERE id>=1000 AND"
customer_records = glueContext.create_dynamic_frame.from_catalog(
 database="customer_db",
 table_name="customer_tbl",
 sample_query = "select name, location from customer_tbl WHERE id>=1000 AND",

 additional_options = {
 "hashpartitions": 36 ,
 "hashfield":"id",
 "enablePartitioningForSampleQuery":True,
 "sampleQuery":sample_query
 }
)
```

## Scala

```
val additionalOptions = Map(
 "hashpartitions" -> "36",
 "hashfield" -> "id",
 "enablePartitioningForSampleQuery" -> "true",
 "sampleQuery" -> "select name, location from customer_tbl WHERE id >= 1000
AND"
)

val customer_records = glueContext.getCatalogSource(
 database="customer_db",
 tableName="customer_tbl").getDynamicFrame()
```

**Note**

`customer_tbl` が Data Catalog と基になるデータストアで異なる名前がある場合、クエリは基になるデータストアに渡されるため、`sample_query` に基になるテーブル名を指定する必要があります。

AWS Glue データカタログと統合せずに JDBC テーブルに対してクエリを実行することもできます。ユーザー名とパスワードをパラメータとしてメソッドに提供する代わりに、`useConnectionProperties` と `connectionName` を指定することで既存の接続の認証情報を再利用できます。この例では、`my_postgre_connection` という接続から認証情報を取得します。

## Python

```
connection_options_dict = {
 "useConnectionProperties": True,
 "connectionName": "my_postgre_connection",
 "dbtable": "customer_tbl",
 "sampleQuery": "select name, location from customer_tbl WHERE id>=1000 AND",
 "enablePartitioningForSampleQuery": True,
 "hashfield": "id",
 "hashpartitions": 36
}

customer_records = glueContext.create_dynamic_frame.from_options(
 connection_type="postgresql",
 connection_options=connection_options_dict
)
```

## Scala

```
val connectionOptionsJson = """
{
 "useConnectionProperties": true,
 "connectionName": "my_postgre_connection",
 "dbtable": "customer_tbl",
 "sampleQuery": "select name, location from customer_tbl WHERE id>=1000 AND",

```

```
 "enablePartitioningForSampleQuery" : true,
 "hashfield" : "id",
 "hashpartitions" : 36
 }
 ""

 val connectionOptions = new JsonOptions(connectionOptionsJson)

 val dyf = glueContext.getSource("postgresql",
connectionOptions).getDynamicFrame()
```

## AWS Glue のプッシュダウンに関する注意事項と制限事項

プッシュダウンは概念上、ストリーミング以外のソースから読み込む場合にも適用できません。AWS Glue はさまざまなソースに対応しています。プッシュダウン機能はソースとコネクタによって異なります。

- Snowflake に接続するときは、`query` オプションを使用できます。AWS Glue 4.0 以降のバージョンの Redshift コネクタにも同様の機能があります。Snowflake から `query` を使用して読み込む方法については、「[the section called “Snowflake から読み込む”](#)」を参照してください。
- DynamoDB ETL リーダーでは、フィルタまたはプッシュダウン述語は使用できません。MongoDB と DocumentDB もこの種の機能をサポートしていません。
- Amazon S3 に保存されているデータをオープンテーブル形式で読み込む場合、Amazon S3 内のファイルのパーティショニング方法ではもはや十分ではありません。オープンテーブル形式を使用してパーティションから読み込みと書き込みを行うには、その形式のドキュメントを参照してください。
- `DynamicFrame` メソッドは Amazon S3 プロジェクションプッシュダウンを実行しません。述語フィルタを通過したファイルからすべての列が読み込まれます。
- AWS Glue で `custom.jdbc` コネクタを操作する場合、プッシュダウンできるかどうかはソースとコネクタによって異なります。該当するコネクタのマニュアルを参照して、AWS Glue のプッシュダウンをサポートしているかどうか、またどのようにサポートしているかを確認してください。

## AWS Glue 向けの Auto Scaling の使用

Auto Scaling は、AWS Glue ETL および AWS Glue バージョン 3.0 以降のストリーミングジョブで使用できます。

Auto Scaling を有効にすると、次の利点が得られます。

- AWS Glue は、ジョブ実行の各ステージでの並行処理またはマイクロバッチに応じて自動的にワーカーをクラスターに追加、またはクラスターから削除します。
- AWS Glue ETL ジョブに割り当てるワーカーの数について実験して判断する必要性を排除します。
- 最大数のワーカーを選択する場合、AWS Glue がワークロードに適したサイズのリソースを選択します。
- ジョブの実行中にクラスターのサイズがどのように変化するかについては、AWS Glue Studio のジョブ実行の詳細ページで CloudWatch メトリクスを確認できます。

AWS Glue ETL およびストリーミングジョブ向けの Auto Scaling は、AWS Glue ジョブのコンピューティングリソースのオンデマンドでのスケールアップとスケールダウンを可能にします。オンデマンドのスケールアップは、ジョブ実行の起動当初に必要なコンピューティングリソースを割り当てるだけでなく、ジョブ実行中の需要に応じて必要なリソースをプロビジョニングするためにも役立ちます。

Auto Scaling は、ジョブの過程全体での AWS Glue ジョブリソースの動的なスケールダウンもサポートします。ジョブ実行中に Spark アプリケーションにより多くのエグゼキューターが要求されると、クラスターに追加されるワーカーが増えます。エグゼキューターがアクティブな計算タスクなしでアイドル状態になると、エグゼキューターと対応するワーカーが削除されます。

Auto Scaling が Spark アプリケーションのコストと使用率に役立つ一般的なシナリオを含める Spark ドライバーは、Amazon S3 内の多数のファイルを一覧表示したり、エグゼキューターが非アクティブのときにロードを実行したり、オーバプロビジョニングのために少数のエグゼキューターだけで実行される Spark ステージ、Spark ステージ間でデータスキューや不均一な計算要求を実行します。

## 要件

Auto Scaling は、AWS Glue バージョン 3.0 以降でのみ利用可能です。Auto Scaling を使用するには、「[移行ガイド](#)」に従って既存のジョブをバージョン 3.0 以降の AWS Glue に移行するか、AWS Glue のバージョン 3.0 以降を使用して新しいジョブを作成します。

Auto Scaling は AWS Glue のジョブと G.1X、G.2X、G.4X、G.8X、または G.025X (Streaming ジョブ用) ワーカータイプで利用可能です。スタンダード DPU はサポートされていません。

## AWS Glue Studio の Auto-Scaling を有効にする

AWS Glue Studio の [Job details] (ジョブ詳細) タブで、タイプに [Spark] または [Spark Streaming] を選択し、[Glue version] (Glue バージョン) には **[Glue 3.0]** もしくは **[Glue 4.0]** を選択します。次に、[Worker type] (ワーカータイプ) の下にチェックボックスが表示されます。

- [Automatically scale the number of workers] (ワーカー数を自動的にスケーリングする) を選択します。
- ワーカーの最大数をセットして、ジョブ実行に投入できるワーカーの最大数を定義します。

Visual | Script | **Job details** | Runs | Data quality | Schedules

### Version Control

**Type**  
The type of ETL job. This is set automatically based on the types of data sources you have selected.

Spark

**Glue version** [Info](#)

Glue 4.0 - Supports spark 3.3, Scala 2, Python 3 ▼

**Language**

Python 3 ▼

**Worker type**  
Set the type of predefined worker that is allowed when a job runs.

G 1X  
(4vCPU and 16GB RAM) ▼

**Automatically scale the number of workers**  
AWS Glue will optimize costs and resource usage by dynamically scaling the number of workers up and down throughout the job run. Requires Glue 3.0 or later.

**Maximum number of workers**  
The number of workers you want AWS Glue to allocate to this job.

10

## AWS CLI または SDK を使用した Auto Scaling の有効化

Auto Scaling を有効にするには、ジョブ実行の AWS CLI から、次の設定で `start-job-run` を実行します。

```
{
 "JobName": "<your job name>",
 "Arguments": {
 "--enable-auto-scaling": "true"
 },
 "WorkerType": "G.2X", // G.1X and G.2X are allowed for Auto Scaling Jobs
 "NumberOfWorkers": 20, // represents Maximum number of workers
 ...other job run configurations...
}
```

ETL ジョブ実行が終了したら、`get-job-run` を呼び出して、実行されたジョブの実際のリソース使用量を DPU 秒単位で確認することもできます。注: 新しいフィールド `[DPUSecods]` は、Auto Scaling が有効になっている AWS Glue 3.0 以降のパッチジョブに対してのみ表示されます。このフィールドは、ストリーミングジョブではサポートされません。

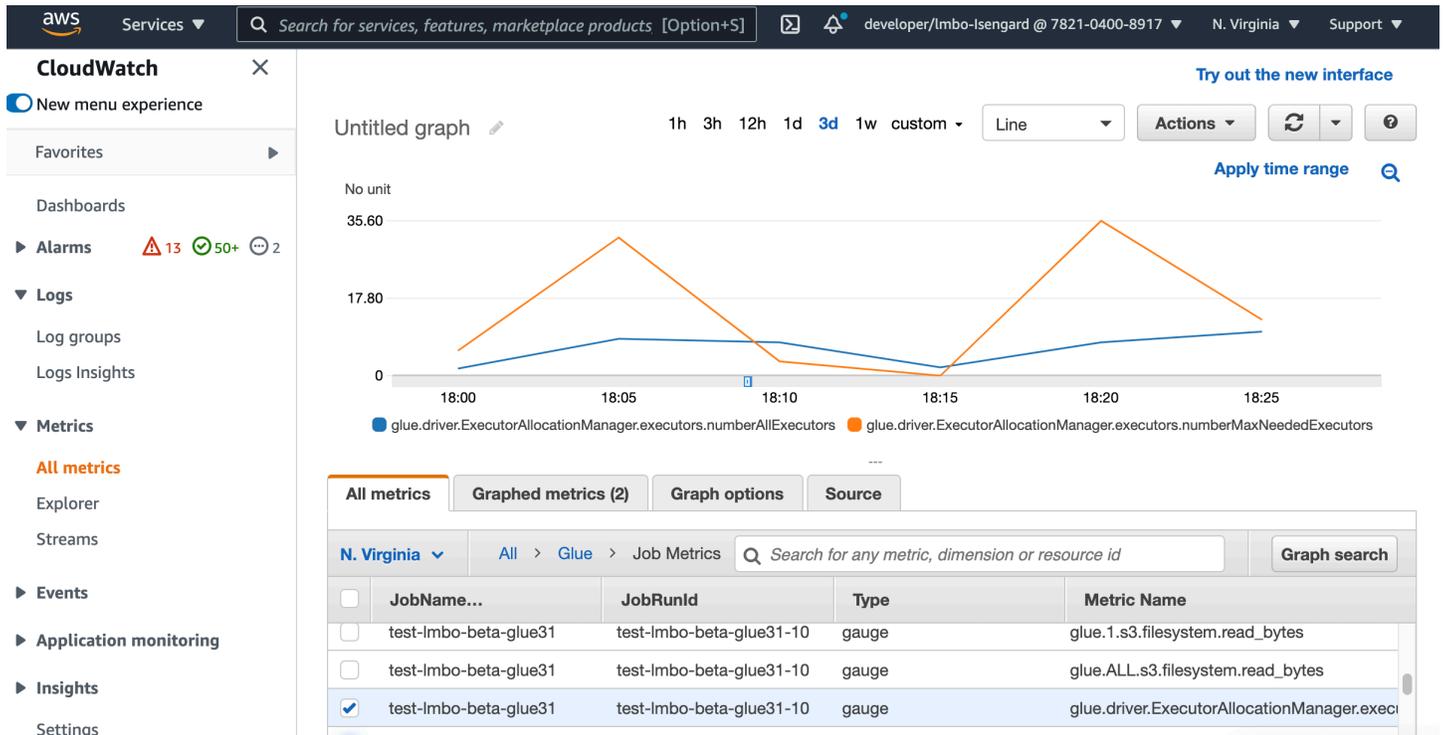
```
$ aws glue get-job-run --job-name your-job-name --run-id jr_xx --endpoint https://
glue.us-east-1.amazonaws.com --region us-east-1
{
 "JobRun": {
 ...
