



開発者ガイド

Deep Learning AMI



Deep Learning AMI: 開発者ガイド

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

Amazon の商標とトレードドレスは、Amazon 以外の製品またはサービスとの関連において、顧客に混乱を招いたり、Amazon の名誉または信用を毀損するような方法で使用することはできません。Amazon が所有しない他の商標はすべてそれぞれの所有者に帰属します。所有者は必ずしも Amazon との提携や関連があるわけではありません。また、Amazon の支援を受けているとはかぎりません。

Table of Contents

| | |
|---|----|
| とは AWS Deep Learning AMI | 1 |
| 本ガイドについて | 1 |
| 前提条件 | 1 |
| 使用例 | 1 |
| 機能 | 2 |
| あらかじめインストールされたフレームワーク | 2 |
| 事前インストールされた GPU ソフトウェア | 3 |
| モデル処理および可視化 | 3 |
| 開始方法 | 4 |
| DLAMI の開始方法 | 4 |
| DLAMI の選択 | 4 |
| CUDA のインストール環境およびフレームワークのバインド | 5 |
| 基本 | 6 |
| Conda | 6 |
| アーキテクチャ | 8 |
| OS | 8 |
| インスタンスの選択 | 9 |
| 料金 | 10 |
| 利用可能なリージョン | 10 |
| GPU | 11 |
| CPU | 12 |
| Inferentia | 12 |
| Trainium | 13 |
| フレームワークサポートポリシー | 14 |
| サポートされるフレームワーク | 14 |
| よくある質問 | 14 |
| どのフレームワークバージョンにセキュリティパッチが適用されますか? | 15 |
| 新しいフレームワークバージョンがリリースされると、AWS はどのようなイメージを公開 しますか? | 15 |
| SageMaker どの画像に新機能や機能が加わったのか? AWS | 16 |
| サポート対象フレームワークの表では、現在のバージョンはどのように定義されていま すか? | 16 |
| サポート対象フレームワークの表に記載されていないバージョンを実行している場合はどう なりますか? | 16 |

| | |
|---|----|
| DLAMI は以前のバージョンのをサポートしていますか? TensorFlow | 16 |
| サポートされるフレームワークバージョン用の最新のパッチ適用済みイメージはどこにあり ますか? | 16 |
| 新しいイメージはどのくらいの頻度でリリースされますか? | 17 |
| ワークロードの実行中にインスタンスにインプレースでパッチが適用されますか? | 17 |
| 新しいパッチが適用されたフレームワークバージョン、または更新されたフレームワーク バージョンが利用可能になった場合はどうなりますか? | 17 |
| フレームワークのバージョンを変更せずに依存関係は更新されますか? | 17 |
| 使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか? | 18 |
| アクティブにメンテナンスされなくなったフレームワークバージョンのイメージにはパッチ が適用されますか? | 19 |
| 古いフレームワークバージョンを使用するにはどうすればよいですか? | 19 |
| up-to-date フレームワークとそのバージョンにおけるサポートの変更に遅れずについてい くにはどうすればいいですか? | 20 |
| Anaconda リポジトリを使用するには商用ライセンスが必要ですか? | 20 |
| DLAMI を起動する | 21 |
| ステップ 1: DLAMI を起動する | 22 |
| DLAMI ID を取得する | 22 |
| Amazon EC2 コンソールから起動する | 23 |
| ステップ 2: DLAMI に接続する | 24 |
| ステップ 3: DLAMI をテストする | 25 |
| ステップ 4: DLAMI インスタンスを管理する | 25 |
| クリーンアップ | 26 |
| Jupyter のセットアップ | 26 |
| Jupyter を保護する | 27 |
| サーバーの起動 | 28 |
| クライアントの設定 | 28 |
| Jupyter ノートブックサーバーへのログイン | 30 |
| DLAMI の使用 | 33 |
| Conda DLAMI | 33 |
| Deep Learning AMI with Conda の概要 | 33 |
| DLAMI にログインする | 34 |
| TensorFlow 環境を開始する | 34 |
| PyTorch Python 3 環境に切り替える | 35 |
| 環境を削除する | 36 |
| Base DLAMI | 36 |

| | |
|--------------------------------------|-----|
| Deep Learning Base AMI の使用 | 36 |
| CUDA のバージョンの設定 | 37 |
| Jupyter ノートブック | 37 |
| インストールされたチュートリアルを操作する | 38 |
| Jupyter で環境を切り替える | 38 |
| チュートリアル | 39 |
| 10 分間チュートリアル | 39 |
| フレームワークのアクティブ化 | 40 |
| Elastic Fabric Adapter | 43 |
| GPU のモニタリングおよび最適化 | 58 |
| AWS 推論 | 68 |
| ARM64 DLAMI | 90 |
| 推論 | 93 |
| モデル提供 | 94 |
| DLAMI のアップグレード | 100 |
| DLAMI のアップグレード | 100 |
| ソフトウェアの更新 | 101 |
| リリース通知 | 102 |
| セキュリティ | 103 |
| データ保護 | 104 |
| Identity and Access Management | 105 |
| アイデンティティを使用した認証 | 105 |
| ポリシーを使用したアクセスの管理 | 108 |
| Amazon EMR での IAM | 111 |
| ログ記録とモニタリング | 111 |
| 使用状況の追跡 | 111 |
| コンプライアンス検証 | 112 |
| 耐障害性 | 112 |
| インフラストラクチャセキュリティ | 113 |
| DLAMI に関する重要な変更点 | 114 |
| よくある質問 | 114 |
| 違いは何か | 114 |
| なぜこの変更が必要なのか? | 115 |
| この変更の影響を受けるのはどの DLAM ですか? | 116 |
| これはあなたにとってどのような意味がありますか? | 116 |
| 新しい DLAMI をいつ使い始めるべきですか? | 117 |

| | |
|------------------------------------|---------|
| 新しい DLAMI では機能が低下することはありますか? | 117 |
| DLC についてはどうですか? | 117 |
| 関連情報 | 118 |
| フォーラム | 118 |
| ブログ | 118 |
| よくある質問 | 118 |
| DLAMI のリリースノート | 123 |
| | 123 |
| Base DLAMI | 123 |
| シングルフレーム DLAMI | 124 |
| マルチフレームワーク DLAMI | 125 |
| DLAMI の廃止に関する通知 | 126 |
| ドキュメント履歴 | 128 |
| AWS 用語集 | 132 |
| | cxxxiii |

とは AWS Deep Learning AMI

AWS Deep Learning AMI のユーザーガイドへようこそ。

AWS Deep Learning AMI (DLAMI) は、クラウドでの深層学習のためのワンストップソリューションです。このカスタマイズされたマシンイメージは、小さな CPU 専用インスタンスから最新の高機能マルチ GPU インスタンスまで、さまざまなインスタンスタイプでほとんどの Amazon EC2 リージョンで利用できます。これは、[NVIDIA CUDA](#)、[NVIDIA cuDNN](#)、および最も人気が高い深層学習フレームワークの最新リリースを使用するように事前に設定されて提供されます。

本ガイドについて

このガイドでは、DLAMI を起動して使用方法について説明します。また、トレーニングと推論の両方でのディープラーニングの一般的なユースケースについてもいくつか説明します。その他、目的に応じた適切な AMI の選択や、好まれるインスタンスの種類についても取り上げます。DLAMI には、各フレームワークのチュートリアルがいくつか付属しています。また、分散トレーニング、デバッグ、Inferentia と AWS Trainium AWS の使用、その他の主要な概念に関するチュートリアルもあります。ブラウザでチュートリアルを実行するように Jupyter を設定する手順を確認できます。

前提条件

DLAMI を正常に実行するには、コマンドラインツールと基本的な Python についての知識が必要です。各フレームワークの使用方法に関するチュートリアルがフレームワーク自体に用意されていますが、このガイドでは、各フレームワークの有効化の方法を示すとともに、作業を開始するためのチュートリアルを探ることができます。

DLAMI の使用例

深層学習に関する学習: DLAMI は、機械学習および深層学習フレームワークの習得や教育にうってつけです。これにより、各フレームワークのインストール環境のトラブルシューティングや、同じコンピュータ上でフレームワークを協調させる作業に悩まなくて済むようになります。DLAMI には Jupyter ノートブックが付属しています。そのため、機械学習や深層学習が初めての方でも、フレームワークが提供するチュートリアルを簡単に実行することができます。

アプリ開発: アプリの開発を担当しており、深層学習によりアプリで AI の最新テクノロジーを利用することに興味をお持ちの方にとって、DLAMI は、申し分のない試験台になります。各フレームワー

クには、ディープラーニングを開始する方法についてのチュートリアルが付属しています。これらのチュートリアルの多くには Model Zoo があり、ニューラルネットワークを自分で作成したり、モデルトレーニングを実行したりせずに、簡単にディープラーニングを試すことができます。また、画像検出アプリケーションをわずか数分で作成する方法や、chatbot 用の音声認識アプリを作成する方法が例として示されています。

機械学習およびデータ分析: データサイエンティストの方、またはディープラーニングでデータを処理することに関心をお持ちの方は、フレームワークの多くで R と Spark がサポートされていることに気付くでしょう。スケーラブルなパーソナライズ用データ処理システムと予測システムを構築するまで、シンプルな回帰を実行する方法についてのチュートリアルが用意されています。

研究: 研究者が新しいフレームワークを試したり、新しいモデルをテストしたり、新しいモデルをトレーニングしたりする場合、DLAMI とのスケールアップ AWS 機能は、複数のトレーニングノードの面倒なインストールと管理の負担を軽減できます。

Note

より多くの GPU (最大 8 個) を持つ大きなインスタンスにインスタンスタイプをアップグレードすることを最初に選択する場合もあるかもしれませんが、DLAMI インスタンスのクラスターを作成して、水平方向にスケールアップすることも可能です。クラスター構築の詳細については、[関連情報](#) を参照してください。

DLAMI の機能

あらかじめインストールされたフレームワーク

現在、DLAMI には、主に 2 つのタイプがあります。また、それぞれについて、オペレーティングシステム (OS) とソフトウェアのバージョンに関連する他のバリエーションがあります。

- [Deep Learning AMI with Conda](#) - conda パッケージ、および独立した Python 環境を使用して個別にインストールされるフレームワーク
- [Deep Learning Base AMI](#) - フレームワークはインストールされず、[NVIDIA CUDA](#) と他の依存関係のみがインストールされます

Deep Learning AMI with Conda は、conda 環境を使用して各フレームワークを分離します。そのため、各フレームワークを自由に切り替えることができ、依存関係の競合を心配する必要はありません。

これは、Deep Learning AMI with Conda でサポートされているフレームワークの全リストです。

- PyTorch
- TensorFlow 2

Note

では、MXnet、CNTK、Caffe、Caffe2、Theano、Chainer、または Keras はサポートされなくなりました AWS Deep Learning AMI。

事前インストールされた GPU ソフトウェア

CPU のみのインスタンスを使用する場合でも、DLAMI では、[NVIDIA CUDA](#) と [NVIDIA cuDNN](#) が使用されます。インストールされたソフトウェアは、インスタンスタイプに関係なく同じです。GPU 専用ツールは、GPU が 1 つ以上あるインスタンスでのみ機能することに注意してください。詳細は、詳細な「[DLAMI のインスタンスタイプを選択する](#)」を参照してください。

CUDA のインストールに関する詳細については、「[CUDA のインストール環境およびフレームワークのバインド](#)」を参照してください。

モデル処理および可視化

Deep Learning AMI with Conda には TensorFlow、用のモデルサーバーと TensorBoard、モデルの視覚化用の [gcpcli](#) がインストールされています。

- [TensorFlow サービス](#)

開始方法

DLAMI の開始方法

このガイドには、最適な DLAMI を選択する方法やユースケースと予算に適したインスタンスタイプを選択する方法についてのヒントのほか、興味深いカスタムセットアップを説明している[関連情報](#)が記載されています。

Amazon EC2 AWS を初めて使用する場合は、 から始めます[Deep Learning AMI with Conda](#)。Amazon EC2 や Amazon EMR、Amazon EFS、Amazon S3 AWS などの他のサービスに精通していて、分散トレーニングや推論が必要なプロジェクトにこれらのサービスを統合することに関心がある場合は、ユースケースに合っている[関連情報](#)かどうかを確認してください。

[DLAMI の選択](#) を参照し、アプリケーションに最適なインスタンスタイプの概念を理解することをお勧めします。

もう 1 つのオプションは、このクイックチュートリアルである「 [を起動する AWS Deep Learning AMI \(10 分 \)](#)」です。

次のステップ

[DLAMI の選択](#)

DLAMI の選択

さまざまな DLAMI オプションを提供しています。ユースケースに適切な DLAMI を選択できるように、開発されたハードウェアのタイプまたは機能別にイメージをグループ化しています。トップレベルのグループは次のとおりです。

- DLAMI タイプ: CUDA、Base、シングルフレームワーク、マルチフレームワーク (Conda DLAMI)
- コンピューティングアーキテクチャ: x86 ベースと Arm64-based [AWS Graviton](#)
- プロセッサタイプ: [GPU 対 CPU 対 Inferentia](#)
- SDK: [CUDA](#) と [AWS Neuron](#)
- OS: Amazon Linux、Ubuntu

このガイドの残りのトピックでは、より詳細な情報を説明します。

トピック

- [CUDA のインストール環境およびフレームワークのバインド](#)
- [Deep Learning Base AMI](#)
- [Deep Learning AMI with Conda](#)
- [DLAMI アーキテクチャオプション](#)
- [DLAMI のオペレーティングシステムのオプション](#)

次回の予定

[Deep Learning AMI with Conda](#)

CUDA のインストール環境およびフレームワークのバインド

深層学習は、非常に最先端のものですが、各フレームワークは、"安定した" バージョンを提供しています。これらの安定したバージョンは、CUDA または cuDNN の最新の実装および機能とともに動作しない場合があります。ユースケースと必要な機能に基づいて、フレームワークを選択できます。よくわからない場合は、最新の Deep Learning AMI with Conda を使用してください。CUDA を使用するすべてのフレームワークの公式pipバイナリがあり、各フレームワークでサポートされている最新バージョンを使用します。最新のバージョンが必要な場合で、深層学習環境をカスタマイズするには、Deep Learning Base AMI を使用してください。

「[安定性とリリース候補](#)」で当社のガイドを参照してください。

DLAMI with CUDA の選択

[Deep Learning Base AMI](#) には、使用可能なすべての CUDA バージョンシリーズがあります。

[Deep Learning AMI with Conda](#) には、使用可能なすべての CUDA バージョンシリーズがあります。

Note

MXNet、CNTK、Caffe、Caffe2、Theano、Chainer、または Keras Conda 環境は に含まれなくなりました AWS Deep Learning AMI。

特定のフレームワークのバージョン番号については、[DLAMI のリリースノート](#)を参照してください。

この DLAMI タイプを選択するか、次回の予定オプションを使ってさまざまな DLAMI の詳しい情報を確認してください。

いずれかのバージョンの CUDA を選択してから、そのバージョンを含む DLAMI の詳細なリストを付録で確認するか、次回の予定オプションで別の DLAMI の詳細を確認してください。

次回の予定

[Deep Learning Base AMI](#)

関連トピック

- CUDA のバージョンを切り替える方法については、「[Deep Learning Base AMI の使用](#)」チュートリアルを参照してください。

Deep Learning Base AMI

Deep Learning Base AMI は、深層学習のための真っ白なキャンバスのようなものです。特定のフレームワークをインストールするまでに必要なすべてが含まれており、CUDA のバージョンも選択できます。

Base DLAMI を選択する理由

最先端のディープラーニングプロジェクトへの参加を目指すプロジェクト協力者は、この AMI グループを有効活用できます。また、環境を展開する際に、最新の NVIDIA ソフトウェアがインストールされ動作していることを確信できるため、インストールするフレームワークとバージョンの選択に集中できます。

この DLAMI タイプを選択するか、次回の予定オプションを使ってさまざまな DLAMI の詳しい情報を確認してください。

次回の予定

[DLAMI with Conda](#)

関連トピック

- [Deep Learning Base AMI の使用](#)

Deep Learning AMI with Conda

Conda DLAMI は conda 仮想環境を使用し、マルチフレームワークまたは単一フレームワークの DLAMIs のいずれかで存在します。これらの環境は、インストールした異なるフレームワークの独立

性が維持され、フレームワーク間の切り替えが合理化されるように設定されます。これは、DLAMI が提供するすべてのフレームワークを学習したり、試したりするのに最適です。ほとんどのユーザーにとって、新しい Deep Learning AMI with Conda は申し分のないものとなります。

フレームワークからの最新バージョンで頻繁に更新され、最新の GPU ドライバとソフトウェアが導入されることとなります。これらは通常、AWS Deep Learning AMI ほとんどのドキュメントではと呼ばれます。これらの DLAMIs Ubuntu20.04、Amazon Linux 2 オペレーティングシステムをサポートしています。オペレーティングシステムのサポートは、アップストリーム OS からのサポートによって異なります。

安定性とリリース候補

Conda AMI は、各フレームワークの最新の正式リリース版の最適化バイナリを使用します。リリース候補および実験的機能は対象外です。この最適化は、C5、および C4 CPU インスタンスタイプにおけるトレーニングおよび推論を高速化する高速化テクノロジー (MKL DNN など) をサポートするフレームワークに依存します。また、バイナリは高度な Intel 指示セット (AVX、AVX-2、SSE4.1、SSE4.2 などを含む) をサポートするようにコンパイルされています。これによって、Intel CPU アーキテクチャ上のベクターおよび浮動小数点オペレーションが高速化されます。さらに、GPU インスタンスタイプでは、CUDA および cuDNN は最新の公式リリースがサポートするバージョンで更新されます。

Deep Learning AMI with Conda は、Amazon EC2 インスタンスにフレームワークの最初のアクティベーションから最適化されたフレームワークのバージョンを自動的にインストールします。詳細については、「[Deep Learning AMI with Conda の使用](#)」を参照してください。

カスタムまたは最適化されたビルドオプションを使用してソースからインストールする場合は、[Deep Learning Base AMI](#) の方が適しています。

Python 2 の廃止

Python オープンソースコミュニティは、Python 2 のサポートを 2020 年 1 月 1 日に正式に終了しました。TensorFlow および PyTorch コミュニティは、Python TensorFlow 2 をサポートする最後のリリースが 2.1 および PyTorch 1.4 リリースであることを発表しました。Python 2 Conda 環境を含む DLAMI (v26、v25 など) の以前のリリースは、引き続き利用できます。ただし、以前に公開された DLAMI のバージョンで Python 2 Conda 環境の更新を提供するのは、それらのバージョンのオープンソースコミュニティによって公開されたセキュリティ修正がある場合のみです。TensorFlow および PyTorch フレームワークの最新バージョンを使用する DLAMI リリースには、Python 2 Conda 環境は含まれていません。

CUDA サポート

特定の CUDA バージョン番号は [GPU DLAMIリリースノート](#) で確認できます。

次回の予定

[DLAMI アーキテクチャオプション](#)

関連トピック

- Deep Learning AMI with Conda の使用に関するチュートリアルについては、[Deep Learning AMI with Conda の使用](#) のチュートリアルを参照してください。

DLAMI アーキテクチャオプション

AWS Deep Learning AMIは、x86 ベースまたは Arm64-based [AWS Graviton2](#) アーキテクチャのいずれかで提供されます。

ARM64 GPU DLAMI の開始方法については、「」を参照してください[ARM64 DLAMI](#)。使用可能なインスタンスタイプの詳細については、[DLAMI のインスタンスタイプを選択する](#) を参照してください。

次回の予定

[DLAMI のオペレーティングシステムのオプション](#)

DLAMI のオペレーティングシステムのオプション

DLAMI は、以下のオペレーティングシステムで提供されています。

- Amazon Linux 2
- Ubuntu 20.04
- Ubuntu 22.04

古いバージョンのオペレーティングシステムは、廃止された DLAMI でも使用できます。DLAMI の廃止について詳しくは、「[DLAMI の廃止](#)」を参照してください。

DLAMI を選択する前に、必要なインスタンスタイプを評価し、AWS リージョンを特定してください。

次回の予定

[DLAMI のインスタンスタイプを選択する](#)

DLAMI のインスタンスタイプを選択する

一般的に、DLAMI のインスタンスタイプを選択するときは、次の点を考慮してください。

- 深層学習を初めて導入する場合は、GPU が 1 個搭載のインスタンスがニーズに合っている可能性があります。
- 予算重視の場合は、CPU のみのインスタンスを使用できます。
- 深層学習モデルの推論に高いパフォーマンスとコスト効率を最適化する場合は、Inferentia AWS チップでインスタンスを使用できます。
- Arm64-based CPU アーキテクチャを備えた高性能 GPU インスタンスをお探しの場合は、G5g インスタンスタイプを使用できます。
- 推論と予測用に事前トレーニング済みモデルを実行することに興味がある場合は、[Amazon Elastic Inference](#) を Amazon EC2 インスタンスにアタッチできます。Amazon Elastic Inference では、GPU の一部を利用したアクセラレーターにアクセスできます。
- 大容量の推論サービスの場合は、大容量のメモリを搭載した 1 つの CPU インスタンス、またはそのようなインスタンスのクラスターが、より適したソリューションになることがあります。
- 大量のデータまたは大きなバッチサイズを持つ大規模なモデルを使用している場合は、より多くのメモリを持つより大きなインスタンスが必要になります。モデルを GPU のクラスターに配布することもできます。バッチサイズを小さくした場合、メモリの少ないインスタンスを使用すると改善されることがあります。これは、精度とトレーニング速度に影響を与える可能性があります。
- 大規模で高レベルのノード間通信を必要とする NVIDIA Collective Communications Library (NCCL) を使用した機械学習アプリケーションの実行に関心がある場合は、[Elastic Fabric Adapter \(EFA\)](#) を使用できます。

インスタンスの詳細については、[EC2 インスタンスタイプ](#)を参照してください。

次のトピックでは、インスタンスタイプに関する考慮事項について説明します。

Important

ディープラーニング AMI には、ドライバ、ソフトウェアやツールキット (NVIDIA Corporation によって開発、所有あるいは提供) が含まれています。これらの NVIDIA ドライ

バ、ソフトウェア、ツールキットは、NVIDIA ハードウェアが含まれている Amazon EC2 インスタンスでのみ使用することにユーザーが同意するものとします。

トピック

- [DLAMI の価格設定](#)
- [DLAMI で利用可能なリージョン](#)
- [推奨 GPU インスタンス](#)
- [推奨 CPU インスタンス](#)
- [推奨 Inferentia インスタンス](#)
- [推奨 Trainium インスタンス](#)

DLAMI の価格設定

DLAMI に含まれている深層学習フレームワークは無料で、それぞれに独自のオープンソースライセンスが付与されています。DLAMI に含まれているソフトウェアは無料ですが、基盤となる Amazon EC2 インスタンスハードウェアの料金を支払う必要があります。

Amazon EC2 インスタンスタイプには無料とされているものがあります。そうした無料のインスタンスのいずれかで DLAMI を実行できます。これは、そのインスタンスの容量だけを使う場合は、DLAMI の使用が完全に無料という意味です。CPU コア数を増やす、ディスク容量を増やす、RAM を増やす、GPU を 1 つ以上使用するなど、より強力なインスタンスが必要な場合は、無料利用枠のインスタンスクラス以外のインスタンスが必要になります。

インスタンスの選択と料金の詳細については、[Amazon EC2 の料金](#)を参照してください。

DLAMI で利用可能なリージョン

各リージョンでサポートされるインスタンスタイプの範囲はそれぞれ異なり、多くの場合、インスタンスタイプのコストもリージョンごとにわずかに違います。DLAMI を利用できないリージョンはありますが、選択したリージョンに DLAMI をコピーすることは可能です。詳細については、[AMI のコピー](#)を参照してください。リージョンの選択リストを確認し、必ず自社または自社の顧客に近いリージョンを選択してください。複数の DLAMI を使用する予定があり、クラスターを作成する可能性がある場合は、クラスター内のすべてのノードに同じリージョンを使う必要があります。

リージョンの詳細については、[EC2 のリージョン](#)を参照してください。

次回の予定

[推奨 GPU インスタンス](#)

推奨 GPU インスタンス

GPU インスタンスは、深層学習の大半の目的に推奨されます。新しいモデルのトレーニングは CPU インスタンス上よりも GPU インスタンス上の方が高速に実行できます。複数の GPU インスタンスがある場合や、トレーニングを複数の GPU インスタンスに分散した場合は、ほぼ直線的な拡張性を得ることができます。

以下のインスタンスタイプで DLAMI がサポートされています。GPU インスタンスタイプのオプションとその使用方法の詳細については、[EC2 インスタンスタイプ](#)を参照し、高速コンピューティングを選択してください。

Note

モデルのサイズは、インスタンスを選択する際の要因となります。モデルがインスタンスの使用可能な RAM を超えている場合は、アプリケーション用に十分なメモリを持つ別のインスタンスタイプを選択します。

- [Amazon EC2 P3 インスタンス](#)には、NVIDIA Tesla V100 GPU が最大 8 個搭載されます。
- [Amazon EC2 P4 インスタンス](#)には、NVIDIA Tesla A100 GPU が最大 8 個搭載されます。
- [Amazon EC2 P5 インスタンス](#)には、最大 8 つの NVIDIA Tesla H100 GPUs。
- [Amazon EC2 G3 インスタンス](#)には、NVIDIA Tesla M60 GPU が最大 4 個搭載されます。
- [Amazon EC2 G4 インスタンス](#)には、NVIDIA T4 GPU が最大 4 個搭載されます。
- [Amazon EC2 G5 インスタンス](#)には、NVIDIA A10G GPU が最大 8 個搭載されます。
- [Amazon EC2 G6 インスタンス](#)には、最大 8 つの NVIDIA L4 GPUs。
- [Amazon EC2 G5g インスタンス](#)には Arm64-based [AWS Graviton2 プロセッサ](#) があります。

DLAMI インスタンスでは、GPU プロセスをモニタリングおよび最適化するためのツールが提供されています。GPU プロセスのモニタリングの詳細については、[GPU のモニタリングおよび最適化](#)を参照してください。

G5g インスタンスの操作に関する具体的なチュートリアルについては、[ARM64 DLAMI](#)を参照してください。

次回の予定

[推奨 CPU インスタンス](#)

推奨 CPU インスタンス

予算に制約がある場合、ディープラーニングを学習する目的の場合、または予測サービスの実行が唯一の目的である場合は、CPU のカテゴリに数多くの低価格な選択肢があります。一部のフレームワークでは、インテルの MKL DNN を利用しています。これにより、C5 でのトレーニングと推論が高速化されます (一部のリージョンでは使用できません)。CPU インスタンスタイプの詳細については、[EC2 インスタンスタイプ](#)を参照し、コンピューティング最適化を選択してください。

Note

モデルのサイズは、インスタンスを選択する際の要因となります。モデルがインスタンスの使用可能な RAM を超えている場合は、アプリケーション用に十分なメモリを持つ別のインスタンスタイプを選択します。