 "GlueVersion": "3.0",
 "DPUSecods": 386.0
 }
}
```

同じ設定で [AWS Glue SDK](#) を使用して、Auto Scaling でジョブ実行を設定することもできます。

## Amazon CloudWatch メトリクスを使用した Auto Scaling のモニタリング

Auto Scaling を有効にすると、CloudWatch の実行者メトリクスが AWS Glue 3.0 以降のジョブで利用できるようになります。メトリクスを使用して、Auto Scaling で有効化された Spark アプリケーションでのエグゼキューターの需要と最適化された使用状況をモニタリングできます。詳しくは、「[Amazon CloudWatch メトリクスを使用した AWS Glue のモニタリング](#)」を参照してください。

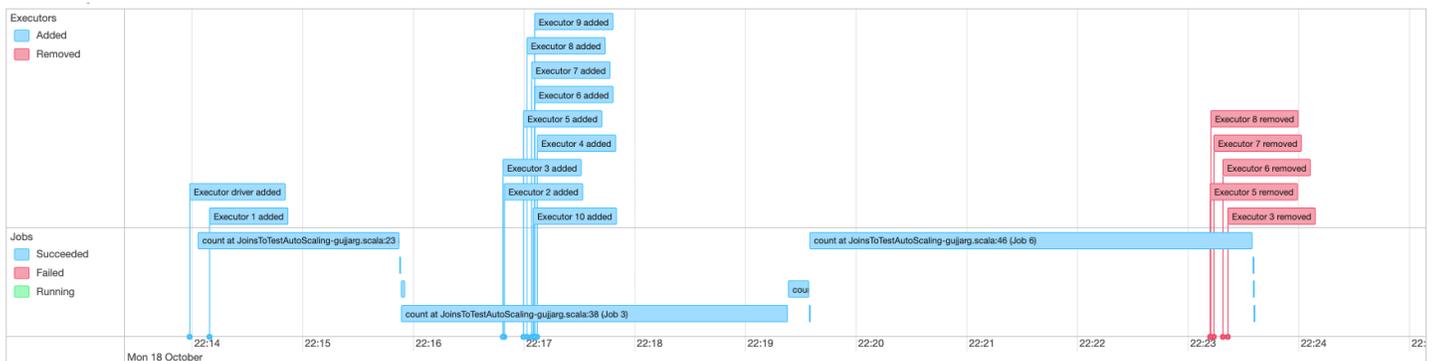
- `glue.driver.executoraAllocationManager.executors.numberallExecutor`
- `glue.driver.executorAllocationManager.executors.numbermaxNeededExecutors`



これらのメトリクスの詳細については、「[DPU の容量計画のモニタリング](#)」を参照してください。

## Spark UI による Auto Scaling のモニタリング

Auto Scaling を有効にすると、Glue Spark UI を使用して、AWS Glue ジョブの需要に基づいた動的なスケールアップとスケールダウンで追加または削除されるエグゼキューターをモニタリングすることもできます。詳しくは、「[AWS Glue ジョブ用の Apache Spark ウェブ UI の有効化](#)」を参照してください。



## Auto Scaling ジョブ実行の DPU 使用量モニタリング

[\[AWS Glue Studio ジョブ実行ビュー\]](#) を使用して、Auto Scaling ジョブの DPU 使用量をチェックします。

1. AWS Glue Studio ナビゲーションペインから [モニタリング] を選択します。[Monitoring] (モニタリング) ページが表示されます。
2. ジョブ実行チャートまでスクロールダウンします。
3. チェックしたいジョブ実行に移動して、[DPU hours] (DPU 時間) 列を下にスクロールして、そのジョブ実行の使用量をチェックします。

## 制限事項

AWS Glue ストリーミングの Auto Scaling は現在、ForEachBatch 外で作成された静的 DataFrame とのストリーミング DataFrame の結合をサポートしていません。ForEachBatch 内で作成された静的 DataFrame は、期待通りに動作します。

## 実行に上限を設定してワークロードをパーティション化する

通常、Spark アプリケーションのエラーは、Spark スクリプトが非効率的なこと、大規模な変換を分散メモリ内で実行すること、およびデータセットの異常によって発生します。ドライバーやエグゼキュターがメモリ不足の問題を引き起こす原因はたくさんあります。例えば、データの偏り、オブジェクトのリストが多すぎることで、または大きなデータのシャッフルなどです。これらの問題は、Sparkで大量のバックログデータを処理しているときにしばしば発生します。

AWS Glue を使用すると、OOM の問題を解決し、ワークロードのパーティション化によって ETL 処理を容易に行うことができます。ワークロードのパーティション化を有効にすると、各 ETL ジョブ実行では、データセットサイズまたはこのジョブ実行で処理されるファイル数に上限を設定して、未処理のデータのみを選択します。残りのデータは、後で実行されるジョブで処理されます。例えば、1000 個のファイルを処理する必要がある場合、ファイル数を 500 に設定し、2 つのジョブ実行に分けることができます。

ワークロードのパーティション化は、Amazon S3 のデータソースでのみサポートされています。

## ワークロードのパーティション化の有効化

スクリプトでオプションを手動で設定するか、カタログテーブルのプロパティを追加することで、実行に上限を設定できます。

スクリプトで実行に上限を設定してワークロードのパーティション化を有効にするには、次の手順に従います。

1. データの再処理を回避するには、新しいジョブまたは既存のジョブでジョブブックマークを有効にします。詳細については、「[Tracking Processed Data Using Job Bookmarks](#)」を参照してください。
2. スクリプトを修正し、AWS Glue getSource API の追加のオプションに上限の制限を設定します。また、ジョブブックマークの変換コンテキストを設定して、state 要素を保存します。例:

### Python

```
glueContext.create_dynamic_frame.from_catalog(
 database = "database",
 table_name = "table_name",
 redshift_tmp_dir = "",
 transformation_ctx = "datasource0",
 additional_options = {
 "boundedFiles" : "500", # need to be string
 # "boundedSize" : "1000000000" unit is byte
 }
)
```

### Scala

```
val datasource0 = glueContext.getCatalogSource(
 database = "database", tableName = "table_name", redshiftTmpDir = "",
 transformationContext = "datasource0",
 additionalOptions = JsonOptions(
 Map("boundedFiles" -> "500") // need to be string
 //"boundedSize" -> "1000000000" unit is byte
)
)
.getDynamicFrame()
```

```
val connectionOptions = JsonOptions(
 Map("paths" -> List(baseLocation), "boundedFiles" -> "30")
)
val source = glueContext.getSource("s3", connectionOptions, "datasource0", "")
```

Data Catalog テーブルで実行に上限を設定してワークロードのパーティション化を有効にするには

1. Data Catalog のテーブル構造の `parameters` フィールドに、キーと値のペアを設定します。詳細については、「[Viewing and Editing Table Details](#)」を参照してください。
2. データセットサイズの上限または処理されるファイル数を次のように設定します。
  - `boundedSize` をデータセットのターゲットサイズ (バイト単位) に設定します。ジョブの実行は、テーブルからターゲットサイズに達すると停止します。
  - `boundedFiles` をターゲットファイル数に設定します。ジョブの実行は、ターゲット数のファイル进行处理すると停止します。

#### Note

単一の上限のみがサポートされているので、`boundedSize` または `boundedFiles` のどちらか一方だけを設定します。

## ジョブを自動的に実行する AWS Glue トリガーの設定

実行に上限を設定すると、AWS Glue トリガーを設定してジョブを自動的に実行し、シーケンシャル実行でデータを段階的にロードすることができます。AWS Glue コンソールに移動してトリガーを作成し、スケジュール時刻を設定し、ジョブにアタッチします。その後、自動的に次のジョブ実行がトリガーされ、データの新しいバッチが処理されます。

また、AWS Glue ワークフローを使用して、複数のジョブを調整して、異なるパーティションのデータを並行して処理することもできます。詳細については、「[AWS Glue Triggers](#)」および「[AWS Glue Workflows](#)」を参照してください。

ユースケースとオプションの詳細については、ブログ「[Optimizing Spark applications with workload partitioning in AWS Glue](#)」を参照してください。

# AWS Glue の既知の問題

AWS Glue には以下の既知の問題があります。

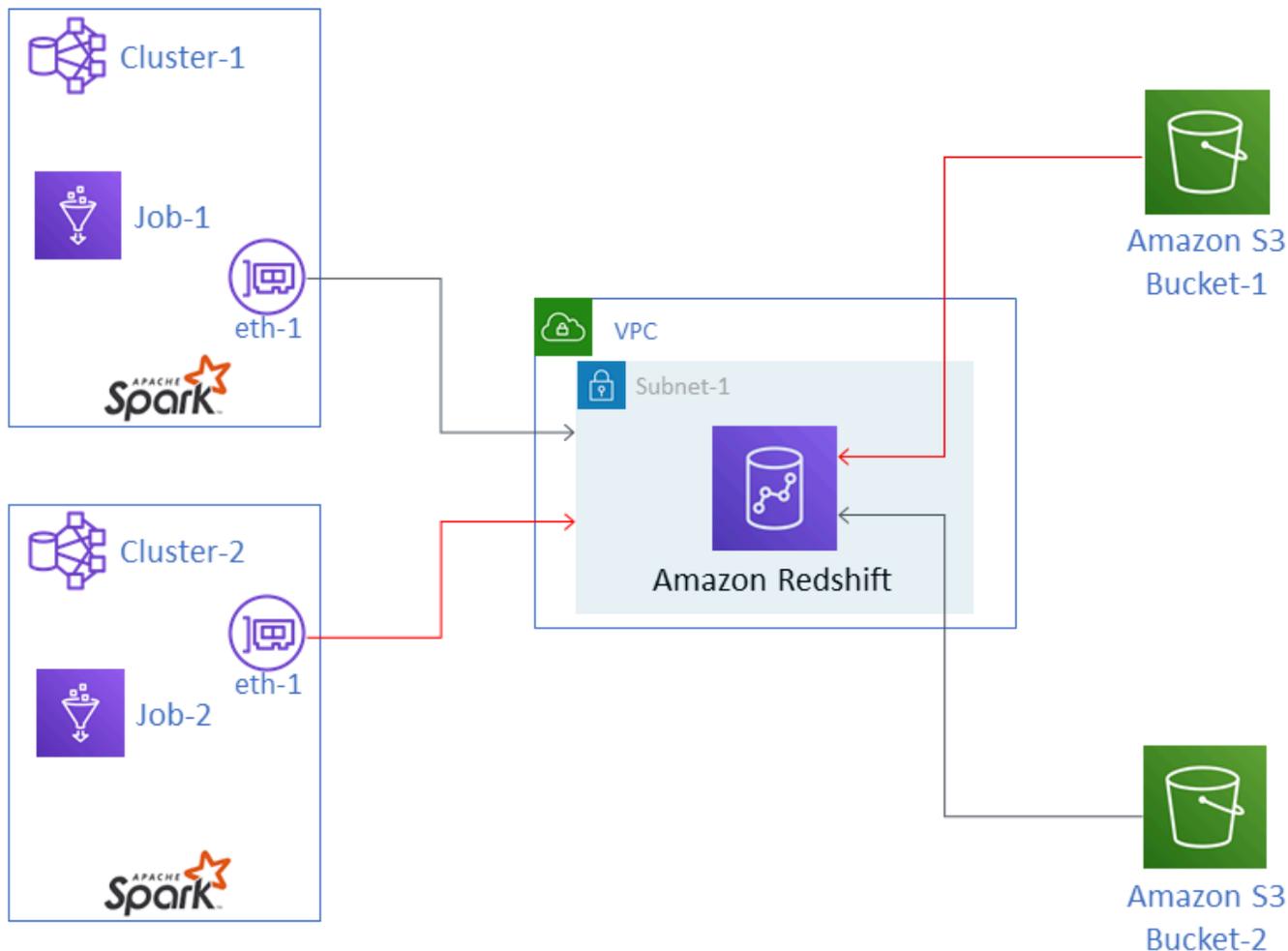
トピック

- [クロスジョブデータアクセスの防止](#)

## クロスジョブデータアクセスの防止

単一の AWS アカウントに 2 つの AWS Glue Spark ジョブがあり、それぞれが別の AWS Glue Spark クラスターで実行されている状況を考慮します。ジョブは、AWS Glue 接続を使用して、同じ Virtual Private Cloud (VPC) 内のリソースにアクセスしています。この状況では、一方のクラスターで実行されているジョブが、もう一方のクラスターで実行されているジョブのデータにアクセスできる可能性があります。

次の図は、この状況の例を示しています。



この図では、AWS Glue Job-1 は Cluster-1 で実行され、Job-2 は Cluster-2 で実行されています。どちらのジョブも、VPC の Subnet-1 に存在する Amazon Redshift の同じインスタンスを使用しています。Subnet-1 は、パブリックサブネットであることもプライベートサブネットであることもあります。

Job-1 は、Amazon Simple Storage Service (Amazon S3) Bucket-1 からデータを変換し、データを Amazon Redshift に書き込んでいます。Job-2 は Bucket-2 のデータで同じ処理を行っています。Job-1 は、Bucket-1 へのアクセスを許可する AWS Identity and Access Management (IAM) ロール Role-1 (非表示) を使用しています。Job-2 は、Bucket-2 へのアクセスを許可する Role-2 (非表示) を使用しています。

この2つのジョブにはネットワークパスがあり、相互のクラスターと通信し、相互のデータにアクセスできるようになっています。たとえば、Job-2 は Bucket-1 のデータにアクセスできます。この図では、これは赤色のパスとして示されています。

このような状況を回避するため、Job-1 および Job-2 に異なるセキュリティ設定をアタッチすることをお勧めします。セキュリティ設定をアタッチすることで、データへのクロスジョブアクセスは AWS Glue が作成する証明書によってブロックされます。セキュリティ設定は、ダミー設定にすることができます。つまり、Amazon S3 データ、Amazon CloudWatch データ、ジョブのブックマークの暗号化を有効にすることなく、セキュリティ設定を作成できます。3つの暗号化オプションはすべて無効にできます。

セキュリティ設定の詳細については、「[the section called “AWS Glue によって書き込まれたデータの暗号化”](#)」を参照してください。

セキュリティ設定をジョブにアタッチするには

1. <https://console.aws.amazon.com/glue/> で AWS Glue コンソール を開きます。
2. ジョブの [Configure the job properties (ジョブプロパティの設定)] ページで、[セキュリティ設定、スクリプトライブラリおよびジョブパラメータ] セクションを展開します。
3. リストでセキュリティ設定を選択します。

# AWS Glue のドキュメント履歴

変更	説明	日付
<a href="#">AWS Glue 使用プロファイルのサポート</a>	管理者は、デベロッパー、テスター、製品チームなど、アカウント内のさまざまなクラスのユーザーの AWS Glue 使用状況プロファイルを作成できます。この柔軟性により、管理者はユーザーのクラスごとに異なる使用量とコスト制御を適用できます。詳細については、 <a href="#">AWS Glue 「使用プロファイルの設定」</a> を参照してください。	2024 年 6 月 18 日
<a href="#">AWS Glue for Spark の Salesforce コネクタのサポート</a>	Salesforce の新しい AWS Glue コネクタに関する情報を追加しました。この機能を使用すると、AWS Glue for Spark を使用して AWS Glue 4.0 以降のバージョンで Salesforce との間で読み書きを行うことができます。詳細については、「 <a href="#">Salesforce への接続</a> 」を参照してください。	2024 年 5 月 22 日
<a href="#">AWS Glue (GA) での Amazon Q データ統合</a>	の Amazon Q データ統合 AWS Glue は、データエンジニアと ETL デベロッパー AWS Glue が自然言語を使用してデータ統合ジョブを構築できるようにする の新しい生成 AI 機能です。エンジニア	2024 年 4 月 30 日

とデベロッパーは、ジョブの作成、問題のトラブルシューティング、と AWS Glue データ統合に関する質問に回答するように Q に依頼できます。詳細については、「[AWS Glue の Amazon Q データ統合](#)」を参照してください。

この機能には、`AwsGlueSessionUserRestrictedPolicy`、`AwsGlueSessionUserRestrictedNotebookServiceRole`、`AwsGlueSessionUserRestrictedServiceRole` の更新が含まれます。詳細については、「[AWS Glue AWS 管理ポリシーの更新](#)」を参照してください。

## [での Amazon Q データ統合 AWS Glue \(プレビュー\)](#)

2024 年 1 月 30 日

の Amazon Q データ統合  
AWS Glue は、データエンジニアと ETL デベロッパー  
AWS Glue が自然言語を使用してデータ統合ジョブを構築できるようにするの新しい生成 AI 機能です。エンジニアとデベロッパーは、ジョブの作成、問題のトラブルシューティング、と AWS Glue データ統合に関する質問に回答するように Q に依頼できます。詳細については、「[AWS Glue の Amazon Q データ統合](#)」を参照してください。  
この機能には、`AwsGlueSessionUserRestrictedNotebookPolicy` AWS マネージドポリシーの更新が含まれます。詳細については、「[AWS Glue AWS 管理ポリシーの更新](#)」を参照してください。

## [AWS Glue ストリーミングのドキュメントの更新](#)

AWS Glue ストリーミングの新しいコンテンツと再編成されたコンテンツを含む新しい章を追加しました。このコンテンツでは AWS Glue、ストリーミングが とどのように連携するか、リアルタイムデータ処理の特性、ストリーミングジョブのモニタリング方法について説明します。詳細については、「[AWS Glue Streaming](#)」を参照してください。

2023 年 12 月 27 日

## [詳細な機密データ検出の使用をサポート](#)

Detect Sensitive Data は、定義したエンティティ、または AWS Glue によって事前定義されたエンティティを検出、マスキング、削除する機能を提供します。さらに、詳細なアクションを使用することで、エンティティごとに特定のアクションを適用できます。詳細については、「[詳細な機密データ検出の使用](#)」を参照してください。

2023 年 11 月 26 日

### [オブザー AWS Glue バビリティメトリクスによるジョブのモニタリングのサポート](#)

AWS Glue オブザーバビリティメトリクスを使用して、AWS Glue for Apache Spark ジョブの内部で何が起きているかに関するインサイトを生成し、問題の優先順位付けと分析を改善できます。詳細については、「[AWS Glue オブザーバビリティメトリクスを使用したモニタリング](#)」を参照してください。

2023 年 11 月 26 日

### [AWS Glue Data Quality での異常検出のサポート](#)

AWS Glue Data Quality の異常検出は、時間の経過とともにデータ統計に機械学習 (ML) アルゴリズムを適用して、ルールでは検出が難しい異常パターンや隠れたデータ品質問題を検出します。詳細については、「[AWS Glue Data Quality での異常検出](#)」を参照してください。

2023 年 11 月 26 日

## [デフォルトの Spark UI ログ動作を更新](#)

Spark UI ログを生成する Spark ジョブは、AWS Glue コンソールで Spark UI をサポートするために、別のファイル名パターンで書き込まれるようになりました。これにより CloudWatch ログの動作は変更されません。ジョブ設定を更新することで、従来の動作に戻すことができます。詳細については、「[Apache Spark ウェブ UI を使用したジョブのモニタリング](#)」を参照してください。

2023 年 11 月 17 日

## [AWS Glue for Spark での新しいデータソースのサポート](#)

Amazon OpenSearch Service、Azure SQL、Azure Cosmos for NoSQL、SAP HANA Teradata Vantage、Vertica への接続が内でネイティブにサポートされるようになりました AWS Glue。さらに、これらのデータソースへの接続と MongoDB が AWS Glue Studio ビジュアルエディタで使用できるようになりました。詳細については、[for Spark のサポートについては AWS Glue for Spark の ETL の接続タイプとオプション、AWS Glue Studio ビジュアルエディタでの使用については AWS Glue](#) 「接続の追加」を参照してください。

2023 年 11 月 17 日

## [カラム統計生成のサポート](#)

追加のデータパイプラインを設定せずに、Parquet、ORC、JSON、ION、CSV、XML などのデータ形式の AWS Glue Data Catalog テーブルの列レベルの統計を計算できます。詳細については、「[カラム統計の操作](#)」を参照してください。

2023 年 11 月 16 日

## [Iceberg テーブルのデータ圧縮のサポート](#)

Amazon Athena や Amazon EMR などの AWS 分析サービス、および AWS Glue ETL ジョブによる読み取りパフォーマンスを向上させるために、Data Catalog は、Data Catalog の Iceberg テーブルに対してマネージド圧縮 (小さな Amazon S3 オブジェクトを大きなオブジェクトに圧縮するプロセス) を提供します。詳細については、「[Iceberg テーブルの最適化](#)」を参照してください。

2023 年 11 月 13 日

## [ジョブ実行待機動作を更新する](#)

標準の Spark および Python シェルのジョブ実行は、すぐに FAILED に移行するのではなく、特定の状況において WAITING に移行するようになりました。詳細については、「[AWS Glue ジョブ実行ステータス](#)」を参照してください。

2023 年 11 月 8 日

[AWS Glue Studio ユーザーガイドを AWS Glue デベロッパーガイドに統合](#)

AWS Glue Studio ユーザーガイドがデベロッパーガイドに移され、AWS Glue Studio、AWS Glue コンソール、および AWS Glue Studio プログラムによるアクセスに関する単一の統合ユーザーガイドが作成されました。

2023 年 10 月 25 日

[AWSGlueServiceNote bookRole AWS 管理ポリシーの更新](#)

AWSGlueServiceNote bookRole AWS 管理ポリシーのマイナー更新に関する情報を追加しました。詳細については、[AWS Glue 「AWS 管理ポリシーの更新」](#)を参照してください。

2023 年 10 月 9 日

[AWS Glue Studio が 5 つの新しい組み込み変換をサポート](#)

AWS Glue Studio は、レコードマッチング、null 行の削除、JSON 列の解析、JSON パスの抽出、正規表現抽出の 5 つの新しい組み込み変換をサポートしています。詳細については、[AWS Glue 「マネージドデータ変換ノードの編集」](#)を参照してください。

2023 年 8 月 11 日

[AWSGlueServiceRole AWS 管理ポリシーの更新](#)

AWSGlueServiceRole AWS 管理ポリシーのマイナー更新に関する情報を追加しました。詳細については、[AWS Glue 「AWS 管理ポリシーの更新」](#)を参照してください。

2023 年 8 月 4 日

## [Apache Hudi テーブルのクローリングをサポート](#)

AWS Glue を使用して Amazon S3 バケット内の Hudi テーブルをクロールし、Hudi テーブルを に登録する方法について説明します AWS Glue Data Catalog。詳細については、「[クロール可能なデータストア](#)」および「[クローラーのプロパティ](#)」を参照してください。

2023 年 7 月 21 日

## [AWSGlueConsoleFullAccess AWS 管理ポリシーの更新](#)

AWSGlueConsoleFullAccess AWS 管理ポリシーのマイナー更新に関する情報を追加しました。詳細については、[AWS Glue 「AWS 管理ポリシーの更新」](#)を参照してください。

2023 年 7 月 14 日

## [Apache Iceberg テーブルのクローリングをサポート](#)

AWS Glue を使用して Amazon S3 バケット内の Iceberg テーブルをクロールし、Iceberg テーブルを に登録する方法について説明します AWS Glue Data Catalog。詳細については、「[クロール可能なデータストア](#)」および「[クローラーのプロパティ](#)」を参照してください。

2023 年 7 月 7 日

## [Ray AWS Glue での のサポート](#)

AWS Glue ジョブ AWS Glue をバックアップできる新しいエンジンである Ray でに関する情報を追加しました。Spark コンテンツ AWS Glue で既存のを再編成して、あいまいさを解き明かしました。

2023 年 5 月 30 日

## [AWS Glue データ品質 \(GA\) のサポート](#)

AWS Glue Data Quality が一般公開されました。AWS Glue Data Quality は、データの品質を評価およびモニタリングするのに役立ちます。Data Catalog で AWS Glue Data Quality を使用方法については、「[Data AWS Glue Quality](#)」を参照してください。のデータ AWS Glue 品質の詳細については AWS Glue Studio、「[によるデータ品質の評価 AWS Glue Studio](#)」を参照してください。

2023 年 5 月 24 日

## [Apache Spark ジョブの大規模ワーカータイプのサポート](#)

Apache Spark ジョブ用に G.4X と G.8X のワーカータイプの使用がサポートされるようになりました。これらのワーカータイプは、ワークロードに含まれる変換、集約、結合、クエリへの要求が非常に厳しいジョブに適しています。詳細については、「[AWS Glue でジョブを追加する](#)」を参照してください。

2023 年 5 月 8 日

## [テーブルをクローलするときのパーティションインデックス作成のサポート](#)

クローラーが検出したテーブルのパーティションインデックスの作成をクローラーがどのようにサポートするかについての情報を追加しました。詳細については、「[クローラー設定オプションの設定](#)」を参照してください。

2023 年 4 月 24 日

[リソースの使用状況メトリクスのサポート](#)

Amazon でのサービスのリソース使用状況の表示とアラームの設定に関する情報を追加しました CloudWatch。詳細については、「[AWS Glue のモニタリング](#)」を参照してください。

2023 年 4 月 7 日

[AWSGlueConsoleFullAccess AWS 管理ポリシーの更新](#)

AWSGlueConsoleFullAccess AWS 管理ポリシーのマイナー更新に関する情報を追加しました。詳細については、[AWS Glue 「AWS 管理ポリシーの更新」](#)を参照してください。

2023 年 3 月 28 日

[AWS SDK AWS Glueで 使用するためのガイダンスと例を追加](#)

AWS Glue デベロッパーガイドには、AWS SDK AWS Glueでの使用に役立つ情報を提供する 2 つの新しいセクションがあります。詳細については、[AWS 「SDK AWS Glue で使用する」](#)および「[SDK AWS Glue を使用するためのコード例 AWS SDKs](#)」を参照してください。

2023 年 2 月 23 日

[を使用した IAM のドキュメントの更新 AWS Glue](#)