- [Amazon EC2 C5 インスタンス](#)には最大 72 個のインテル vCPU があります。C5 インスタンスは、科学モデリング、バッチ処理、分散分析、ハイパフォーマンスコンピューティング (HPC)、機械学習/深層学習による推論などに優れています。

次回の予定

[推奨 Inferentia インスタンス](#)

推奨 Inferentia インスタンス

AWS Inferentia インスタンスは、深層学習モデル推論ワークロードに高いパフォーマンスとコスト効率を提供するように設計されています。具体的には、Inf2 インスタンスタイプは AWS Inferentia チップと [AWS Neuron SDK](#) を使用します。これは、TensorFlow や などの一般的な機械学習フレームワークと統合されています PyTorch。

お客様は Inf2 インスタンスを使用して、検索、レコメンデーションエンジン、コンピュータビジョン、音声認識、自然言語処理、パーソナライゼーション、不正検出などの大規模な機械学習推論アプリケーションをクラウド内で低コストで実行できます。

Note

モデルのサイズは、インスタンスを選択する際の要因となります。モデルがインスタンスの使用可能な RAM を超えている場合は、アプリケーション用に十分なメモリを持つ別のインスタンスタイプを選択します。

- [Amazon EC2 Inf2 インスタンス](#)には、最大 16 個の AWS Inferentia チップと 100 Gbps のネットワークスループットがあります。

Inferentia DLAMI AWS の開始方法の詳細については、「」を参照してください [DLAMI を備えた AWS Inferentia チップ](#)。DLAMIs

次回の予定

[推奨 Trainium インスタンス](#)

推奨 Trainium インスタンス

AWS Trainium インスタンスは、深層学習モデル推論ワークロードに高いパフォーマンスとコスト効率を提供するように設計されています。具体的には、Trn1 インスタンスタイプは AWS Trainium チップと [AWS Neuron SDK](#) を使用します。これは、TensorFlow や などの一般的な機械学習フレームワークと統合されています PyTorch。

お客様は Trn1 インスタンスを使用して、検索、レコメンデーションエンジン、コンピュータビジョン、音声認識、自然言語処理、パーソナライゼーション、不正検出などの大規模な機械学習推論アプリケーションをクラウド内で低コストで実行できます。

Note

モデルのサイズは、インスタンスを選択する際の要因となります。モデルがインスタンスの使用可能な RAM を超えている場合は、アプリケーション用に十分なメモリを持つ別のインスタンスタイプを選択します。

- [Amazon EC2 Trn1 インスタンス](#)には、最大 16 AWS Trainium チップと 100 Gbps のネットワークスループットがあります。

フレームワークサポートポリシー

[AWS Deep Learning AMI \(DLAMI\)](#) は、深層学習ワークロードのイメージ構成を簡素化し、最新のフレームワーク、ハードウェア、ドライバー、ライブラリ、オペレーティングシステムで最適化されています。このページでは、DLAMI のフレームワークサポートポリシーについて詳しく説明します。使用可能な DLAMI のリストについては、「[DLAMI のリリースノート](#)」を参照してください。

サポートされるフレームワーク

以下の [AWS Deep Learning AMI フレームワークサポートポリシーの表](#) を参照して、現在サポートされているフレームワークとバージョンを確認してください。

「End of patch」セクションで、オリジンフレームワークのメンテナンスチームが現在サポートしている現在のバージョンを AWS がサポートする期間を確認してください。フレームワークとバージョンは、シングルフレームワーク DLAMI またはマルチフレームワーク DLAMI で利用できます。

Note

フレームワークバージョン x.y.z では、x はメジャーバージョン、y はマイナーバージョン、z はパッチバージョンを指します。たとえば、TensorFlow 2.6.5 の場合、メジャーバージョンは 2、マイナーバージョンは 6、パッチバージョンは 5 です。

特定のイメージの詳細については、リリースノートを参照してください。

- [シングルフレームワーク DLAMI リリースノート](#)
- [マルチフレームワーク DLAMI リリースノート](#)

よくある質問

- [どのフレームワークバージョンにセキュリティパッチが適用されますか？](#)
- [新しいフレームワークバージョンがリリースされると、AWS はどのようなイメージを公開しますか？](#)
- [SageMaker どの画像に新機能や機能が加わったのか？AWS](#)
- [サポート対象フレームワークの表では、現在のバージョンはどのように定義されていますか？](#)

- [サポート対象フレームワークの表に記載されていないバージョンを実行している場合はどうなりますか？](#)
- [DLAMI は以前のバージョンのをサポートしていますか？ TensorFlow](#)
- [サポートされるフレームワークバージョン用の最新のパッチ適用済みイメージはどこにありますか？](#)
- [新しいイメージはどのくらいの頻度でリリースされますか？](#)
- [ワークロードの実行中にインスタンスにインプレースでパッチが適用されますか？](#)
- [新しいパッチが適用されたフレームワークバージョン、または更新されたフレームワークバージョンが利用可能になった場合はどうなりますか？](#)
- [フレームワークのバージョンを変更せずに依存関係は更新されますか？](#)
- [使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか？](#)
- [アクティブにメンテナンスされなくなったフレームワークバージョンのイメージにはパッチが適用されますか？](#)
- [古いフレームワークバージョンを使用するにはどうすればよいですか？](#)
- [up-to-date フレームワークとそのバージョンにおけるサポートの変更に遅れずについていくにはどうすればいいですか？](#)
- [Anaconda リポジトリを使用するには商用ライセンスが必要ですか？](#)

どのフレームワークバージョンにセキュリティパッチが適用されますか？

[AWS Deep Learning AMI フレームワークサポートポリシーの表](#)で、「Supported」とラベルが付いているフレームワークバージョンにはセキュリティパッチが適用されます。

新しいフレームワークバージョンがリリースされると、AWS はどのようなイメージを公開しますか？

およびの新しいバージョンがリリースされた直後に、新しい DLAMI を公開します TensorFlow 。 PyTorch これには、メジャーバージョン、メジャーマイナーバージョン、フレームワークのバージョンが含まれます major-minor-patch 。 また、新しいバージョンのドライバーやライブラリが利用可能になった場合も、イメージが更新されます。 イメージメンテナンスの詳細については、「[使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか？](#)」を参照してください。

SageMakerどの画像に新機能や機能が加わったのか？AWS

新機能は通常、PyTorch および用の DLAMIs の最新バージョンでリリースされます。TensorFlow SageMaker AWS新機能や機能の詳細については、特定のイメージのリリースノートを参照してください。使用可能な DLAMI のリストについては、「[DLAMI のリリースノート](#)」を参照してください。イメージメンテナンスの詳細については、「[使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか?](#)」を参照してください。

サポート対象フレームワークの表では、現在のバージョンはどのように定義されていますか？

[AWS Deep Learning AMIフレームワークSupport ポリシー表の現在のバージョン](#)は、GitHubで利用できる最新のフレームワークバージョンを指します。AWS各最新リリースには、DLAMI のドライバー、ライブラリ、関連パッケージの更新が含まれています。イメージメンテナンスの詳細については、「[使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか?](#)」を参照してください。

サポート対象フレームワークの表に記載されていないバージョンを実行している場合はどうなりますか？

[AWS Deep Learning AMI フレームワークサポートポリシーの表](#)に記載されていないバージョンを実行している場合は、最新のドライバー、ライブラリ、関連パッケージがインストールされていない可能性があります。up-to-date 新しいバージョンについては、選択した最新の DLAMI を使用して、サポートされているフレームワークのいずれかにアップグレードすることをお勧めします。使用可能な DLAMI のリストについては、「[DLAMI のリリースノート](#)」を参照してください。

DLAMI は以前のバージョンのをサポートしていますか？ TensorFlow

いいえ。フレームワークSupport [ポリシー表に記載されているように](#)、[GitHub AWS Deep Learning AMI最初のリリースから365日後にリリースされた各フレームワークの最新のメジャーバージョンの最新のパッチバージョンをサポートします](#)。詳細については、[サポート対象フレームワークの表に記載されていないバージョンを実行している場合はどうなりますか?](#)を参照してください。

サポートされるフレームワークバージョン用の最新のパッチ適用済みイメージはどこにありますか？

最新のフレームワークバージョンの DLAMI を使用するには、[DLAMI ID](#) を取得し、その ID を使用して [EC2 コンソール](#)を使用して DLAMI を起動します。AWS Deep Learning AMI ID を取得す

るサンプル AWS CLI コマンドについては、[AWS Deep Learning AMI カタログ](#)の「Deep Learning Frameworks」セクションを参照してください。AWS CLI AMI ID クエリは、[シングルフレームワーク DLAMI リリースノート](#)にも含まれています。選択するフレームワークバージョンは、[AWS Deep Learning AMI フレームワークサポートポリシーの表](#)で「Supported」とラベルが付いている必要があります。

新しいイメージはどのくらいの頻度でリリースされますか？

最新のパッチバージョンを提供することは、AWS の最優先事項です。パッチを適用したイメージを、できるだけ早く定期的に作成しています。新たにパッチが適用されたフレームワークのバージョン (例: TensorFlow 2.9 ~ TensorFlow 2.9.1) と新しいマイナーリリースバージョン (例: TensorFlow 2.9 から TensorFlow 2.10 まで)、できるだけ早い機会に利用できるようにしてください。の既存のバージョンが新しいバージョンの CUDA TensorFlow とともにリリースされると、新しい CUDA バージョンをサポートするそのバージョン用の新しい DLAMI がリリースされます。

TensorFlow

ワークロードの実行中にインスタンスにインプレースでパッチが適用されますか？

いいえ。DLAMI のパッチ更新は「インプレース」更新ではありません。

新しい EC2 インスタンスをオンにし、ワークロードとスクリプトを移行してから、以前のインスタンスをオフにする必要があります。

新しいパッチが適用されたフレームワークバージョン、または更新されたフレームワークバージョンが利用可能になった場合はどうなりますか？

イメージのリリースノートページを定期的に確認してください。新しいパッチが適用されたフレームワーク、または更新されたフレームワークが利用可能になった場合は、そのフレームワークにアップグレードすることをお勧めします。使用可能な DLAMI のリストについては、「[DLAMI のリリースノート](#)」を参照してください。

フレームワークのバージョンを変更せずに依存関係は更新されますか？

AWS はフレームワークのバージョンを変更せずに依存関係を更新します。ただし、依存関係の更新によって互換性が失われる場合、AWS は別のバージョンでイメージを作成します。更新された依存関係の情報については、必ず「[DLAMI のリリースノート](#)」を確認してください。

使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか？

DLAMI イメージはイミュータブルです。一度作成されると変更されません。フレームワークバージョンの有効なサポートが終了する主な理由は 4 つあります。

- [フレームワークバージョン \(パッチ\) のアップグレード](#)
- [AWS セキュリティパッチ](#)
- [パッチ終了日 \(エー징アウト\)](#)
- [依存関係 end-of-support](#)

Note

バージョンパッチのアップグレードやセキュリティパッチは頻繁に行われるため、DLAMI のリリースノートページを頻繁に確認し、変更があった場合はアップグレードすることをお勧めします。

フレームワークバージョン (パッチ) のアップグレード

TensorFlow 2.7.0 に基づく DLAMI ワークロードがあり、バージョン 2.7.1 TensorFlow 以降をリリースしている場合は GitHub、2.7.1 で新しい DLAMI AWS をリリースします。TensorFlow 2.7.0 の以前のイメージは、2.7.1 の新しいイメージがリリースされるとアクティブにメンテナンスされなくなります。TensorFlow TensorFlow 2.7.0 の DLAMI には、これ以上のパッチは適用されません。その後、TensorFlow 2.7 の DLAMI リリースノートページが最新の情報で更新されます。マイナーパッチごとの個別のリリースノートページはありません。

パッチアップグレードにより作成された新しい DLAMI には、新しい [AMI ID](#) が割り当てられます。

AWS セキュリティパッチ

2.7.0 のイメージに基づくワークロードがあり、AWSセキュリティパッチを作成すると、TensorFlow 2.7.0 用の DLAMI の新しいバージョンがリリースされます。TensorFlow TensorFlow 2.7.0 の以前のバージョンのイメージはアクティブにメンテナンスされなくなりました。詳細については、「[ワークロードの実行中にインスタンスにインプレースでパッチが適用されますか?](#)」を参照してください。最新の DLAMI を検索する手順については、「[サポートされるフレームワークバージョン用の最新のパッチ適用済みイメージはどこにありますか?](#)」を参照してください。

パッチアップグレードにより作成された新しい DLAMI には、新しい [AMI ID](#) が割り当てられます。

パッチ終了日 (エージングアウト)

DLAMI は、リリース日から 365 日後にパッチの終了日を迎えます。GitHub

[マルチフレームワーク DLAMI](#) の場合、いずれかのフレームワークバージョンが更新されると、更新されたバージョンを含む新しい DLAMI が必要になります。古いフレームワークバージョンの DLAMI は、アクティブにメンテナンスされなくなります。

Important

フレームワークのメジャー更新がある場合は例外とします。たとえば、TensorFlow 1.15 が TensorFlow 2.0 にアップデートされた場合、GitHub リリース日から 2 年間、またはオリジンフレームワークのメンテナンスチームがサポートを終了してから 6 か月間 (どちらか早い方)、TensorFlow 1.15 の最新バージョンを引き続きサポートします。

依存関係 end-of-support

Python 3.6 を搭載した TensorFlow 2.7.0 DLAMI イメージ上でワークロードを実行していて、そのバージョンの Python が対象になっている場合 end-of-support、Python 3.6 に基づくすべての DLAMI イメージはアクティブにメンテナンスされなくなります。同様に、Ubuntu 16.04 のような OS バージョンが対象になっている場合 end-of-support、Ubuntu 16.04 に依存するすべての DLAMI イメージはアクティブにメンテナンスされなくなります。