で IAM を使用するための情報を再編成し、追加しました AWS Glue。詳細については、「[AWS GlueのIdentity and Access Management](#)」を参照してください。

2023 年 2 月 15 日

### [AWS Glue バージョン 4.0 でストリーミング ETL ジョブをサポート](#)

Glue バージョン 4.0 でのストリーミング ETL ジョブの実行のサポートに関する情報と、Kafka クラスターまたは Amazon Managed Streaming for Apache Kafka クラスター、および Amazon Kinesis Data Streams に接続するための新しいオプションに関する情報が追加されました。詳細については、「[AWS Glue でストリーミング ETL ジョブを追加する](#)」および「[AWS Glue での ETL の接続タイプとオプション](#)」を参照してください。

2023 年 2 月 8 日

### [MongoDB Atlas データソースのクローリングをサポート](#)

を使用して MongoDB Atlas データソース AWS Glue をクローリングする方法に関する情報を追加しました。詳細については、「[クローリング可能なデータストア](#)」、「[MongoDB と MongoDB Atlas の接続プロパティ](#)」、「[MongoDB または MongoDB Atlas 接続の使用](#)」を参照してください。

2023 年 2 月 6 日

## [ネイティブの Delta Lake コネクタを使用した Delta Lake テーブルのクローリングをサポート](#)

AWS Glue を使用して、ネイティブの Delta Lake コネクタを使用して Delta Lake テーブルをクローリングする方法に関する情報を追加しました。この機能を使用すると、AWS クエリエンジンを使用して Delta トランザクションログを直接クエリし、タイムトラベルや ACID 保証などの機能を使用できます。また、Amazon S3 トランザクションファイルから Delta Lake メタデータを Data Catalog に同期して、Lake Formation のクエリに対する列許可を有効にすることができます。詳細については、「[Delta Lake データストアの設定オプションを指定する方法](#)」および「[Delta Lake テーブルのクエリを実行する](#)」を参照してください。

2022 年 12 月 15 日

## [AWS Glue データ品質のサポート \(プレビュー\)](#)

AWS Glue Data Quality (プレビュー) のサポートが利用可能になりました。AWS Glue Data Quality は、AWS Glue 3.0 の使用時にデータの品質を評価およびモニタリングするのに役立ちます。Data Catalog で AWS Glue Data Quality を使用する方法については、[AWS Glue 「Data Quality \(プレビュー\)」](#) を参照してください。のデータ AWS Glue 品質の詳細については AWS Glue Studio、[「によるデータ品質の評価 AWS Glue Studio」](#) を参照してください。

2022 年 11 月 30 日

## [新しい機能が追加されパフォーマンスが向上した、新しい Amazon Redshift Spark コネクタのサポート](#)

AWS Glue ETL ジョブで使用できる、新しい Amazon Redshift Spark コネクタおよび JDBC ドライバーがサポートされるようになりました。これにより、データの取り込みおよび変換パイプラインの一部として Amazon Redshift のデータを読み書きする Apache Spark アプリケーションを構築できます。詳細については、[「Amazon Redshift との間でのデータの移動」](#) を参照してください。

2022 年 11 月 29 日

### [AWS Glue バージョン 4.0 のサポート。](#)

AWS Glue バージョン 4.0 のサポートに関する情報を追加しました。機能には、Apache Hudi、Delta Lake、Apache Iceberg でのオープンデータレイクフレームワークのネイティブサポートや、Amazon S3 ベースの Cloud Shuffle Storage Plugin (Apache Spark プラグイン) のネイティブサポートが含まれます。これにより、Amazon S3 を使用してシャッフルと伸縮自在なストレージ容量を実現できます。詳細については、[AWS Glue リリースノート](#)および「[AWS Glue ジョブの AWS Glue バージョン 4.0 への移行](#)」を参照してください。

2022 年 11 月 28 日

### [AWS Glue Studio がカスタムビジュアル変換機能をリリースしました。](#)

カスタムビジュアル変換機能により、お客様がビジネス固有の ETL ロジックをチーム間で定義、再利用、共有できるようになりました。詳細については、「[カスタムビジュアル変更機能](#)」を参照してください。

2022 年 11 月 28 日

### [AWS Glue クローラーを使用した JDBC データストアのメタデータの公開をサポート](#)

AWS Glue クローラーを使用して、コメントや rawtype などのメタデータを JDBC データストアのデータカタログに公開することができるようになりました。詳細については、「[クローラーによって Data Catalog テーブルに設定されたパラメータ](#)」、「[クローラープロパティ](#)」、および「[構造](#)」を参照してください。 [JdbcTarget](#)

2022 年 11 月 18 日

### [Snowflake データストアのクローリングをサポート](#)

AWS Glue を使用して、Snowflake のテーブルとビューをクロールしたり、メタデータをテーブルエントリとしてデータカタログに公開できるようになりました。Amazon S3 での Snowflake の外部テーブルの場合、クローラーは Amazon S3 の場所と外部テーブルのファイル形式のタイプもクロールし、テーブルパラメータとして入力します。詳細については、「[クロール可能なデータストア](#)」、「[AWS Glue 接続プロパティ](#)」および「[クローラーによって設定されたデータカタログテーブルのパラメータ](#)」を参照してください。

2022 年 11 月 18 日

### [Spark アプリケーションのシャッフル管理が改善](#)

Cloud Shuffle Storage Plugin for Apache Spark が利用可能になりました。詳細については、「[AWS Glue Spark シャッフルマネージャーと Amazon S3](#)」および「[Cloud Shuffle Storage Plugin for Apache Spark](#)」を参照してください。

2022 年 11 月 15 日

### [Amazon S3 イベント通知のクローラを高速化するときのデータカタログターゲットに対するサポートの追加](#)

Amazon S3 ターゲットに対する既存のサポートに加えて、Amazon S3 イベント通知を使用したデータカタログターゲットのクローラの高速化に対するサポートが利用可能になりました。詳細については、「[Amazon S3 イベント通知を使用したクローラの高速化](#)」を参照してください。

2022 年 10 月 13 日

### [クローラーが作成できるテーブルの最大数の指定をサポート](#)

クローラーが作成できるテーブルの最大数を指定できるようになりました。詳細については、「[How to specify the maximum number of tables the crawler is allowed to create](#)」(クローラーが作成できるテーブルの最大数を指定する方法)を参照してください。

2022 年 9 月 6 日

### [AWS Glue 内の Python シェルジョブでの Python 3.9 のサポート](#)

AWS Glue 内の Python シェルジョブで Python 3.9 との互換性のあるスクリプトを実行したり、事前にパッケージ済みのライブラリセットを選択できるようになりました。詳細については、「[AWS Glue の Python シェルジョブ](#)」を参照してください。

2022 年 8 月 11 日

### [緊急でないジョブや時間的制約のない AWS Glue ジョブを予備の容量で実行するためのサポート](#)

実稼働前のジョブ、テスト、1 回限りのデータの読み込みなど、緊急性のないジョブに対する柔軟なジョブ実行を設定できるようになりました。詳細については、「[AWS Glue でジョブを追加する](#)」を参照してください。

2022 年 8 月 9 日

### [ストリーミングジョブ向けの新しいワーカータイプのサポート](#)

少量のストリーミングジョブ用に G.025X ワーカータイプを使用できるようになりました。詳細については、「[AWS Glue でジョブを追加する](#)」を参照してください。

2022 年 7 月 14 日

### [AWS Glue 接続での Kafka SASL の使用のサポート](#)

AWS Glue 接続で Kafka SASL の使用がサポートされるようになりました。詳細については、「[AWS Glue Kafka connection properties for client authentication](#)」(クライアント認証用の Kafka 接続プロパティ)を参照してください。

2022 年 7 月 5 日

[Protobuf スキーマ用の Apache kafka コネクタサポート](#)

現在、Protobuf スキーマ用の Apache Kafka コネクタサポートがご利用可能です。詳細については、「[AWS Glue Schema Registry](#)」を参照してください。

2022 年 6 月 9 日

[AWS Glue ジョブ \(GA\) の Auto Scaling のサポート](#)

AWS Glue バージョン 3.0 のジョブに Auto Scaling を使用して、コンピューティングリソースを動的にスケールングする方法に関する情報を追加しました。詳細については、「[AWS Glue の Auto Scaling を利用する](#)」を参照してください。

2022 年 4 月 14 日

[AWS Glue 開発とAWS Glue ジョブスクリプトのテストに関するドキュメントの更新](#)

Docker での開発手順を含む、AWS Glue 向けとして利用可能な開発とテストの方法に関する情報を再編成して追加しました。詳細については、「[AWS Glue ETL ライブラリを使用した ETL スクリプトのローカルでの開発とテスト](#)」を参照してください。

2022 年 3 月 14 日

[AWS Glue スキーマレジストリでサポートされるデータ形式に protocol buffers \(protobuf\) を追加](#)

サポートされているデータ形式として (AVRO および JSON に加えて) Protobuf に関する情報を追加しました。詳細については、「[AWS Glue Schema Registry](#)」を参照してください。

2022 年 2 月 25 日

### [Delta Lake テーブルのクロール リングのサポート](#)

を使用して Delta Lake テーブル AWS Glue をクロールする方法についての情報を追加しました。詳細は、「[Delta Lake データストアの設定オプションを指定する方法](#)」を参照してください。

2022 年 2 月 24 日

### [AWS Glue ジョブインサイトの サポート](#)

AWS Glue ジョブインサイトを使用してジョブのデバッグと最適化を簡素化する情報を追加しました AWS Glue 。詳細については、「[AWS Glue ジョブインサイトを使用したモニタリング](#)」を参照してください。

2022 年 2 月 8 日

### [VPC エンドポイントを使用した Amazon S3 backedデータ カタログテーブルのクロール のサポート](#)

セキュリティ、監査、またはコントロールのために、Amazon Virtual Private Cloud 環境 (Amazon VPC) でのみアクセスするように Amazon S3 backed データカタログテーブルを設定する方法に関する情報を追加しました。詳細については、「[Crawling an Amazon S3 Data Store or Amazon S3 backedData Catalog tables using a VPC Endpoint](#)」を参照してください。

2022 年 2 月 3 日

## [Lake Formation ガバメント テーブルのSupport](#)

ACID トランザクション、自動データ圧縮、および時間移動クエリを AWS Glue サポートする Lake Formation 管理テーブルのサポートについての情報を追加しました。詳細については、「[AWS Glue API](#) および [AWS Lake Formation デベロッパーガイド](#)」を参照してください。

2021 年 11 月 30 日

## [インタラクティブセッション とノートブックに新しい AWS マネージドポリシーが追加 されました](#)

IAM の新しい マネージドポリシーでは、 [をインタラクティブセッションとノートブック AWS Glue で使用するためのセキュリティを強化しました。](#)詳細については、「[AWS GlueのAWS マネージドポリシー](#)」を参照してください。

2021 年 11 月 30 日

## [Glue スキーマレジストリがス トリーミングジョブでサポー トされるようになりました](#)

Glue スキーマレジストリの一部であるテーブルにアクセスするストリーミングジョブを作成できます。詳細については、「[AWS Glue Schema Registry とストリーミング ETL ジョブを追加 AWS Glue](#)」を参照してください。

2021 年 11 月 15 日

## [新しい機械学習機能のサポート](#)

インクリメンタルマッチングやマッチスコアリングなど、機械学習変換の新機能に関する情報を追加しました。詳細については、「[インクリメンタルマッチの検索](#)」と「[マッチの信頼度スコアを使用した一致の質の推定](#)」を参照してください。

2021 年 10 月 31 日

## [\(プライベートプレビュー\) AWS Glue フレックスジョブをサポート](#)

柔軟な実行クラスを持つ AWS Glue Spark ジョブの設定に関する情報が追加されました。これは、開始時間と完了時間が異なる時間的な制約のあるジョブに適しています。詳細については、「[AWS Glue でジョブを追加する](#)」を参照してください。

2021 年 10 月 29 日

## [Amazon S3 イベント通知を使用したクロールの高速化をサポート](#)

Amazon S3 イベント通知を使用したクロールの高速化に関する追加情報。詳細については、「[Amazon S3 イベント通知を使用したクロールの高速化](#)」を参照してください。

2021 年 10 月 15 日

### [アクセスコントロールおよび VPC に関連するセキュリティ設定オプションの追加](#)

AWS Glue および VPC の設定で新しいアクセス制御アクセス許可を設定する方法についての情報を追加しました。詳細については、[AWS 「のタグ AWS Glue」](#)、[「条件キーまたはコンテキストキーを使用して設定を制御するアイデンティティベースのポリシー \(IAM ポリシー\)」](#)、および[「VPC を通過するすべての呼び出しの設定」](#)を参照してください。 [AWS](#)

2021 年 10 月 13 日

### [VPC エンドポイントポリシーのサポートが追加されました](#)

AWS Glue で Virtual Private Cloud (VPC) エンドポイントポリシーのサポートに関する情報を追加しました。詳細については、[「AWS Glue とインターフェイス VPC エンドポイント \(AWS PrivateLink\)」](#)を参照してください。

2021 年 10 月 11 日

### [Glue Studio が中国で利用可能になりました](#)

AWS Glue Studio が北京および寧夏リージョンで利用可能になりました。

2021 年 10 月 11 日

### [AWS Glue Studio が、インタラクティブなジョブ編集におけるノートブックのオーサリング機能をリリースしました。](#)

ノートブックはコードの記述と実行、結果の可視化、およびインサイトの共有に役立ちます。通常データサイエンティストは、ノートブックを使用して実験やデータ探索のタスクを行います。詳細については、[「ノートブックの使用」](#)を参照してください。

2021 年 10 月 1 日

[ストリーミングソースへの直接アクセスが可能になりました](#)

ビジュアルエディタで ETL ジョブにデータソースを追加する場合、データカタログデータベースとテーブルを使用せずに、データストリームにアクセスするための情報を指定できます。

2021 年 9 月 30 日

[AWS Glue バージョンのサポートポリシーのドキュメント化](#)

AWS Glue バージョンサポートポリシーと、特定の期間終了フェーズ AWS Glue バージョンに関する情報を追加しました。詳細については、「[AWS Glue バージョンサポートのポリシー](#)」を参照してください。

2021 年 9 月 24 日

[カスタムコネクタをデータプレビューで使用できるようになりました](#)

カスタムコネクタを使用してデータソースノードを編集する場合、[データプレビュー] タブを選択してデータセットをプレビューできます。詳細については、「[カスタムコネクタ](#)」を参照してください。

2021 年 9 月 24 日

### [AWS Glue インタラクティブセッションのサポート \(プライベートプレビュー\)](#)

(プライベートプレビュー)  
AWS Glue インタラクティブセッションを使用して、任意の Jupyter Notebook からクラウドで Spark ワークロードを実行する方法に関する情報を追加しました。インタラクティブセッションは、AWS Glue 2.0 以降を使用する場合に AWS Glue 抽出、変換、ロード (ETL) コードを開発するために推奨される方法です。詳細については、「[Jupyter Notebook の AWS Glue インタラクティブセッションのセットアップと実行](#)」を参照してください。

2021 年 8 月 24 日

### [ブループリントからのワークフロー作成のサポート \(GA\)](#)

ブループリントでの一般的な抽出、変換、ロード (ETL) ユースケースのコーディングと、その後のブループリントからのワークフローの作成に関する情報を追加しました。データアナリストが複雑な ETL プロセスを簡単に作成して実行できるようにしています。詳細については、「[Performing Complex ETL Activities Using blueprints and Workflows in AWS Glue](#)」を参照してください。

2021 年 8 月 23 日

## [AWS Glue バージョン 3.0 のサポート。](#)

AWS Glue バージョン 3.0 のサポートについての情報を追加しました。Apache Spark ETL ジョブを実行するための Apache Spark 3.0 エンジンのアップグレード、およびその他の最適化とアップグレードをサポートしていません。詳細については、[AWS Glue リリースノート](#)および「[Migrating AWS Glue jobs to AWS Glue version 3.0](#)」を参照してください。このリリースの特徴には、他に、AWS Glue シャッフルマネージャー、SIMD ベクトル化された CSV リーダー、カタログパーティションの述語などがあります。詳細については、「[AWS Glue Spark shuffle manager with Amazon S3](#)」、「[Format Options for ETL Inputs and Outputs in AWS Glue](#)」、および「[Server-side filtering using catalog partition predicates](#)」を参照してください。

2021 年 8 月 18 日

## [AWS GovCloud \(US\) Region](#)

AWS Glue Studio がで利用可能になりました AWS GovCloud (US) Region

2021 年 8 月 18 日

### [AWS Glue Studio で Python シェルのオーサリングが使用可能になりました](#)

作成する新しいジョブとして、Python シェルジョブが選択できるようになりました。詳細については、[ジョブの作成プロセスの開始](#) および [Editing Python shell jobs in AWS Glue Studio](#) を参照してください。

2021 年 8 月 13 日

### [Amazon EventBridge イベントでワークフローを開始するためのサポート](#)

イベント駆動型アーキテクチャで AWS Glue がイベントを使用する方法についての情報を追加しました。詳細については、「[Amazon EventBridge イベントでAWS Glueワークフローを開始する](#)」および「[ワークフローを開始したEventBridge イベントを表示する](#)」を参照してください。

2021 年 7 月 14 日

### [AWS Glue スキーマレジストリでサポートされるデータ形式として JSON を追加](#)

サポートされているデータ形式として JSON に関する情報を追加しました (AVRO に加えて)。詳細については、「[AWS Glue Schema Registry](#)」を参照してください。

2021 年 6 月 30 日

### [データカタログテーブルなしで AWS Glue ストリーミングジョブを作成](#)

[create\\_data\\_frame\\_from\\_options](#) Python 関数または [getSource](#) for Scala スクリプトは、データカタログテーブルを必要とせず、データストリームを直接参照するストリーミング ETL ジョブの作成をサポートします。

2021 年 6 月 15 日

### [AWS Glue 機械学習変換で AWS Key Management Service キーのサポートを開始](#)

コンソール、CLI、または AWS Glue APIs を使用して AWS Glue Machine Learning 変換を設定するときに、セキュリティ設定または AWS KMS キーを指定できます。詳細については、「[Using Data Encryption with Machine Learning Transforms](#)」および「[AWS Glue Machine Learning API](#)」を参照してください。

2021 年 6 月 15 日

### [AWSGlueConsoleFullAccess AWS 管理ポリシーの更新](#)

管理ポリシーのマイナー更新に関する情報を追加しました AWSGlueConsoleFull Access AWS。詳細については、[AWS Glue 「AWS 管理ポリシーの更新」](#)を参照してください。

2021 年 6 月 10 日

### [ジョブの作成および編集中にジョブのデータセットを表示する](#)

ジョブダイアグラム内のノードで、新規に追加された [Data preview] (データのプレビュー) タブを使用すると、そのノードによって処理されるデータのサンプルを表示できます。詳細については、「[Using data previews in the visual job editor](#)」を参照してください。

2021 年 6 月 7 日

### [クローラー出力のテーブルの場所を示す値の指定をサポート](#)

クローラーの出力設定時、テーブルの場所を示す値の指定に関する情報を追加しました。詳細については、「[How to specify the table location](#)」を参照してください。

2021 年 6 月 4 日

[Amazon S3 データストアのクローリング時、データセット内のファイルのサンプルのクローリングをサポート](#)

Amazon S3 をクローリングするときにファイルのサンプルをクローリングする方法に関する情報を追加しました。詳細については、「[クローラーのプロパティ](#)」を参照してください。

2021 年 5 月 10 日

[AWS Glue 最適化された parquet ライターのサポート](#)

にAWS Glue最適化された parquet ライターを使用して、parquet分類でテーブルを作成または更新 DynamicFrames する方法について説明します。詳細については、「[Creating Tables, Updating Schema, and Adding New Partitions in theデータカタログfrom AWS Glue ETL Jobs](#)」および「[Format Options for ETL Inputs and Outputs in AWS Glue](#)」を参照してください。

2021 年 5 月 4 日

## [Kafka クライアント認証パスワードのサポート](#)

Apache Kafka ストリームプロデューサーを使用した AWS Glue でサポートされる SSL クライアント証明書認証での ETL ジョブのストリーミング方法に関する情報を追加しました。Apache Kafka クラスターへの AWS Glue 接続を定義する際に、カスタム証明書を指定できるようになりました。これは、認証時に AWS Glue で使用されます。詳細については、「[AWS Glue Connection Properties](#)」および「[Connection API](#)」を参照してください。