アクティブにメンテナンスされなくなったフレームワークバージョンのイメージにはパッチが適用されますか？

いいえ。アクティブにメンテナンスされていないイメージには新しいリリースは適用されません。

古いフレームワークバージョンを使用するにはどうすればよいですか？

古いフレームワークバージョンの DLAMI を使用するには、[DLAMI ID](#) を取得し、その ID を使用して [EC2 コンソール](#) を使用して DLAMI を起動します。AMI ID を取得する AWS CLI コマンドについては、[AWS Deep Learning AMI カタログ](#) の「Deep Learning Frameworks」セクションを参照してください。AWSCLI AMI ID クエリは、[シングルフレームワーク DLAMI リリースノート](#) にも含まれています。

up-to-date フレームワークとそのバージョンにおけるサポートの変更に遅れずについていくにはどうすればいいですか？

DLAMI リリースノートの「[フレームワークSupport ポリシー](#)」の表を参照して、[DLAMI up-to-date AWS Deep Learning AMI フレームワークとバージョンを把握してください](#)。

Anaconda リポジトリを使用するには商用ライセンスが必要ですか？

Anaconda は特定のユーザー向けの商用ライセンスモデルに移行しました。アクティブにメンテナンスされている DLAMI は、Anaconda チャンネルから公開されているオープンソースバージョンの Conda ([conda-forge](#)) に移行されました。

DLAMI を起動および設定する

このページを表示している場合は、起動する AMI について既に十分考慮済みかと思えます。そうでない場合は、[DLAMI とそれに関連するハードウェア、フレームワーク、および ID の取得を見つけます](#)[DLAMI のリリースノート](#)。

さらに、選択するインスタンスタイプとリージョンを把握しておく必要があります。そうでない場合、「[DLAMI のインスタンスタイプを選択する](#)」を参照してください

Note

この例では、デフォルトのインスタンスタイプとして p3.16xlarge を使用します。このインスタンスタイプは、予定しているものに置き換えてください。

Important

Elastic Inference を使用する予定の場合は、DLAMI を起動する前に、[Elastic Inference の設定](#)を完了する必要があります。

トピック

- [ステップ 1: DLAMI を起動する](#)
- [ステップ 2: DLAMI に接続する](#)
- [ステップ 3: DLAMI をテストする](#)
- [ステップ 4: DLAMI インスタンスを管理する](#)
- [クリーンアップ](#)
- [Jupyter ノートブックサーバーの設定](#)

ステップ 1: DLAMI を起動する

Note

このチュートリアルでは、Deep Learning AMI (Ubuntu 18.04) に固有の内容に言及する可能性があります。異なる DLAMI を選択している場合でも、このガイドに従うことができます。

1. [DLAMI の ID を検索します。](#)
2. [DLAMI から Amazon EC2 インスタンスを起動します。](#)

Amazon EC2 コンソールを使用します。「[Amazon EC2 コンソールから起動する](#)」で説明している手順に従ってください。

Tip

CLI オプション: AWS CLI を使用して DLAMI を起動する場合は、AMI の ID、リージョンとインスタンスタイプ、およびセキュリティトークン情報が必要です。AMI とインスタンスの ID があることを確かめてください。AWS CLI をまだセットアップしていない場合は、まず [AWS コマンドラインインターフェイスのインストール](#) のガイドを使用して設定を行います。

3. これらのいずれかのオプションのステップを完了したら、インスタンスの準備完了を待機します。この待機時間は通常わずか数分です。インスタンスの状態は [EC2 コンソール](#) で確認できます。

DLAMI ID を取得する

各 AMI には一意の識別子 (ID) があります。AWS コマンドラインインターフェイス (AWS CLI) を使用して、選択した DLAMI の ID をクエリできます。がまだ AWS CLI インストールされていない場合は、「[AWS CLI の開始方法](#)」を参照してください。

Note

注意: すべての DLAMI とそれに関連するプロセッサ/アクセラレーター、オペレーティングシステム、コンピューティングアーキテクチャ、推奨される Amazon EC2 インスタンスファミリー

ミラー、サポートステータス、ID 取得クエリは、[AWS Deep Learning AMI カタログ](#)で確認できます。追加情報 (ドライバー、Python バージョン、Amazon EBS タイプ) については、「[DLAMI のリリースノート](#)」で DLAMI リリースノートも参照してください。

1. AWS 認証情報が設定されていることを確認します。

```
aws configure
```

2. 次のコマンドを使用して、DLAMI の ID を取得するか、AWS Deep Learning AMI カタログで指定されたクエリを検索します。

```
aws ec2 describe-images --region us-east-1 --owners amazon \  
--filters 'Name=name,Values=Deep Learning AMI (Ubuntu 18.04) Version ???.?' \  
'Name=state,Values=available' \  
--query 'reverse(sort_by(Images, &CreationDate))[1].ImageId' --output text
```

Note

特定のフレームワークのリリースバージョンを指定したり、バージョン番号を疑問符に置き換えて最新のリリースを取得したりできます。

3. 出力は次の例に類似したものになります:

```
ami-094c089c38ed069f2
```

この DLAMI ID をコピーし、q を押してプロンプトを終了します。

次のステップ


[Amazon EC2 コンソールから起動する](#)

Amazon EC2 コンソールから起動する

Note


Elastic Fabric Adapter (EFA) を使用してインスタンスを起動するには、[この手順](#)を参照してください。

1. [EC2 コンソール](#)を開きます。
2. 最上部のナビゲーションバーで現在のリージョンを確認します。これが目的のリージョンでない場合は AWS リージョン、先に進む前にこのオプションを変更してください。詳細については、[EC2 リージョン](#)を参照してください。
3. [Launch Instance] (インスタンスの起動) を選択します。
4. インスタンスの名前を入力し、適切な DLAMI を選択します。
 - a. [自分の全 AMI] で既存の DLAMI を検索するか、[クイックスタート] を選択します。
 - b. DLAMI ID で検索します。オプションを参照し、目的の項目を選択します。
5. インスタンスタイプを選択します。DLAMI の推奨インスタンスファミリーは、AWS Deep Learning AMI カタログにあります。DLAMI インスタンスタイプに関する推奨事項については、「[インスタンスの選択](#)」を参照してください。

 Note

[Elastic Inference](#) (EI) を使用する場合は、[Configure Instance Details] (インスタンスの詳細の設定) をクリックし、[Add an Amazon EI accelerator] (Amazon EI アクセラレーターの追加) を選択してから、Amazon EI アクセラレーターのサイズを選択します。

6. [Launch Instance] (インスタンスの起動) を選択します。

 Tip

スクリーンショット付きのチュートリアルについては、「[Deep Learning AMI を使用した AWS Deep Learning の開始方法](#)」を参照してください。

次のステップ

[ステップ 2: DLAMI に接続する](#)

ステップ 2: DLAMI に接続する

クライアント (Windows、MacOS、または Linux) から起動した DLAMI に接続します。詳細については、「Amazon EC2 ユーザーガイド」の「[Linux インスタンスに接続する](#)」を参照してください。

Amazon EC2

ログイン後に Jupyter のセットアップを実行する場合は、SSH ログインコマンドをコピーしていつでも参照できるようにしておきます。Jupyter ウェブページに接続する際に、このコマンドのバリエーションを使用します。

次のステップ

[ステップ 3: DLAMI をテストする](#)

ステップ 3: DLAMI をテストする

DLAMI のバージョンに応じて、さまざまなテストのオプションを利用できます。

- [Deep Learning AMI with Conda](#) – [Deep Learning AMI with Conda の使用](#) に移動します。
- [Deep Learning Base AMI](#) - 目的のフレームワークのインストールに関するドキュメントを参照してください。

Jupyter ノートブックの作成、チュートリアルの実行、または Python によるコーディングの開始も可能です。詳細については、「[Jupyter ノートブックサーバーの設定](#)」を参照してください。

ステップ 4: DLAMI インスタンスを管理する

オペレーティングシステムやその他のインストールされているソフトウェアは、パッチや更新が利用可能になり次第適用して、常に最新の状態を保ってください。

Amazon Linux または Ubuntu を使用している場合は、DLAMI にログインしたときに、更新が利用可能であれば通知され、更新の手順が表示されます。Amazon Linux のメンテナンスの詳細については、[インスタンスソフトウェアの更新](#)を参照してください。Ubuntu インスタンスの場合は、[Ubuntu の公式ドキュメント](#)を参照してください。

Windows では、Windows Update でソフトウェアとセキュリティの更新を定期的にチェックします。必要に応じて、更新が自動的に適用されるようにします。

Important

Meltdown と Spectre の脆弱性と、オペレーティングシステムにパッチを適用して対処する方法の詳細については、「[セキュリティ情報 AWS-2018 年 13 月](#)」を参照してください。

クリーンアップ

DLAMI が不要になった場合は、停止または終了させることによって継続的な料金の発生を避けることができます。インスタンスを停止してもそのまま保持されるため、後で再開できます。設定、ファイル、およびその他の非揮発性情報は Amazon S3 のボリュームに格納されます。インスタンスの停止中は S3 の料金がわずかに課金されますが、停止状態にある間はコンピューティングリソースへの課金は停止されます。再度インスタンスを起動すると、そのボリュームがマウントされ、データを利用できるようになります。インスタンスを終了した場合は、インスタンスは消去され、再度起動することはできません。実際にはデータがまだ S3 に存在するため、それ以降の変更を防止するためには、ボリュームも削除する必要があります。詳細については、「[Amazon EC2 ユーザーガイド](#)」の「[インスタンスの終了](#)」を参照してください。Amazon EC2

Jupyter ノートブックサーバーの設定

Jupyter ノートブックサーバーを使用すると、DLAMI インスタンスから Jupyter ノートブックを作成および実行できます。Jupyter ノートブックを使用すると、AWS インフラストラクチャを使用し、DLAMI に組み込まれたパッケージにアクセスしながら、トレーニングと推論のための機械学習 (ML) の実験を作成および実行できます。Jupyter ノートブックの詳細については、[Jupyter ノートブックのドキュメント](#)を参照してください。

Jupyter ノートブックサーバーを設定するには、次の作業を実行する必要があります。

- Amazon EC2 DLAMI インスタンスで Jupyter ノートブックサーバーを設定する。
- Jupyter ノートブックサーバーに接続できるようにクライアントを設定する。Windows クライアント、macOS クライアント、および Linux クライアントの設定手順が用意されています。
- Jupyter ノートブックサーバーにログインして、セットアップをテストする。

Jupyter を設定する手順を完了するには、以下のトピックの指示に従ってください。Jupyter ノートブックサーバーを設定したら、DLAMI に同梱されているノートブックの例を実行する方法について、[Jupyter ノートブックチュートリアルを実行する](#)を参照してください。

トピック

- [Jupyter サーバーの保護](#)
- [Jupyter ノートブックサーバーの起動](#)
- [Jupyter サーバーに接続するためのクライアントの設定](#)
- [Jupyter ノートブックサーバーへのログインによるテスト](#)

Jupyter サーバーの保護

ここでは、SSL およびカスタムパスワードで Jupyter を設定します。

Amazon EC2 インスタンスに接続した後、次の手順を実行します。

Jupyter サーバーの設定

1. Jupyter にはパスワードユーティリティが用意されています。次のコマンドを実行して、プロンプトに任意のパスワードを入力します。

```
$ jupyter notebook password
```

出力は次のようになります。

```
Enter password:  
Verify password:  
[NotebookPasswordApp] Wrote hashed password to /home/ubuntu/.jupyter/  
jupyter_notebook_config.json
```

2. 自己署名 SSL 証明書を作成します。プロンプトに従って、該当するローカリティを入力します。プロンプトを空白のままにする場合は、. を入力する必要があります。この回答は証明書の機能に影響を及ぼしません。

```
$ cd ~  
$ mkdir ssl  
$ cd ssl  
$ openssl req -x509 -nodes -days 365 -newkey rsa:2048 -keyout mykey.key -out  
mycert.pem
```

Note

サードパーティが署名した定期的 SSL 証明書を作成して、ブラウザでセキュリティ警告が発生しないようにすることもできます。このプロセスは非常に複雑になります。詳細については、[Jupyter ドキュメント](#)を参照してください。

次のステップ

[Jupyter ノートブックサーバーの起動](#)

Jupyter ノートブックサーバーの起動

これで、インスタンスにログインし、前のステップで作成した SSL 証明書を使用して次のコマンドを実行して Jupyter サーバーを開始できます。

```
$ jupyter notebook --certfile=~/.ssl/mycert.pem --keyfile ~/.ssl/mykey.key
```

サーバーが起動していれば、クライアントコンピュータから SSH トンネルを通じてサーバーに接続できます。サーバーを実行すると、サーバーが実行中であることを確認するメッセージが Jupyter から出力されます。トンネルを作成するまでは機能しないため、このとき表示される localhost URL 経由でサーバーにアクセスできるという吹き出しは無視してください。

Note

Jupyter ウェブインターフェイスを使用してフレームワークを切り替えるとき、Jupyter は自動的に環境に切り替えを処理します。これについての詳細は、「[Jupyter で環境を切り替える](#)」を参照してください。

次のステップ

[Jupyter サーバーに接続するためのクライアントの設定](#)

Jupyter サーバーに接続するためのクライアントの設定

Jupyter ノートブックサーバーに接続するようにクライアントを設定した後、サーバー上のワークスペース内にノートブックを作成して、そのノートブックにアクセスしたり、サーバー上でディープラーニングコードを実行したりできます。