2021 年 4 月 28 日

## [ストリーミング ETL ジョブで、別のアカウントの Amazon Kinesis Data Streams からのデータの使いをサポート](#)

ストリーミング ETL ジョブの作成による別のアカウントの Amazon Kinesis Data Streams からのデータの使いに関する情報を追加しました。詳細については、「[AWS Glue でストリーミング ETL ジョブを追加する](#)」を参照してください。

2021 年 3 月 30 日

## [SQL 変換が使用可能になりました](#)

SQL の変換ノードを使用して、SQL クエリの形式で独自の変換を記述できます。詳細については、「[Using a SQL query to transform data](#)」を参照してください。

2021 年 3 月 23 日

### [ブループリントからのワークフロー作成のサポート \(公開レビュー\)](#)

(公開レビュー) ブループリントでの一般的な抽出、変換、ロード (ETL) のユースケースのコーディングと、その後のブループリントからのワークフローの作成に関する情報を追加しました。データアナリストが複雑な ETL プロセスを簡単に作成して実行できるようにしています。詳細については、「[Performing Complex ETL Activities Using blueprints and Workflows in AWS Glue](#)」を参照してください。

2021 年 3 月 22 日

### [データターゲットでコネクタが使用可能になりました](#)

データターゲットにカスタムまたは AWS Marketplace コネクタを使用することがサポートされるようになりました。詳細については、「[Authoring jobs with custom connectors](#)」を参照してください。

2021 年 3 月 15 日

### [AWS Glue 機械学習変換の列重要度メトリクスのサポート](#)

AWS Glue 機械学習変換使用時の列重要度メトリクスの表示に関する情報が追加されました。詳細については、「[Working with Machine Learning Transforms on the AWS Glue Console](#)」を参照してください。

2021 年 2 月 5 日

### [AWS Glue Studio でジョブのスケジュールリングが使用可能になりました](#)

AWS Glue Studio で時間ベースのスケジュールをジョブ実行用に定義できます。コンソールを使用して基本的なスケジュールを作成したり、Unix ライクの [cron](#) 構文を使用して、より複雑なスケジュールを定義したりできます。詳細については、「[Schedule job runs](#)」を参照してください。

2020 年 12 月 21 日

### [AWS Glue カスタムコネクタがリリースされました](#)

AWS Glue の Custom Connectors を使用すると、AWS Marketplace内でコネクタを検索およびサブスクライブできます。また、Apache Spark データソース、Athena フェデレーティッドクエリ、および JDBC API 用に構築されたコネクタをプラグインするための、AWS Glue Spark ランタイムインターフェイスもリリースしました。詳細については、「[Using connectors and connections with AWS Glue Studio](#)」を参照してください。

2020 年 12 月 21 日

### [AWS Glue バージョン 2.0でストリーミング ETL ジョブをサポート](#)

Glue バージョン 2.0 でのストリーミング ETL ジョブの実行に関するサポートに関する情報を追加しました。詳細については、「[AWS Glue でストリーミング ETL ジョブを追加する](#)」を参照してください。

2020 年 12 月 18 日

### [実行に上限を設定したワークロードのパーティション化のサポート](#)

ワークロードのパーティション化を有効にして、データセットサイズまたは ETL ジョブ実行で処理されるファイル数の上限を設定する方法に関する情報を追加しました。詳細については、「[Workload Partitioning with Bounded Execution](#)」を参照してください。

2020 年 11 月 23 日

### [拡張パーティション管理のサポート](#)

新しい API を使用して、既存のテーブルに対してパーティションインデックスを追加または削除する方法に関する情報を追加しました。詳細については、「[Working with Partition Indexes](#)」を参照してください。

2020 年 11 月 23 日

### [AWS Glue スキーマレジストリのサポート](#)

AWS Glue スキーマレジストリを使用して、スキーマを一元的に検出し、コントロールし、進化させる方法に関する情報を追加しました。詳細については、「[AWS Glue Schema Registry](#)」を参照してください。

2020 年 11 月 19 日

### [ストリーミング ETL ジョブでの grok 入力フォーマットのサポート](#)

ログファイルなどのストリーミングソースへの Grok パターンの適用に関する情報が追加されました。詳細については、「[Applying Grok Patterns to Streaming Sources](#)」を参照してください。

2020 年 11 月 17 日

### [AWS Glue コンソールでのワークフローへのタグの追加をサポート](#)

AWS Glue コンソールを使用してワークフローを作成するときのタグの追加に関する情報を追加しました。詳細については、「[Creating and Building Out a Workflow Using the AWS Glue Console](#)」を参照してください。

2020 年 10 月 27 日

### [増分クローラー実行のサポート](#)

前回の実行以降に追加された Amazon S3 フォルダのみをクローラーする増分クローラー実行のサポートに関する情報を追加しました。詳細については、「[Incremental Crawls](#)」を参照してください。

2020 年 10 月 21 日

### [ストリーミング ETL データソースのスキーマ検出のサポート。Avro ストリーミング ETL データソースとセルフマネージド型の kafka のサポート](#)

AWS Glue でのストリーミング抽出、変換、ロード (ETL) ジョブで、受信レコードのスキーマを自動的に検出し、レコードごとにスキーマの変更を処理できるようになりました。自己管理型 Kafka データソースがサポートされるようになりました。ストリーミング ETL ジョブで、データソースの Avro 形式がサポートされるようになりました。詳細については、「[Streaming ETL in AWS Glue](#)」、「[Defining Job Properties for a Streaming ETL Job](#)」、および「[Notes and Restrictions for Avro Streaming Sources](#)」を参照してください。

2020 年 10 月 7 日

### [MongoDB および DocumentDB データソースのクローリングのサポート](#)

MongoDB および Amazon DocumentDB (MongoDB 互換) データソースをクローリングするためのサポートに関する情報を追加しました。詳細については、「[Defining Crawlers](#)」を参照してください。

2020 年 10 月 5 日

### [FIPS コンプライアンスのサポート](#)

AWS Glue を使用してデータにアクセスするときに FIPS 140-2 検証済みの暗号化モジュールを必要とするお客様のための FIPS エンドポイントに関する情報を追加しました。詳細については、「[FIPS Compliance](#)」を参照してください。

2020 年 9 月 23 日

### [AWS Glue Studio Glue Studio での、ジョブを作成およびモニタリングするための使いやすいビジュアルインターフェイスの提供](#)

シンプルなグラフベースのインターフェイスを使用して、データを移動および変換するジョブを作成し、AWS Glue で実行できるようになりました。その後、AWS Glue Studio Glue Studio のジョブ実行ダッシュボードを使用して ETL の実行をモニタリングし、ジョブが意図したとおりに動作していることを確認できます。詳細については、[AWS Glue Studio ユーザーガイド](#) を参照してください。

2020 年 9 月 23 日

### [クエリのパフォーマンスを向上させるためのテーブルインデックス作成のサポート](#)

テーブルからパーティションのサブセットを取得できるようにするテーブルインデックスの作成に関する情報が追加されました。詳細については、「[Working with Partition Indexes](#)」を参照してください。

2020 年 9 月 9 日

### [Apache Spark ETL ジョブを実行する際のスタートアップ時間の削減を AWS Glue バージョン 2.0 でサポートしました。](#)

AWS Glue バージョン 2.0 のサポートについての情報を追加しました。このバージョンでは、Apache Spark ETL ジョブを実行するインフラストラクチャがアップグレードされ、スタートアップ時間の短縮、ログ記録の変更、ジョブレベルでの追加の Python モジュール指定のサポートが行われます。詳細については、「[AWS Glue Release Notes](#)」および「[Running Spark ETL Jobs with Reduced Startup Times](#)」を参照してください。

2020 年 8 月 10 日

### [同時ワークフロー実行数の制限のサポート。](#)

特定のワークフローの同時ワークフロー実行数を制限する方法に関する情報を追加しました。詳細については、「[Creating and Building Out a Workflow Using the AWS Glue Console](#)」を参照してください。

2020 年 8 月 10 日

### [VPC エンドポイントを使用した Amazon S3 データストアのクロールのサポート](#)

セキュリティ、監査、またはコントロールのために、Amazon Virtual Private Cloud 環境 (Amazon VPC) でのみアクセスするように Amazon S3 データストアを設定する方法に関する情報を追加しました。詳細については、「[Crawling an Amazon S3 Data Store using a VPC Endpoint](#)」を参照してください。

2020 年 8 月 7 日

### [ワークフロー実行再開のサポート](#)

1 つ以上のノード (ジョブまたはクローラー) が正常に完了しなかったために部分的にしか完了しなかったワークフロー実行を再開する方法に関する情報を追加しました。詳細については、「[Repairing and Resuming a Workflow Run](#)」を参照してください。

2020 年 7 月 27 日

### [kafka 接続でのプライベート CA 証明書の有効化をサポート。AWS Glue](#)

AWS Glue の Kafka 接続のプライベート CA 証明書を有効にすることをサポートする新しい接続オプションに関する情報を追加しました。詳細については、「[Connection Types and Options for ETL in AWS Glue](#)」および「[Special Parameters Used by AWS Glue](#)」を参照してください。

2020 年 7 月 20 日

### [別のアカウント内の DynamoDB データ読み取りの サポート](#)

別の AWS アカウントの DynamoDB テーブルからのデータの読み取りの AWS Glue サポートに関する情報を追加しました。詳細については、「[Reading from DynamoDB Data in Another Account](#)」を参照してください。

2020 年 7 月 17 日

### [AWS Glue バージョン 1.0 以 降での DynamoDB ライター接 続のサポート](#)

DynamoDB ライター、および DynamoDB の読み書き用の新しい接続オプションまたは更新された接続オプションのサポートに関する情報を追加しました。詳細については、「[Connection Types and Options for ETL in AWS Glue](#)」を参照してください。

2020 年 7 月 17 日

### [リソースリンクのサポート AWS Glue および and Lake Formation を用いたアカウント 横断アクセス制御のサポート](#)

リソースリンクと呼ばれる新しいデータカタログオブジェクトに関するコンテンツ、および AWS Glue と AWS Lake Formation の両方でアカウント間のデータカタログリソース共有を管理する方法に関するコンテンツを追加しました。詳細については、「[Granting Cross-Account Access](#)」および「[Table Resource Links](#)」を参照してください。

2020 年 7 月 7 日

### [DynamoDB データストアをクローリングするときのレコードのサンプリングのサポート](#)

DynamoDB データストアのクローリング時に設定できる新しいプロパティに関する情報を追加しました。詳細については、「[クローラーのプロパティ](#)」を参照してください。

2020 年 6 月 12 日

### [ワークフロー実行停止のサポート。](#)

特定のワークフローのワークフロー実行を停止する方法に関する情報を追加しました。詳細については、「[ワークフロー実行の停止](#)」を参照してください。

2020 年 5 月 14 日

### [Spark ストリーミング ETL ジョブのサポート](#)

ストリーミングデータソースを使用した抽出/変換/ロード (ETL、Extract/Transform/Load) ジョブの作成に関する情報を追加しました。詳細については、「[AWS Glue でストリーミング ETL ジョブを追加する](#)」を参照してください。

2020 年 4 月 27 日

### [ETL ジョブの実行後のテーブルの作成、スキーマの更新、データカタログでの新規パーティション追加のサポート](#)

テーブルの作成、スキーマの更新、およびデータカタログでの ETL ジョブの結果を確認するための新しいパーティションの追加を有効にする方法についての情報が追加されました。詳細については、「[Creating Tables, Updating Schema, and Adding New Partitions in the データカタログ from AWS Glue ETL Jobs](#)」を参照してください。

2020 年 4 月 2 日

[AWS Glue で、ETL 入力および出力として Apache Avro データ形式のバージョン指定のサポート](#)

Apache Avro データ形式のバージョンを AWS Glue の ETL 入力および出力として指定する方法についての情報を追加しました。デフォルトバージョン 1.7。version フォーマットオプションを使用して、Avro バージョン 1.8 を指定して、論理読み取り / 書き込みを有効にすることができます。詳細については、「[Format Options for ETL Inputs and Outputs in AWS Glue](#)」を参照してください。

2020 年 3 月 31 日

[Parquet データを Amazon S3 に書き込むための EMRFS S3 最適化コミッターのサポート](#)

AWS Glue ジョブの作成または更新時に、EMRFS S3 向けに最適化されたコミッターが Parquet データを Amazon S3 に書き込むための新しいフラグを設定する方法に関する情報を追加しました。詳細については、「[AWS Glue で使用される特別なパラメータ](#)」を参照してください。

2020 年 3 月 30 日

### [リソースタグによって管理される AWS リソースとしての機械学習変換のサポート](#)

AWS リソースタグを使用して、の機械学習変換へのアクセスを管理および制御する方法に関する情報を追加しましたAWS Glue。のジョブ、トリガー、エンドポイント、クローラ、機械学習変換に AWS リソースタグを割り当てることができますAWS Glue。詳細については、「[AWS Tags in AWS Glue](#)」を参照してください。

2020 年 3 月 2 日

### [上書きできないジョブ引数のサポート](#)

トリガーまたはジョブの実行時に上書きできない特殊なジョブパラメータのサポートに関する情報を追加しました。詳細については、「[AWS Glue でジョブを追加する](#)」を参照してください。

2020 年 2 月 12 日

### [Amazon S3 でデータセットを操作するための新しい変換のサポート](#)

Amazon S3 でデータセットを操作する Apache Spark アプリケーションの新しい変換 (マージ、パーティ、遷移) と Amazon S3 ストレージクラスの除外に関する情報を追加しました。Python のこれらの変換のサポートの詳細については、[mergeDynamicFrame](#) 「」および[Amazon S3でのデータセットの使用](#)」を参照してください。Scala については、[mergeDynamicFrames](#) 「」および[AWS Glue 「Scala GlueContext APIs](#)」を参照してください。

2020 年 1 月 16 日

### [ETL ジョブからの新しいパーティション情報を使用したデータカタログ更新のサポート](#)

抽出、変換、ロード (ETL) スクリプトをコーディングして、新しいパーティション情報 AWS Glue Data Catalog で更新する方法に関する情報を追加しました。この機能を使用すると、新しいパーティションを表示するためにジョブの完了後にクローラーを再実行する必要がなくなります。詳細については、「[Updating the Data Catalog with New Partitions](#)」を参照してください。

2020 年 1 月 15 日

<a href="#">新しいチュートリアル: SageMaker ノートブックの使 用</a>	Amazon SageMaker Notebook を使用して ETL および機械 学習スクリプトを開発する方 法を示すチュートリアルを追 加しました。 <a href="#">「チュートリ アル: 開発エンドポイントで Amazon SageMaker ノート ブックを使用する」</a> を参照し てください。	2020 年 1 月 3 日
<a href="#">MongoDB および Amazon DocumentDB (MongoDB 互換) からの読み取りのサポート</a>	MongoDB と Amazon DocumentDB (MongoDB 互換) との読み書き用の新しい接続 タイプと接続オプションに関 する情報を追加しました。詳 細については、「 <a href="#">Connection Types and Options for ETL in AWS Glue</a> 」を参照してくださ い。	2019 年 12 月 17 日
<a href="#">さまざまな修正と説明</a>	全体にわたって修正と説明を 追加しました。「既知の問 題」の章からエントリを削除 しました。データカタログの 暗号化設定を指定し、セキ ュリティ設定を作成する ときに、AWS Glue は対称カス タマーマスターキー (CMK) のみをサポートすることを 知らせる警告が追加されま した。AWS Glue が Amazon DynamoDB への書き込みをサ ポートしていないことを示す 注意を追加しました。	2019 年 12 月 9 日

## [カスタム JDBC ドライバーのサポート](#)

MySQL バージョン 8 や Oracle Database バージョン 18 など、AWS Glue がネイティブでサポートしていない JDBC ドライバーを使用したデータソースおよびターゲットへの接続に関する情報を追加しました。詳細については、「[JDBC connectionType の値](#)」を参照してください。

2019 年 11 月 25 日

## [SageMaker ノートブックをさまざまな開発エンドポイントに接続するためのサポート](#)

SageMaker ノートブックをさまざまな開発エンドポイントに接続する方法に関する情報を追加しました。新しい開発エンドポイントに切り替えるための新しいコンソールアクションと、新しい SageMaker IAM ポリシーを記述するための更新。詳細については、「[AWS Glue コンソールでのノートブックの操作](#)」および「[Amazon Notebooks の IAM SageMaker ポリシーの作成](#)」を参照してください。

2019 年 11 月 21 日

## [機械学習変換での AWS Glue バージョンのサポート](#)

機械学習変換と互換性のある AWS Glue のバージョンを示すため、機械学習変換での AWS Glue バージョンの定義に関する情報を追加しました。詳細については、「[Working with Machine Learning Transforms on the AWS Glue Console](#)」を参照してください。

2019 年 11 月 21 日

## [ジョブブックマークの巻き戻しのサポート](#)

ジョブのブックマークを以前のジョブ実行に巻き戻すことに関する情報を追加しました。その結果、後続のジョブ実行ではブックマークされたジョブ実行からのデータだけが再処理されます。2つのブックマーク間でジョブを実行できる `job-bookmark-pause` オプションの2つの新しいサブオプションについて説明しました。詳細については、「[Tracking Processed Data Using Job Bookmarks](#)」および「[Special Parameters Used by AWS Glue](#)」を参照してください。

2019 年 10 月 22 日

### [データストアに接続するためのカスタム JDBC 証明書のサポート](#)

AWS Glue データソースまたはターゲットへの SSL 接続用のカスタム JDBC 証明書の AWS Glue サポートに関する情報を追加しました。詳細については、「[Working with Connections on the AWS Glue Console](#)」を参照してください。

2019 年 10 月 10 日

### [Python Wheel のサポート](#)

Python シェルジョブの依存関係として wheel ファイル (egg ファイルと併用) の AWS Glue サポートに関する情報を追加しました。詳細については、「[独自の Python ライブラリの提供](#)」を参照してください。

2019 年 9 月 26 日

### [AWS Glue での開発エンドポイントのバージョンニングのサポート](#)

開発エンドポイントでの Glue version の定義に関する情報を追加しました。Glue version により、AWS Glue がサポートする Apache Spark および Python のバージョンが決定されます。詳細については、「[開発エンドポイントの追加](#)」を参照してください。

2019 年 9 月 19 日

## [Spark UI を使用した AWS Glue モニタリングのサポート](#)

Apache Spark UI を使用して、AWS Glue ジョブシステムで実行されている AWS Glue ETL ジョブと AWS Glue 開発エンドポイントの Spark アプリケーションのモニタリングとデバッグに関する情報を追加しました。詳細については、「[Monitoring AWS Glue Using Spark UI](#)」を参照してください。

2019 年 9 月 19 日

## [パブリック AWS Glue ETL ライブラリを使用したローカル ETL スクリプト開発のサポートの強化](#)

AWS Glue ETL ライブラリのコンテンツを更新して、AWS Glue バージョン 1.0 がサポートされたことを反映しました。詳細については、「[Developing and Testing ETL Scripts Locally Using the AWS Glue ETL Library](#)」を参照してください。

2019 年 9 月 18 日

## [ジョブ実行時の Amazon S3 ストレージクラス除外のサポート](#)

Amazon S3 からファイルまたはパーティションを読み取る AWS Glue ETL ジョブを実行する際の Amazon S3 ストレージクラスの除外に関する情報を追加しました。詳細については、「[Excluding Amazon S3 Storage Classes](#)」を参照してください。

2019 年 8 月 29 日

## [パブリック AWS Glue ETL ライブラリを使用したローカル ETL スクリプト開発のサポート](#)

ネットワーク接続を必要とせずに Python および Scala ETL スクリプトをローカルで開発およびテストする方法に関する情報を追加しました。詳細については、「[Developing and Testing ETL Scripts Locally Using the AWS Glue ETL Library](#)」を参照してください。