設定情報については、次のいずれかのリンクを選択してください。

トピック

- [Windows クライアントの設定](#)
- [Linux または macOS クライアントの設定](#)

Windows クライアントの設定

Prepare

次の情報が手元に揃っていることを確認します。これらの情報は、SSH トンネルをセットアップする際に必要です。

- Amazon EC2 インスタンスのパブリック DNS 名。パブリック DNS 名は、EC2 コンソールで確認できます。
- プライベートキーファイルのキーペア。キーペアへのアクセスの詳細については、[Linux インスタンス用 Amazon EC2 ユーザーガイド](#)の「Amazon EC2 のキーペア」を参照してください。

Windows クライアントから Jupyter Norebooks を使用する

Windows クライアントから Amazon EC2 インスタンスへの接続に関する以下のガイドを参照してください。

1. [インスタンスへの接続に関するトラブルシューティング](#)
2. [PuTTY を使用した Windows から Linux インスタンスへの接続](#)

実行中の Jupyter サーバーにトンネルを作成するとき、Windows クライアントに Git Bash をインストールして、Linux/macOS クライアントの説明に従うアプローチが推奨されます。ただし、ポートがマッピングされた SSH トンネルを開くためにどのアプローチを使用することもできます。詳細については、[Jupyter のドキュメント](#)を参照してください。

次のステップ

[Linux または macOS クライアントの設定](#)

Linux または macOS クライアントの設定

1. ターミナルを開きます。
2. ローカルポート 8888 に対するすべてのリクエストをリモート Amazon EC2 インスタンスのポート 8888 に転送するために、次のコマンドを実行します。Amazon EC2 インスタンスにアクセスするキーの場所と Amazon EC2 インスタンスのパブリック DNS 名を置き換えて、コマンドを更新します。Amazon Linux AMI の場合、ユーザー名は ubuntu の代わりに ec2-user であることに注意してください。

```
$ ssh -i ~/mykeypair.pem -N -f -L 8888:localhost:8888 ubuntu@ec2-###-##-##-###.compute-1.amazonaws.com
```

このコマンドを実行すると、Jupyter ノートブックサーバーを実行しているリモート Amazon EC2 インスタンスとクライアントの間にトンネルが開通します。

次のステップ

Jupyter ノートブックサーバーへのログインによるテスト

Jupyter ノートブックサーバーへのログインによるテスト

Jupyter ノートブックサーバーにログインする準備ができました。

次のステップでは、ブラウザを使用してサーバーへの接続をテストします。

1. ブラウザのアドレスバーに次の URL を入力するか、次のリンクをクリックしてください。 <https://localhost:8888>
2. 自己署名 SSL 証明書を使用すると、ブラウザは警告を表示し、このウェブサイトの閲覧を続行しないように求められます。

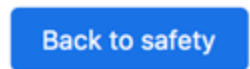


Your connection is not private

Attackers might be trying to steal your information from **localhost** (for example, passwords, messages, or credit cards). [Learn more](#)

NET::ERR_CERT_AUTHORITY_INVALID

Help improve Safe Browsing by sending some [system information and page content](#) to Google.
[Privacy policy](#)



この設定は自身が行ったものであるため、安全に続行できます。ブラウザによっては、「高度な設定」、「詳細の表示」などのボタンが表示されることがあります。



Your connection is not private

Attackers might be trying to steal your information from **localhost** (for example, passwords, messages, or credit cards). [Learn more](#)

NET::ERR_CERT_AUTHORITY_INVALID

Help improve Safe Browsing by sending some [system information and page content](#) to Google.
[Privacy policy](#)

Hide advanced

Back to safety

This server could not prove that it is **localhost**; its security certificate is not trusted by your computer's operating system. This may be caused by a misconfiguration or an attacker intercepting your connection.

[Proceed to localhost \(unsafe\)](#)

これをクリックして「localhostに進む」リンクをクリックします。接続に成功すると、Jupyter ノートブックサーバーのウェブページが表示されます。この時点で、前に設定したパスワードを入力するように求められます。

これで、DLAMI 上で実行中の Jupyter ノートブックサーバーにアクセスできるようになりました。新しいノートブックを作成することも、用意されている [チュートリアル](#) を実行することもできます。

DLAMI の使用

トピック

- [Deep Learning AMI with Conda の使用](#)
- [Deep Learning Base AMI の使用](#)
- [Jupyter ノートブックチュートリアルを実行する](#)
- [チュートリアル](#)

以降のセクションでは、Deep Learning AMI with Conda を使用して環境を切り替える方法、各フレームワークからサンプルコードを実行する方法、および Jupyter を実行して、さまざまなノートブックチュートリアルを試す方法について説明します。

Deep Learning AMI with Conda の使用

トピック

- [Deep Learning AMI with Conda の概要](#)
- [DLAMI にログインする](#)
- [TensorFlow 環境を開始する](#)
- [PyTorch Python 3 環境に切り替える](#)
- [環境を削除する](#)

Deep Learning AMI with Conda の概要

Conda は、Windows、macOS、および Linux で稼働するオープンソースのパッケージ管理システムおよび環境管理システムです。Conda では、パッケージとその依存関係を迅速にインストール、実行、および更新できます。また、ローカルコンピュータ上で簡単に環境を作成、保存、ロードし、環境を切り替えることができます。

Deep Learning AMI with Conda は、深層学習環境を簡単に切り替えられるように設定されています。次に示す手順では、conda での基本的なコマンドを説明します。これらのコマンドは、フレームワークの基本インポートが機能していることや、フレームワークでいくつかの簡単な操作を実行できることの確認にも役立ちます。その後、DLAMI で提供されているより詳細なチュートリアルや、各フレームワークのプロジェクトサイトにあるフレームワークの例に進むことができます。

DLAMI にログインする

サーバーにログインすると、サーバーの「その日のメッセージ (MOTD)」が表示され、さまざまなディープラーニングフレームワークを切り替えるための各種の Conda コマンドが示されます。以下に MOTD の例を示します。新しいバージョンの DLAMI がリリースされる事により、場合によっては特定の MOTD が異なる場合があります。

```
=====
AMI Name: Deep Learning OSS Nvidia Driver AMI (Amazon Linux 2) Version 77
Supported EC2 instances: G4dn, G5, G6, Gr6, P4d, P4de, P5
  * To activate pre-built tensorflow environment, run: 'source activate
tensorflow2_p310'
  * To activate pre-built pytorch environment, run: 'source activate
pytorch_p310'
  * To activate pre-built python3 environment, run: 'source activate python3'

NVIDIA driver version: 535.161.08

CUDA versions available: cuda-11.7 cuda-11.8 cuda-12.0 cuda-12.1 cuda-12.2

Default CUDA version is 12.1

Release notes: https://docs.aws.amazon.com/dlami/latest/devguide/appendix-ami-release-notes.html
AWS Deep Learning AMI Homepage: https://aws.amazon.com/machine-learning/amis/
Developer Guide and Release Notes: https://docs.aws.amazon.com/dlami/latest/devguide/what-is-dlami.html
Support: https://forums.aws.amazon.com/forum.jspa?forumID=263
For a fully managed experience, check out Amazon SageMaker at https://aws.amazon.com/sagemaker
=====
```

TensorFlow 環境を開始する

Note

初めて Conda 環境を起動する際には、ロードするまで辛抱して待機してください。Deep Learning AMI with Conda は、EC2 インスタンスにフレームワークの最初のアクティベーションから最適化されたフレームワークのバージョンを自動的にインストールします。継続する遅延はありません。

1. Python 3 の TensorFlow 仮想環境をアクティブ化します。

```
$ source activate tensorflow2_p310
```

2. iPython ターミナルを起動します。

```
(tensorflow2_p310)$ ipython
```

3. クイック TensorFlow プログラムを実行します。

```
import tensorflow as tf
hello = tf.constant('Hello, TensorFlow!')
sess = tf.Session()
print(sess.run(hello))
```

"Hello, Tensorflow!" が表示されます。

次回の予定

[Jupyter ノートブックチュートリアルを実行する](#)

PyTorch Python 3 環境に切り替える

iPython コンソールをまだ使用している場合は、`quit()` を使用して、環境を切り替える準備を整えます。

- Python 3 の PyTorch 仮想環境をアクティブ化します。

```
$ source activate pytorch_p310
```

一部の PyTorch コードをテストする

インストールをテストするには、Python を使用して配列を作成および出力する PyTorch コードを記述します。

1. iPython ターミナルを起動します。

```
(pytorch_p310)$ ipython
```

2. をインポートします PyTorch。

```
import torch
```

サードパーティー製パッケージに関する警告メッセージが表示される場合があります。このメッセージは無視できます。

3. 各要素をランダムに初期化して、5x3 行列を作成します。配列を出力します。

```
x = torch.rand(5, 3)
print(x)
```

結果を確認します。

```
tensor([[0.3105, 0.5983, 0.5410],
        [0.0234, 0.0934, 0.0371],
        [0.9740, 0.1439, 0.3107],
        [0.6461, 0.9035, 0.5715],
        [0.4401, 0.7990, 0.8913]])
```

環境を削除する

DLAMI にスペースが足りない場合、使用していない Conda パッケージをアンインストールすることを選択できます。

```
conda env list
conda env remove --name <env_name>
```

Deep Learning Base AMI の使用

Deep Learning Base AMI の使用

Base AMI には、独自でカスタマイズしたディープラーニング環境をデプロイするための GPU ドライバの基盤プラットフォームとアクセラレーションライブラリが含まれています。デフォルトでは、AMI は任意の 1 つの NVIDIA CUDA バージョン環境で設定されます。CUDA のさまざまなバージョン間で切り替えを行うこともできます。これを行う方法については、次の手順を参照してください。

CUDA のバージョンの設定

NVIDIA の `nvcc` プログラムを実行して、CUDA バージョンを検証できます。

```
nvcc --version
```

次の Bash コマンドを使用して、特定の CUDA バージョンを選択して検証できます。

```
sudo rm /usr/local/cuda
sudo ln -s /usr/local/cuda-12.0 /usr/local/cuda
```

詳細については、「[Base DLAMI リリースノート](#)」を参照してください。

Jupyter ノートブックチュートリアルを実行する

各深層学習プロジェクトにはチュートリアルと例が付属しています。ほとんどの場合、それらはすべての DLAMI で実行することができます。[Deep Learning AMI with Conda](#) を選択した場合は、厳選されたチュートリアルが既にセットアップされ、すぐに試すことができるという利点があります。

⚠ Important

DLAMI にインストールされている Jupyter ノートブックチュートリアルを実行するには、[Jupyter ノートブックサーバーの設定](#)の内容を実行する必要があります。

Jupyter サーバー実行されたら、ウェブブラウザからチュートリアルを実行できます。Deep Learning AMI with Conda を実行している場合や、Python 環境をセットアップしている場合は、Jupyter ノートブックのインターフェイスから Python カーネルを切り替えることができます。フレームワーク固有のチュートリアルを実行する前に、適切なカーネルを選択してください。これに関するその他の例が Deep Learning AMI with Conda のユーザー向けに提供されています。

ℹ Note

多くのチュートリアルでは追加の Python モジュールが必要になりますが、それらが DLAMI にセットアップされていない可能性があります。"xyz module not found" などのエラーが発生した場合は、DLAMI にログインし、上記を参考に環境をアクティブ化して、必要なモジュールをインストールしてください。

i Tip

ディープラーニングのチュートリアルや例の多くは、1つ以上の GPU を必要とします。GPU を使用しないインスタンスタイプの場合、例のコードを実行するために、コードを一部変更することが必要になる可能性があります。

インストールされたチュートリアルを操作する

Jupyter サーバーにログインしてチュートリアルディレクトリが表示されると (Deep Learning AMI with Conda 上のみ)、各フレームワークの名前別に整理されたチュートリアルのフォルダを確認できます。フレームワークのリストが表示されない場合は、現在の DLAMI で、そのフレームワークのチュートリアルが提供されていないことを示しています。フレームワークの名前をクリックし、リストされているチュートリアルを確認して、チュートリアルをクリックして起動します。

Deep Learning AMI with Conda でノートブックを初めて実行すると、使用する環境を指定するように求められます。選択のためのリストがプロンプトで表示されます。それぞれの環境は、以下のパターンに従って命名されています。

Environment (conda_framework_python-version)

たとえば、"Environment (conda_mxnet_p36)" と表示された場合は、環境に MXNet と Python 3 が存在しています。別のバリエーションとして "Environment (conda_mxnet_p27)" があります。この場合は環境に MXNet と Python 2 が存在しています。

i Tip

アクティブ化される CUDA のバージョンについて不明な点がある場合は、最初に DLAMI にログインする際の MOTD で確認できます。

Jupyter で環境を切り替える

別のフレームワークのチュートリアルを試す場合は、現在実行中のカーネルを検証する必要があります。この情報は、Jupyter のインターフェイスの右上、ログアウトボタンの真下に表示されます。開いているノートブックでカーネルを変更するには、Jupyter のメニュー項目を [Kernel]、[Change Kernel] の順にクリックして、実行中のノートブックに適合する環境をクリックします。

カーネルに変更が加わると、これまでに実行したすべての項目の状態が消去されるため、この時点ですべてのセルを再度実行することが必要になります。

Tip

フレームワークの切り替えは、興味深くさまざまな情報も入手できる操作ですが、メモリ不足が発生する可能性があります。エラーが表示されるようになったら、実行中の Jupyter サーバーが表示されているターミナルウィンドウを確認してください。ここには役立つメッセージとエラーログがあり、エラー out-of-memoryが表示されることがあります。この問題を修正するには、Jupyter サーバーのホームページにアクセスして、バックグラウンドで実行中でメモリを占有していると思われる各チュートリアルで、[Running] タブ、[Shutdown] の順にクリックします。