2019 年 8 月 28 日

## [既知の問題](#)

AWS Glue の既知の問題に関する情報を追加しました。詳細については、「[AWS Glue の既知の問題](#)」を参照してください。

2019 年 8 月 28 日

## [AWS Glueでの機械学習変換のサポート](#)

カスタム変換を作成するために AWS Glue によって提供される機械学習機能に関する情報を追加しました。これらの変換は、ジョブの作成時に作成できません。機械学習変換の詳細については、「[AWS Glue の機械学習変換](#)」を参照してください。

2019 年 8 月 8 日

## [共有 Amazon Virtual Private Cloud のサポート](#)

共有 Amazon Virtual Private Cloud の AWS Glue サポートに関する情報を追加しました。詳細については、「[共有 Amazon VPC](#)」を参照してください。

2019 年 8 月 6 日

## [AWS Glue でのバージョンニングのサポート](#)

ジョブプロパティでの Glue version の定義に関する情報を追加しました。AWS Glue バージョンにより、AWS Glue がサポートする Apache Spark および Python のバージョンが決定されます。詳細については、「[AWS Glue でジョブを追加する](#)」を参照してください。

2019 年 7 月 24 日

## [開発エンドポイントの追加設定オプションのサポート](#)

メモリを大量に使用するワークロードがある開発エンドポイントの設定オプションに関する情報を追加しました。エグゼキュターあたりのメモリを増強した新しい 2 つの設定からも選択できます。詳細については、「[AWS Glue コンソールでの開発エンドポイントの使用](#)」を参照してください。

2019 年 7 月 24 日

## [ワークフローを使用した ETL \(抽出、変換、ロード\) アクティビティの実行のサポート](#)

ワークフローと呼ばれる新しい構成を使用して、複雑なマルチジョブの抽出、変換、ロード (ETL) アクティビティを設計し、これを AWS Glue で単一のエンティティとして実行および追跡する機能に関する情報を追加しました。詳細については、「[Performing Complex ETL Activities Using Workflows in AWS Glue](#)」を参照してください。

2019年6月20日

### [Python シェルジョブでの Python 3.6 のサポート](#)

Python シェルジョブでの Python 3.6 のサポートに関する情報を追加しました。Python 3.6 または Python 2.7 をジョブプロパティとして指定できます。詳細については、「[Adding Python Shell Jobs in AWS Glue](#)」を参照してください。

2019 年 6 月 5 日

### [Virtual Private Cloud \(VPC\) エンドポイントのサポート](#)

VPC のインターフェイスエンドポイントを介した AWS Glue への直接接続に関する情報を追加しました。VPC インターフェイスエンドポイントを使用すると、AWS ネットワーク内で VPC と AWS Glue 間の通信が完全かつ安全に実施されます。詳細については、「[Using AWS Glue with VPC Endpoints](#)」を参照してください。

2019 年 6 月 4 日

### [AWS Glue ジョブのリアルタイムの連続ログ記録のサポート。](#)

ドライバーログ、各エグゼキューターログ、Spark ジョブの進行状況バー CloudWatch など、のリアルタイム Apache Spark ジョブログの有効化と表示に関する情報を追加しました。詳細については、「[AWS Glue ジョブの連続ログ記録](#)」を参照してください。

2019 年 5 月 28 日

### [クローラーのソースとして既存のデータカタログテーブルをサポート](#)

既存のデータカタログテーブルのリストをクローラーのソースとして指定することに関する情報を追加しました。新しいデータが利用可能になると、クローラーは、テーブルスキーマの変更を検出し、テーブル定義を更新して、新しいパーティションを登録できます。詳細については、「[クローラーのプロパティ](#)」を参照してください。

2019 年 5 月 10 日

### [メモリ大量使用ジョブ向け追加設定オプションのサポート](#)

メモリを集中的に使用する作業での Apache Spark ジョブの設定オプションに関する情報を追加しました。エグゼキューターあたりのメモリを増強した新しい 2 つの設定からも選択できます。詳細については、「[AWS Glue でジョブを追加する](#)」を参照してください。

2019 年 4 月 5 日

### [CSV カスタム分類子のサポート](#)

カスタム CSV 分類子を使用して、さまざまな種類の CSV データのスキーマを推測することに関する情報を追加しました。詳細については、「[カスタム分類子の書き込み](#)」を参照してください。

2019 年 3 月 26 日

## [AWS リソースタグのサポート](#)

AWS リソースへのアクセスを管理および制御するのに役立つAWS Glueリソースタグの使用に関する情報を追加しました。のジョブ、トリガー、エンドポイント、クローラーに AWS リソースタグを割り当てることができます AWS Glue。詳細については、「[AWS Tags in AWS Glue](#)」を参照してください。

2019 年 3 月 20 日

## [Spark SQL ジョブ用データカタログのサポート](#)

を外部 Apache Hive メタストア AWS Glue Data Catalog として使用するようにAWS Glueジョブと開発エンドポイントを設定する方法に関する情報を追加しました。これにより、ジョブおよび開発エンドポイントは AWS Glue Data Catalogに格納されているテーブルに対して Apache Spark SQL クエリを直接実行できます。詳細については、「[AWS Glue Data Catalog Support for Spark SQL Jobs](#)」を参照してください。

2019 年 3 月 14 日

## [Python シェルジョブのサポート](#)

Python シェルジョブと新しいフィールド [最小キャパシティー] についての情報を追加しました。詳細については、「[Adding Python Shell Jobs in AWS Glue](#)」を参照してください。

2019 年 1 月 18 日

### [データベースおよびテーブル 変更通知のサポート](#)

データベース、テーブル、およびパーティション API 呼び出しへの変更対して生成されるイベントに関する情報を追加しました。これらの CloudWatch イベントに応答するように Events でアクションを設定できます。詳細については、[CloudWatch 「イベントAWS Glueによる自動化」](#)を参照してください。

2019 年 1 月 16 日

### [接続パスワード暗号化のサポート](#)

接続オブジェクトで使用されるパスワードの暗号化についての情報を追加しました。詳細については、「[接続パスワードの暗号化と復号](#)」を参照してください。

2018 年 12 月 11 日

### [リソースレベルのアクセス許可とリソースに基づくポリシーのサポート](#)

AWS Glue でリソースレベルのアクセス許可とリソースに基づくポリシーを使用することに関する情報を追加しました。詳細については、「[AWS Glue でのセキュリティ](#)」のトピックを参照してください。

2018 年 10 月 15 日

### [SageMaker ノートブックのサポート](#)

AWS Glue 開発エンドポイントでの SageMaker ノートブックの使用に関する情報を追加しました。詳細については、「[ノートブックの管理](#)」を参照してください。

2018 年 10 月 5 日

## [暗号化のサポート](#)

AWS Glue で暗号化を使用することに関する情報を追加しました。詳細については、「[Encryption at Rest](#)」、「[Encryption in Transit](#)」、および「[Setting Up Encryption in AWS Glue](#)」を参照してください。

2018 年 8 月 24 日

## [Apache Spark ジョブメトリクスのサポート](#)

ETL のデバッグとプロファイリングを強化するため、Apache Spark メトリクスの使用に関する情報を追加しました。読み取りおよび書き込みされたバイト数、ドライバーとエグゼキュターのメモリ使用量と CPU 負荷、エグゼキュター間のデータシャッフルなどのランタイムメトリクスを AWS Glue コンソールから簡単に追跡できます。詳細については、[CloudWatch 「メトリクスAWS Glueを使用したモニタリング」](#)、「[ジョブのモニタリングとデバッグ](#)」、「[AWS Glue 「コンソールでのジョブの操作」](#)」を参照してください。

2018 年 7 月 13 日

<a href="#">データソースとしての DynamoDB のサポート</a>	DynamoDB をクローलし、ETL ジョブのデータソースとして使用する方法に関する情報を追加しました。詳細については、「 <a href="#">クローラーを使用してテーブルを分類する</a> 」と「 <a href="#">接続パラメータ</a> 」を参照してください。	2018 年 7 月 10 日
<a href="#">ノートブックサーバー作成手順の更新</a>	開発エンドポイントに関連付けられた Amazon EC2 インスタンスでノートブックサーバーを作成する方法に関する情報を更新しました。詳細については、「 <a href="#">開発エンドポイントに関連付けられているノートブックサーバーを作成する</a> 」を参照してください。	2018 年 7 月 9 日
<a href="#">更新を RSS で今すぐ入手可能</a>	RSS フィードにサブスクライブすると、AWS Glue デベロッパーガイドの更新に関する通知を受け取れるようになりました。	2018 年 6 月 25 日
<a href="#">ジョブ遅延通知のサポート</a>	ジョブ実行時の遅延しきい値の設定に関する情報を追加しました。詳細については、「 <a href="#">AWS Glue でジョブを追加する</a> 」を参照してください。	2018 年 5 月 25 日
<a href="#">新しい列を追加するようにクローラーを設定する</a>	クローラーの新しい設定オプションに関する情報を追加しました MergeNewColumns。詳細については、「 <a href="#">クローラーの設定</a> 」を参照してください。	2018 年 5 月 7 日

[ジョブタイムアウトのサポート](#)

ジョブ実行時のタイムアウトしきい値の設定に関する情報を追加しました。詳細については、「[AWS Glue でジョブを追加する](#)」を参照してください。

2018 年 4 月 10 日

[Scala ETL スクリプトと追加の実行条件に基づくジョブのトリガーのサポート](#)

Scala を ETL プログラミング言語として使用することについての追加の情報。さらに、トリガー API は、(すべての条件に加えて) いずれかの条件が満たされたときの発生をサポートするようになりました。また、ジョブは、(「succeeded」ジョブ実行に加えて)「failed」または「stopped」のジョブ実行に基づいてトリガーすることができます。

2018 年 1 月 12 日

## 以前の更新

次の表に、2018 年 1 月以前の AWS Glue 開発者ガイドの各リリースにおける重要な変更点を示します。

変更	説明	日付
XML データソースと新しいクローラー設定オプションをサポート	XML データソースとパーティション変更の新しいクローラーオプションについての情報を追加しました。	2017 年 11 月 16 日
新しい変換、追加の Amazon RDS データベースエン	マップとフィルターの変換、Amazon RDS Microsoft SQL Server と Amazon RDS Oracle のサポート、および開発エン	2017 年 9 月 29 日

変更	説明	日付
ジンのサポート、 および開発エンド ポイントの機能強 化	ドポイントの新機能に関する情報を追加 しました。	
AWS Glue 初回リ リース	これは AWS Glue デベロッパーガイドの 最初のリリースです。	2017 年 8 月 14 日

# AWS 用語集

最新の AWS 用語については、「AWS の用語集 リファレンス」の [AWS 「用語集」](#) を参照してください。

翻訳は機械翻訳により提供されています。提供された翻訳内容と英語版の間で齟齬、不一致または矛盾がある場合、英語版が優先します。