チュートリアル

Deep Learning AMI with Conda のソフトウェアの使用方法のチュートリアルを以下に示します。

トピック

- [10 分間チュートリアル](#)
- [フレームワークのアクティブ化](#)
- [Elastic Fabric Adapter を使用した分散トレーニング](#)
- [GPU のモニタリングおよび最適化](#)
- [DLAMI を備えた AWS Inferentia チップ](#)
- [ARM64 DLAMI](#)
- [推論](#)
- [モデル提供](#)

10 分間チュートリアル

- [を起動する AWS Deep Learning AMI \(10 分\)](#)
- [Amazon EC2 で DLC を使用して深層学習モデルをトレーニングする \(10 分\)](#)

フレームワークのアクティブ化

以下に示しているのは、Deep Learning AMI with Conda にインストールされている深層学習フレームワークです。フレームワークをクリックすると、そのフレームワークをアクティブ化する方法を参照できます。

トピック

- [PyTorch](#)
- [TensorFlow 2](#)

PyTorch

のアクティブ化 PyTorch

フレームワークの安定版 Conda パッケージがリリースされると、DLAMI でテストされ事前にインストールされます。テストされていない最新のナイトリービルドを実行する場合は、手動で「[PyTorch のナイトリービルドをインストールする \(実験的\)](#)」ことができます。

現在インストールされているフレームワークをアクティブ化するには、使用する Deep Learning AMI with Conda に関する以下の手順に従います。

CUDA と MKL-DNN を使用する Python 3 PyTorch 上では、次のコマンドを実行します。

```
$ source activate pytorch_p310
```

iPython ターミナルを起動します。

```
(pytorch_p310)$ ipython
```

クイック PyTorch プログラムを実行します。

```
import torch
x = torch.rand(5, 3)
print(x)
print(x.size())
y = torch.rand(5, 3)
print(torch.add(x, y))
```

最初のランダムな配列が出力された後、そのサイズが出力され、次に別のランダムな配列が追加で出力されます。

PyTorchのナイトリービルドをインストールする (実験的)

ナイトリービルド PyTorch から をインストールする方法

Deep Learning AMI with Conda の PyTorch Conda 環境のいずれかまたは両方に最新の PyTorch ビルドをインストールできます。

1. • (Python 3 のオプション) - Python 3 PyTorch 環境をアクティブ化します。

```
$ source activate pytorch_p310
```

2. 残りの手順は、pytorch_p310 環境を使用していることを前提としています。現在インストールされている を削除します PyTorch。

```
(pytorch_p310)$ pip uninstall torch
```

3. • (GPU インスタンスのオプション) - CUDA.0 PyTorch で の最新夜間ビルドをインストールします。

```
(pytorch_p310)$ pip install torch_nightly -f https://download.pytorch.org/whl/nightly/cu100/torch_nightly.html
```

- (CPU インスタンスのオプション) - GPU のないインスタンス PyTorch に の最新夜間ビルドをインストールします。 GPUs

```
(pytorch_p310)$ pip install torch_nightly -f https://download.pytorch.org/whl/nightly/cpu/torch_nightly.html
```

4. 最新の夜間ビルドが正常にインストールされたことを確認するには、IPython ターミナルを起動し、 のバージョンを確認します PyTorch。

```
(pytorch_p310)$ ipython
```

```
import torch
print (torch.__version__)
```

出力は 1.0.0.dev20180922 のように表示されます。

5. PyTorch 夜間ビルドが MNIST の例とうまく機能することを確認するには、PyTorchのサンプルリポジトリからテストスクリプトを実行します。

```
(pytorch_p310)$ cd ~
(pytorch_p310)$ git clone https://github.com/pytorch/examples.git pytorch_examples
(pytorch_p310)$ cd pytorch_examples/mnist
(pytorch_p310)$ python main.py || exit 1
```

他のチュートリアル

その他のチュートリアルと例については、フレームワークの公式ドキュメント、[PyTorch ドキュメント](#)、および[PyTorch](#) ウェブサイトを参照してください。

TensorFlow 2

このチュートリアルでは、Deep Learning AMI with Conda (DLAMI on Conda) を実行しているインスタンスで TensorFlow 2 をアクティブ化し、TensorFlow 2 プログラムを実行する方法を示します。

フレームワークの安定版 Conda パッケージがリリースされると、DLAMI でテストされ事前にインストールされます。

アクティブ化 TensorFlow 2

DLAMI with Conda TensorFlow で実行するには

1. TensorFlow 2 を有効にするには、DLAMI with Conda の Amazon Elastic Compute Cloud (Amazon EC2) インスタンスを開きます。
2. CUDA 10.1 と MKL-DNN を使用する Python 3 の TensorFlow 2 と Keras 2 の場合は、次のコマンドを実行します。

```
$ source activate tensorflow2_p310
```

3. iPython ターミナルを起動します。

```
(tensorflow2_p310)$ ipython
```

4. TensorFlow 2 プログラムを実行して、正常に動作していることを確認します。


```
import tensorflow as tf
hello = tf.constant('Hello, TensorFlow!')
tf.print(hello)
```

Hello, TensorFlow! が画面に表示される必要があります。

他のチュートリアル

その他のチュートリアルと例については、[TensorFlow Python API](#) の TensorFlow ドキュメントまたは [TensorFlow](#) ウェブサイトを参照してください。

Elastic Fabric Adapter を使用した分散トレーニング

[Elastic Fabric Adapter](#) (EFA) は、ハイパフォーマンスコンピューティング (HPC) アプリケーションを高速化するために DLAMI インスタンスにアタッチできるネットワークデバイスです。EFA を使用すると、AWS クラウドが提供するスケーラビリティ、柔軟性、伸縮性を備えたオンプレミス HPC クラスターのアプリケーションパフォーマンスを実現できます。

以下のトピックでは、EFA と DLAMI の使用を開始する方法について説明します。

Note

この [Base GPU DLAMI リスト](#) からお使いの DLAMI を選択してください。

トピック

- [EFA AWS Deep Learning AMI を使用したインスタンスの起動](#)
- [DLAMI での EFA の使用](#)

EFA AWS Deep Learning AMI を使用したインスタンスの起動

最新版 Base DLAMI は EFA と共に使用する準備ができており、GPU インスタンス用の必要なドライバー、カーネルモジュール、Libfabric、OpenMPI および [NCCL OFI プラグイン](#) が付属しています。

サポートされている Base DLAMI の CUDA バージョンは、[リリースノート](#) で確認できます。

[Note:] (メモ:)

- EFA 上の mpirun を使用して NCCL アプリケーションを実行する際、EFA がサポートされているインストラクションへの完全なパスを以下のように指定する必要があります。

```
/opt/amazon/openmpi/bin/mpirun <command>
```

- アプリケーションで EFA を使用するには、[DLAMI での EFA の使用](#) に示すように、mpirun コマンドに FI_PROVIDER="efa" を追加します。

トピック

- [EFA 対応のセキュリティグループを準備する](#)
- [インスタンスの起動](#)
- [EFA 添付ファイルの確認](#)

EFA 対応のセキュリティグループを準備する

EFA には、セキュリティグループ自体との間で送受信されるすべてのトラフィックを許可するセキュリティグループが必要です。詳細については、[「EFA ドキュメント」](#)を参照してください。

1. Amazon EC2 コンソール (<https://console.aws.amazon.com/ec2/>) を開きます。
2. ナビゲーションペインで [セキュリティグループ] を選択して、[セキュリティグループの作成] を選択します。
3. [セキュリティグループの作成] ウィンドウで、以下を行います。
 - [セキュリティグループ名] に、EFA-enabled security group のような、分かりやすいセキュリティグループ名を入力します。
 - (オプション) [説明] に、セキュリティグループの簡単な説明を入力します。
 - [VPC] で、EFA 対応のインスタンスを起動する VPC を選択します。
 - [作成] を選択します。
4. 作成したセキュリティグループを選択し、[説明] タブで [グループ ID] をコピーします。
5. [Inbound] (インバウンド) タブおよび [Outbound] (アウトバウンド) タブで、次の手順を実行します。
 - [Edit] を選択します。
 - [Type] で、[All traffic] を選択します。

- [ソース] で [カスタム] を選択します。
 - コピーしたセキュリティグループ ID をフィールドに貼り付けます。
 - [保存] を選択します。
6. 「[Linux インスタンスのインバウンドトラフィックの承認](#)」を参照するインバウンドトラフィックを有効にします。このステップを抜かすと、DLAMI インスタンスと通信できなくなります。

インスタンスの起動

の EFA AWS Deep Learning AMI は現在、以下のインスタンスタイプとオペレーティングシステムでサポートされています。

- P3dn .24xlarge: Amazon Linux 2、Ubuntu 20.04
- P4d .24xlarge: Amazon Linux 2、Ubuntu 20.04
- P5.48xlarge: Amazon Linux 2、Ubuntu 20.04

以下のセクションでは、EFA 対応の DLAMI インスタンスを起動する方法について説明します。EFA 対応のインスタンスの起動の詳細については、[クラスタープレイメントグループで EFA 対応のインスタンスを起動する](#)を参照してください。

1. Amazon EC2 コンソール (<https://console.aws.amazon.com/ec2/>) を開きます。
2. [Launch Instance] (インスタンスの起動) を選択します。
3. AMI の選択ページで、DLAMI [リリースノートページにあるサポートされている DLAMI](#) を選択します。
4. [Choose an Instance Type] (インスタンスタイプの選択) ページで、次のいずれかのサポート対象のインスタンスタイプを選択し、[Next: Configure Instance Details] (次の手順: インスタンスの詳細の設定) を選択します。サポートされているインスタンスのリストについては、このリンクを参照してください。[EFA と MPI の使用を開始する](#)
5. [Configure Instance Details] ページで以下の操作を実行します。
 - [インスタンス数] に、起動する EFA 対応のインスタンスの数を入力します。
 - [ネットワーク] および [サブネット] で、インスタンスを起動する VPC およびサブネットを選択します。
 - [オプション] プレイメントグループで、プレイメントグループにインスタンスを追加を選択します。最適なパフォーマンスを得るには、プレイメントグループ内でインスタンスを起動します。

- [オプション] プレイメントグループ名 で、新しいプレイメントグループ に追加 を選択し、プレイメントグループのわかりやすい名前を入力し、プレイメントグループ戦略 でクラスター を選択します。
 - 必ず、このページで [Elastic Fabric Adapter] を有効にしてください。このオプションが無効になっている場合は、選択したインスタンスタイプに対応するサブネットに変更します。
 - [ネットワークインターフェイス] セクションの [eth0] で、[新しいネットワークインターフェイス] を選択します。必要に応じて、プライマリ IPv4 アドレスと 1 つ以上のセカンダリ IPv4 アドレスを指定できます。関連付けられている IPv6 CIDR ブロックを持つサブネットにインスタンスを起動する場合は、必要に応じて、プライマリ IPv6 アドレスと 1 つ以上のセカンダリ IPv6 アドレスを指定することができます。
 - [次の手順: ストレージの追加] を選択します。
6. [ストレージの追加] ページで、AMI で指定されたボリュームに加えてインスタンスにアタッチするボリューム (例: ルートデバイスのボリューム) を指定し、[Next: Add Tags (次へ: タグの追加)] を選択します。
 7. [Add Tags] ページで、ユーザーフレンドリーな名前などを使ってインスタンスのタグを指定し、[Next: Configure Security Group] を選択します。
 8. セキュリティグループの設定ページで、セキュリティグループの割り当て で既存のセキュリティグループの選択 を選択し、前に作成したセキュリティグループを選択します。
 9. [Review and Launch] を選択します。
 10. [インスタンス作成の確認] ページで設定を確認し、[起動] を選択してキーペアを選択し、インスタンスを起動します。

EFA 添付ファイルの確認

コンソールから

インスタンスを起動したら、AWS コンソールでインスタンスの詳細を確認します。これを行うには、EC2 コンソールでインスタンスを選択し、ページ下のペインにある [Description (説明)] タブを確認します。[Network Interfaces: eth0 (ネットワークインターフェイス: eth0)] というパラメータを探し、eth0 をクリックするとポップアップが開きます。[Elastic Fabric Adapter] が有効になっていることを確認します。

EFA が有効になっていない場合は、次のいずれかの方法でこれを修正できます。

- EC2 インスタンスを終了し、同じ手順で新しいインスタンスを起動します。EFA が添付されていることを確認します。

- 既存のインスタンスに EFA を添付します。
 1. EC2 コンソールで、[Network Interfaces (ネットワークインターフェイス)] に移動します。
 2. [Create a Network Interface (ネットワークインターフェイスを作成)] をクリックします。
 3. インスタンスが入っている同じサブネットを選択します。
 4. 必ず、[Elastic Fabric Adapter] を有効にし、[Create] (作成) をクリックします。
 5. [EC2 Instances (EC2 インスタンス)] タブに戻り、インスタンスを選択します。
 6. [Actions: Instance State] (アクション: インスタンスの状態) に移動し、EFA にアタッチする前にインスタンスを停止します。
 7. [Actions (アクション)] から、[Networking: Attach Network Interface (ネットワーキング: ネットワークインターフェイスの接続)] を選択します。
 8. 先ほど作成したインターフェイスを選択し、[Attach (接続)] をクリックします。
 9. インスタンスを再起動します。

インスタンスから

DLAMI にすでに、次のテストスクリプトが存在します。これを実行して、カーネルモジュールが正しくロードされていることを確認します。

```
$ fi_info -p efa
```

出力は以下のようになります。

```
provider: efa
  fabric: EFA-fe80::e5:56ff:fe34:56a8
  domain: efa_0-rdm
  version: 2.0
  type: FI_EP_RDM
  protocol: FI_PROTO_EFA
provider: efa
  fabric: EFA-fe80::e5:56ff:fe34:56a8
  domain: efa_0-dgrm
  version: 2.0
  type: FI_EP_DGRAM
  protocol: FI_PROTO_EFA
provider: efa;ofi_rxd
  fabric: EFA-fe80::e5:56ff:fe34:56a8
  domain: efa_0-dgrm
```

```
version: 1.0
type: FI_EP_RDM
protocol: FI_PROTO_RXD
```

セキュリティグループ構成の確認

DLAMI にすでに、次のテストスクリプトが存在します。これを実行して、作成したセキュリティグループが正しく設定されていることを確認します。

```
$ cd /opt/amazon/efa/test/
$ ./efa_test.sh
```

出力は以下のようになります。

```
Starting server...
Starting client...
bytes  #sent  #ack  total  time  MB/sec  usec/xfer  Mxfers/sec
64     10    =10   1.2k   0.02s  0.06    1123.55    0.00
256    10    =10   5k     0.00s  17.66   14.50     0.07
1k     10    =10   20k    0.00s  67.81   15.10     0.07
4k     10    =10   80k    0.00s  237.45  17.25     0.06
64k    10    =10   1.2m   0.00s  921.10  71.15     0.01
1m     10    =10   20m    0.01s  2122.41 494.05    0.00
```

応答しなかったり完了しない場合は、セキュリティグループに正しいインバウンド/アウトバウンドルールが設定されていることを確認します。

DLAMI での EFA の使用

次の項では、EFA を使用して、AWS Deep Learning AMIでマルチノードアプリケーションを実行する方法について説明します。

EFA でマルチノードアプリケーションを実行

ノードのクラスター間でアプリケーションを実行するには、次の設定が必要です。

トピック

- [パスワードレス SSH の有効化](#)
- [Hosts ファイルの作成](#)
- [NCCL テスト](#)

パスワードレス SSH の有効化

クラスター内の 1 つのノードをリーダーノードとして選択します。残りのノードはメンバーノードと呼ばれます。

1. リーダーノードで、RSA キーペアを生成します。

```
ssh-keygen -t rsa -N "" -f ~/.ssh/id_rsa
```

2. リーダーノードのプライベートキーの許可を変更します。

```
chmod 600 ~/.ssh/id_rsa
```

3. パブリックキーを ~/.ssh/id_rsa.pub にコピーし、クラスター内のメンバーノード ~/.ssh/authorized_keys の に追加します。
4. これで、プライベート IP を使用して、リーダーノードからメンバーノードに直接ログインできるようになります。

```
ssh <member private ip>
```

5. リーダーノードの ~/.ssh/config ファイルに以下を追加して、リーダーノードでのエージェントの転送 strictHostKey の確認と有効化を無効にします。

```
Host *
    ForwardAgent yes
Host *
    StrictHostKeyChecking no
```

6. Amazon Linux 2 インスタンスでは、リーダーノードで次のコマンドを実行して、設定ファイルへの正しいアクセス許可を付与します。

```
chmod 600 ~/.ssh/config
```

Hosts ファイルの作成

リーダーノードで、クラスター内のノードを識別する Hosts ファイルを作成します。Hosts ファイルには、クラスター内の各ノードのエントリが必要です。~/hosts ファイルを作成し、次のようにプライベート IP を使用して各ノードを追加します。

```
localhost slots=8
```

```
<private ip of node 1> slots=8
<private ip of node 2> slots=8
```

NCCL テスト

Note

これらのテストは、EFA バージョン 1.30.0 と OFI NCCL プラグイン 1.7.4 を使用して実行されています。

以下は、複数のコンピューティングノードで機能とパフォーマンスの両方をテストするために Nvidia が提供する NCCL テストのサブセットです。

サポートされているインスタンス: P3dn, P4, P5

機能テスト

NCCL メッセージ転送マルチノードテスト

nccl_message_transfer は、NCCL OFI プラグインが期待どおりに動作していることを確認するための簡単なテストです。このテストでは、NCCL の接続の確立とデータ転送 API の機能を検証します。EFA で NCCL アプリケーションを実行と同時に、例に示すように mpirun への完全なパスを使用してください。クラスター内のインスタンス数と GPU に基づき、np と N のパラメータを変更します。詳細については、[AWS OFI NCCL ドキュメント](#)を参照してください。

次の nccl_message_transfer テストは汎用 CUDA xx.x バージョン用です。スクリプト内の CUDA バージョンを置き換えることで、Amazon EC2 インスタンスで使用可能な任意の CUDA バージョンに対するコマンドを実行できます。

```
$/opt/amazon/openmpi/bin/mpirun -n 2 -N 1 --hostfile hosts \
-x LD_LIBRARY_PATH=/usr/local/cuda-xx.x/efa/lib:/usr/local/cuda-xx.x/lib:/usr/
local/cuda-xx.x/lib64:/usr/local/cuda-xx.x:$LD_LIBRARY_PATH \
--mca btl tcp,self --mca btl_tcp_if_exclude lo,docker0 --bind-to none \
opt/aws-ofi-nccl/tests/nccl_message_transfer
```

出力は次のようになります。出力を確認して、EFA が OFI プロバイダとして使用されていることを確認します。

```
INFO: Function: nccl_net_ofi_init Line: 1069: NET/OFI Selected Provider is efa (found 4
 nics)
```



```
INFO: Function: nccl_net_ofi_init Line: 1160: NET/OFI Using transport protocol SENDRECV
INFO: Function: configure_ep_inorder Line: 261: NET/OFI Setting
  FI_OPT_EFA_SENDRECV_IN_ORDER_ALIGNED_128_BYTES not supported.
INFO: Function: configure_nccl_proto Line: 227: NET/OFI Setting NCCL_PROTO to "simple"
INFO: Function: main Line: 86: NET/OFI Process rank 1 started. NCCLNet device used on
  ip-172-31-13-179 is AWS Libfabric.
INFO: Function: main Line: 91: NET/OFI Received 4 network devices
INFO: Function: main Line: 111: NET/OFI Network supports communication using CUDA
  buffers. Dev: 3
INFO: Function: main Line: 118: NET/OFI Server: Listening on dev 3
INFO: Function: main Line: 131: NET/OFI Send connection request to rank 1
INFO: Function: main Line: 173: NET/OFI Send connection request to rank 0
INFO: Function: main Line: 137: NET/OFI Server: Start accepting requests
INFO: Function: main Line: 141: NET/OFI Successfully accepted connection from rank 1
INFO: Function: main Line: 145: NET/OFI Send 8 requests to rank 1
INFO: Function: main Line: 179: NET/OFI Server: Start accepting requests
INFO: Function: main Line: 183: NET/OFI Successfully accepted connection from rank 0
INFO: Function: main Line: 187: NET/OFI Rank 1 posting 8 receive buffers
INFO: Function: main Line: 161: NET/OFI Successfully sent 8 requests to rank 1
INFO: Function: main Line: 251: NET/OFI Got completions for 8 requests for rank 0
INFO: Function: main Line: 251: NET/OFI Got completions for 8 requests for rank 1
```

パフォーマンステスト

P4d.24xlarge でのマルチノード NCCL パフォーマンステスト

EFA で NCCL パフォーマンステストを確認するには、公式の [NCCL-Tests Repo](#) で実施可能な標準 NCCL パフォーマンステストを実施します。DLAMI には、CUDA XX.X 用に構築されたこのテストが付属しています。同様に、EFA を使用して独自のスクリプトを実行できます。

独自のスクリプトを作成する場合は、次のガイダンスを参照してください。

- EFA で NCCL アプリケーションを実行中に、例に示すように、mpirun への完全なパスを使用します。
- クラスター内のインスタンスの数と GPU の数に基づいて、パラメータ np と N を変更します。
- NCCL_DEBUG=INFO フラグを追加し、ログの EFA 使用状況が、[Selected Provider is EFA] になっていることを確認します。
- 検証のために解析するようにトレーニングログの場所を設定する

```
TRAINING_LOG="testEFA_$(date +"%N").log"
```

いずれかのメンバーノードでコマンド `watch nvidia-smi` を使用し、GPU の使用状況を監視します。次の `watch nvidia-smi` コマンドは汎用 CUDA `xx.x` バージョン用のコマンドで、インスタンスのオペレーティングシステムによって異なります。スクリプト内の CUDA バージョンを置き換えることで、Amazon EC2 インスタンスで使用可能な任意の CUDA バージョンに対するコマンドを実行できます。

- Amazon Linux 2:

```
$ /opt/amazon/openmpi/bin/mpirun -n 16 -N 8 \
-x NCCL_DEBUG=INFO -x --mca pml ^cm \
-x LD_LIBRARY_PATH=/usr/local/cuda-xx.x/efa/lib:/usr/local/cuda-xx.x/lib:/usr/
local/cuda-xx.x/lib64:/usr/local/cuda-xx.x:/opt/amazon/efa/lib64:/opt/amazon/openmpi/
lib64:$LD_LIBRARY_PATH \
--hostfile hosts --mca btl tcp,self --mca btl_tcp_if_exclude lo,docker0 --bind-to
none \
/usr/local/cuda-xx.x/efa/test-cuda-xx.x/all_reduce_perf -x NCCL_PROTO=simple -b 8 -e
1G -f 2 -g 1 -c 1 -n 100 | tee ${TRAINING_LOG}
```

- Ubuntu 20.04:

```
$ /opt/amazon/openmpi/bin/mpirun -n 16 -N 8 \
-x NCCL_DEBUG=INFO -x --mca pml ^cm \
-x LD_LIBRARY_PATH=/usr/local/cuda-xx.x/efa/lib:/usr/local/cuda-xx.x/lib:/usr/
local/cuda-xx.x/lib64:/usr/local/cuda-xx.x:/opt/amazon/efa/lib:/opt/amazon/openmpi/
lib:$LD_LIBRARY_PATH \
--hostfile hosts --mca btl tcp,self --mca btl_tcp_if_exclude lo,docker0 --bind-to
none \
/usr/local/cuda-xx.x/efa/test-cuda-xx.x/all_reduce_perf -x NCCL_PROTO=simple-b 8 -e
1G -f 2 -g 1 -c 1 -n 100 | tee ${TRAINING_LOG}
```

出力は次のようになります。

```
# nThread 1 nGpus 1 minBytes 8 maxBytes 1073741824 step: 2(factor) warmup iters: 5
iters: 100 agg iters: 1 validation: 1 graph: 0
#
# Using devices
# Rank 0 Group 0 Pid 9591 on ip-172-31-4-37 device 0 [0x10] NVIDIA A100-SXM4-40GB
# Rank 1 Group 0 Pid 9592 on ip-172-31-4-37 device 1 [0x10] NVIDIA A100-SXM4-40GB
# Rank 2 Group 0 Pid 9593 on ip-172-31-4-37 device 2 [0x20] NVIDIA A100-SXM4-40GB
# Rank 3 Group 0 Pid 9594 on ip-172-31-4-37 device 3 [0x20] NVIDIA A100-SXM4-40GB
# Rank 4 Group 0 Pid 9595 on ip-172-31-4-37 device 4 [0x90] NVIDIA A100-SXM4-40GB
```

```
# Rank 5 Group 0 Pid 9596 on ip-172-31-4-37 device 5 [0x90] NVIDIA A100-SXM4-40GB
# Rank 6 Group 0 Pid 9597 on ip-172-31-4-37 device 6 [0xa0] NVIDIA A100-SXM4-40GB
# Rank 7 Group 0 Pid 9598 on ip-172-31-4-37 device 7 [0xa0] NVIDIA A100-SXM4-40GB
# Rank 8 Group 0 Pid 10216 on ip-172-31-13-179 device 0 [0x10] NVIDIA A100-
SXM4-40GB
# Rank 9 Group 0 Pid 10217 on ip-172-31-13-179 device 1 [0x10] NVIDIA A100-
SXM4-40GB
# Rank 10 Group 0 Pid 10218 on ip-172-31-13-179 device 2 [0x20] NVIDIA A100-
SXM4-40GB
# Rank 11 Group 0 Pid 10219 on ip-172-31-13-179 device 3 [0x20] NVIDIA A100-
SXM4-40GB
# Rank 12 Group 0 Pid 10220 on ip-172-31-13-179 device 4 [0x90] NVIDIA A100-
SXM4-40GB
# Rank 13 Group 0 Pid 10221 on ip-172-31-13-179 device 5 [0x90] NVIDIA A100-
SXM4-40GB
# Rank 14 Group 0 Pid 10222 on ip-172-31-13-179 device 6 [0xa0] NVIDIA A100-
SXM4-40GB
# Rank 15 Group 0 Pid 10223 on ip-172-31-13-179 device 7 [0xa0] NVIDIA A100-
SXM4-40GB
ip-172-31-4-37:9591:9591 [0] NCCL INFO Bootstrap : Using ens32:172.31.4.37
ip-172-31-4-37:9591:9591 [0] NCCL INFO NET/Plugin: Failed to find ncclCollNetPlugin_v6
symbol.
ip-172-31-4-37:9591:9591 [0] NCCL INFO NET/Plugin: Failed to find ncclCollNetPlugin
symbol (v4 or v5).
ip-172-31-4-37:9591:9591 [0] NCCL INFO cudaDriverVersion 12020
NCCL version 2.18.5+cuda12.2
...
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Initializing aws-ofi-nccl 1.7.4-aws
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Using CUDA runtime version 11070
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Configuring AWS-specific options
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Using CUDA runtime version 11070
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Configuring AWS-specific options
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Setting provider_filter to efa
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Setting FI_EFA_FORK_SAFE environment
variable to 1
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Disabling NVLS support due to NCCL
version 21602
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Setting provider_filter to efa
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Setting FI_EFA_FORK_SAFE environment
variable to 1
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Disabling NVLS support due to NCCL
version 21602
```

```
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Running on p4d.24xlarge platform,
Setting NCCL_TOPO_FILE environment variable to /opt/aws-ofi-nccl/share/aws-ofi-nccl/
xml/p4d-24x1-topo.xml
```

```
...
```

```
-----some output truncated-----
#                                     out-of-place
#           in-place
#   size      count      type  redop  root   time  algbw  busbw #wrong
#   time      algbw    busbw #wrong
#   (us)      (GB/s)   (GB/s)
#           0           0     float  sum    -1    11.02  0.00  0.00  0
11.04  0.00  0.00  0
#           0           0     float  sum    -1    11.01  0.00  0.00  0
11.00  0.00  0.00  0
#           0           0     float  sum    -1    11.02  0.00  0.00  0
11.02  0.00  0.00  0
#           0           0     float  sum    -1    11.01  0.00  0.00  0
11.00  0.00  0.00  0
#           0           0     float  sum    -1    11.02  0.00  0.00  0
11.02  0.00  0.00  0
#           256         4     float  sum    -1    632.7  0.00  0.00  0
628.2  0.00  0.00  0
#           512         8     float  sum    -1    627.4  0.00  0.00  0
629.6  0.00  0.00  0
#           1024        16     float  sum    -1    632.2  0.00  0.00  0
631.7  0.00  0.00  0
#           2048        32     float  sum    -1    631.0  0.00  0.00  0
634.2  0.00  0.00  0
#           4096        64     float  sum    -1    623.3  0.01  0.01  0
633.6  0.01  0.01  0
#           8192       128     float  sum    -1    635.1  0.01  0.01  0
633.5  0.01  0.01  0
#           16384      256     float  sum    -1    634.8  0.03  0.02  0
637.0  0.03  0.02  0
#           32768      512     float  sum    -1    647.9  0.05  0.05  0
636.8  0.05  0.05  0
#           65536     1024     float  sum    -1    658.9  0.10  0.09  0
667.0  0.10  0.09  0
#           131072    2048     float  sum    -1    671.9  0.20  0.18  0
662.9  0.20  0.19  0
#           262144    4096     float  sum    -1    692.1  0.38  0.36  0
685.1  0.38  0.36  0
```

```

524288      8192      float      sum      -1      715.3      0.73      0.69      0
696.6      0.75      0.71      0
1048576     16384     float      sum      -1      734.6      1.43      1.34      0
729.2      1.44      1.35      0
2097152     32768     float      sum      -1      785.9      2.67      2.50      0
794.5      2.64      2.47      0
4194304     65536     float      sum      -1      837.2      5.01      4.70      0
837.6      5.01      4.69      0
8388608     131072    float      sum      -1      929.2      9.03      8.46      0
931.4      9.01      8.44      0
16777216    262144    float      sum      -1      1773.6     9.46      8.87      0
1772.8     9.46      8.87      0
33554432    524288    float      sum      -1      2110.2     15.90     14.91     0
2116.1     15.86     14.87     0
67108864    1048576   float      sum      -1      2650.9     25.32     23.73     0
2658.1     25.25     23.67     0
134217728   2097152   float      sum      -1      3943.1     34.04     31.91     0
3945.9     34.01     31.89     0
268435456   4194304   float      sum      -1      7216.5     37.20     34.87     0
7178.6     37.39     35.06     0
536870912   8388608   float      sum      -1      13680      39.24     36.79     0
13676      39.26     36.80     0
[ 1073741824 16777216   float      sum      -1      25645      41.87     39.25     0
25497      42.11     39.48     0 ] <- Used For Benchmark
...
# Out of bounds values : 0 OK
# Avg bus bandwidth    : 7.46044

```

検証テスト

EFA テストが有効な結果を返したことを検証するには、次のテストを使用して確認してください。

- EC2 インスタンスメタデータを使用してインスタンスタイプを取得します。

```

TOKEN=$(curl -X PUT "http://169.254.169.254/latest/api/token" -H "X-aws-ec2-metadata-token-ttl-seconds: 21600")
INSTANCE_TYPE=$(curl -H "X-aws-ec2-metadata-token: $TOKEN" -v http://169.254.169.254/latest/meta-data/instance-type)

```

- [パフォーマンステスト](#) を実行
- 以下のパラメータを設定する

```
CUDA_VERSION
```

```
CUDA_RUNTIME_VERSION
NCCL_VERSION
```

- 次に示すように、結果を検証します。

```
RETURN_VAL=`echo $?`
if [ ${RETURN_VAL} -eq 0 ]; then

    # Information on how the version come from logs
    #
    # ip-172-31-27-205:6427:6427 [0] NCCL INFO cudaDriverVersion 12020
    # NCCL version 2.16.2+cuda11.8
    # ip-172-31-27-205:6427:6820 [0] NCCL INFO NET/OFI Initializing aws-ofi-nccl
    1.7.1-aws
    # ip-172-31-27-205:6427:6820 [0] NCCL INFO NET/OFI Using CUDA runtime version
    11060

    # cudaDriverVersion 12020 --> This is max supported cuda version by nvidia
    driver
    # NCCL version 2.16.2+cuda11.8 --> This is NCCL version compiled with cuda
    version
    # Using CUDA runtime version 11060 --> This is selected cuda version

    # Validation of logs
    grep "NET/OFI Using CUDA runtime version ${CUDA_RUNTIME_VERSION}" ${TRAINING_LOG}
    || { echo "Runtime cuda text not found"; exit 1; }
    grep "NET/OFI Initializing aws-ofi-nccl" ${TRAINING_LOG} || { echo "aws-ofi-nccl
    is not working, please check if it is installed correctly"; exit 1; }
    grep "NET/OFI Configuring AWS-specific options" ${TRAINING_LOG} || { echo "AWS-
    specific options text not found"; exit 1; }
    grep "Using network AWS Libfabric" ${TRAINING_LOG} || { echo "AWS Libfabric text
    not found"; exit 1; }
    grep "busbw" ${TRAINING_LOG} || { echo "busbw text not found"; exit 1; }
    grep "Avg bus bandwidth " ${TRAINING_LOG} || { echo "Avg bus bandwidth text not
    found"; exit 1; }
    grep "NCCL version $NCCL_VERSION" ${TRAINING_LOG} || { echo "Text not found: NCCL
    version $NCCL_VERSION"; exit 1; }

    if [[ ${INSTANCE_TYPE} == "p4d.24xlarge" ]]; then
        grep "NET/AWS Libfabric/0/GDRDMA" ${TRAINING_LOG} || { echo "Text not found:
    NET/AWS Libfabric/0/GDRDMA"; exit 1; }
        grep "NET/OFI Selected Provider is efa (found 4 nics)" ${TRAINING_LOG} ||
        { echo "Selected Provider is efa text not found"; exit 1; }
```

```

    grep "aws-ofi-nccl/xml/p4d-24x1-topo.xml" ${TRAINING_LOG} || { echo "Topology
file not found"; exit 1; }
    elif [[ ${INSTANCE_TYPE} == "p4de.24xlarge" ]]; then
        grep "NET/AWS Libfabric/0/GDRDMA" ${TRAINING_LOG} || { echo "Avg bus
bandwidth text not found"; exit 1; }
        grep "NET/OFI Selected Provider is efa (found 4 nics)" ${TRAINING_LOG} ||
{ echo "Avg bus bandwidth text not found"; exit 1; }
        grep "aws-ofi-nccl/xml/p4de-24x1-topo.xml" ${TRAINING_LOG} || { echo
"Topology file not found"; exit 1; }
        elif [[ ${INSTANCE_TYPE} == "p5.48xlarge" ]]; then
            grep "NET/AWS Libfabric/0/GDRDMA" ${TRAINING_LOG} || { echo "Avg bus
bandwidth text not found"; exit 1; }
            grep "NET/OFI Selected Provider is efa (found 32 nics)" ${TRAINING_LOG} ||
{ echo "Avg bus bandwidth text not found"; exit 1; }
            grep "aws-ofi-nccl/xml/p5.48x1-topo.xml" ${TRAINING_LOG} || { echo "Topology
file not found"; exit 1; }
            elif [[ ${INSTANCE_TYPE} == "p3dn.24xlarge" ]]; then
                grep "NET/OFI Selected Provider is efa (found 4 nics)" ${TRAINING_LOG} ||
{ echo "Selected Provider is efa text not found"; exit 1; }
            fi
            echo "***** check_efa_nccl_all_reduce passed for cuda
version ${CUDA_VERSION} *****"
        else
            echo "***** check_efa_nccl_all_reduce failed for cuda
version ${CUDA_VERSION} *****"
        fi

```

- ベンチマークデータにアクセスするには、マルチノード all_reduce テストからのテーブル出力の最終行を解析します。

```

benchmark=$(sudo cat ${TRAINING_LOG} | grep '1073741824' | tail -n1 | awk -F " "
'{{print $12}}' | sed 's/ //' | sed 's/ 5e-07//')
if [[ -z "${benchmark}" ]]; then
    echo "benchmark variable is empty"
    exit 1
fi

echo "Benchmark throughput: ${benchmark}"

```

GPU のモニタリングおよび最適化

次のセクションでは、GPU の最適化とモニタリングオプションについて説明します。このセクションは、モニタリング、監督、事前処理、トレーニングが伴う一般的なワークフローと同様に編成されています。

- [モニタリング](#)
 - [で GPUs モニタリングする CloudWatch](#)
- [最適化](#)
 - [前処理](#)
 - [トレーニング](#)

モニタリング

DLAMI には、いくつかの GPU モニタリングツールがプリインストールされています。このガイドでは、ダウンロードしてインストールするために利用できるツールについても言及されています。

- [で GPUs モニタリングする CloudWatch](#) - GPU 使用状況統計を Amazon に報告するプリインストールされたユーティリティ CloudWatch。
- [nvidia-smi CLI](#) - 全体的な GPU コンピューティングおよびメモリ使用率をモニタリングするユーティリティ。これは (AWS Deep Learning AMI DLAMI) にプリインストールされています。
- [NVML C ライブラリ](#) - GPU モニタリングおよび管理機能に直接アクセスできる C ベースの API。これは、内部の nvidia-smi CLI によって使用され、DLAMI にプリインストールされています。また、それらの言語での開発を容易にするため、Python および Perl がバインドされています。DLAMI にプリインストールされた gpumon.py ユーティリティは、[nvidia-ml-py](#) の pynvml パッケージを使用しています。
- [NVIDIA DCGM](#) - クラスタ管理ツール。開発者ページにアクセスし、このツールをインストールして設定する方法を確認してください。

Tip

NVIDIA の開発者ブログで、DLAMI にインストールされている CUDA ツールの使用方法に関する最新情報を確認してください。

- [Nsight IDE と nvprof を使用して TensorCore 使用率をモニタリングします。](#)

で GPUs モニタリングする CloudWatch

GPU で DLAMI を使用すると、トレーニングや推論中にその使用状況を追跡する方法が必要になることがあります。これは、データパイプラインの最適化や深層学習ネットワークのチューニングに役立ちます。

を使用して GPU メトリクスを設定するには、次の 2 つの方法があります CloudWatch。

- [AWS CloudWatch エージェントでメトリクスを設定する \(推奨\)](#)
- [プリインストールされた gpumon.py スクリプトを使用してメトリクスを設定する](#)

AWS CloudWatch エージェントでメトリクスを設定する (推奨)

DLAMI を [統合 CloudWatch エージェント](#) と統合して GPU メトリクスを設定し、Amazon EC2 高速インスタンスでの GPU コプロセスの使用状況をモニタリングします。

DLAMI で [GPU メトリクス](#) を設定するには、次の 4 つの方法があります。

- [最小限の GPU メトリクスを設定する](#)
- [部分的な GPU メトリクスを設定する](#)
- [使用可能なすべての GPU メトリクスを設定する](#)
- [カスタム GPU メトリクスを設定する](#)

更新とセキュリティパッチの詳細については、「[AWS CloudWatch エージェントのセキュリティパッチ適用](#)」を参照してください。

前提条件

開始するには、インスタンスがメトリクスを にプッシュできるようにする Amazon EC2 インスタンス IAM アクセス許可を設定する必要があります CloudWatch。詳細な手順については、[CloudWatch 「エージェント で使用する IAM ロールとユーザーを作成する」](#) を参照してください。

最小限の GPU メトリクスを設定する

dlami-cloudwatch-agent@minimal systemd サービスを使用して最小限の GPU メトリクスを設定します。このサービスは以下のメトリクスを設定します。

- utilization_gpu
- utilization_memory

事前設定済みの最小限の GPU メトリクス向けの systemd サービスは以下の場所にあります。

```
/opt/aws/amazon-cloudwatch-agent/etc/dlami-amazon-cloudwatch-agent-minimal.json
```

以下のコマンドで systemd サービスを有効にして起動します。

```
sudo systemctl enable dlami-cloudwatch-agent@minimal
sudo systemctl start dlami-cloudwatch-agent@minimal
```

部分的な GPU メトリクスを設定する

dlami-cloudwatch-agent@partial systemd サービスを使用して部分的な GPU メトリクスを設定します。このサービスは以下のメトリクスを設定します。

- utilization_gpu
- utilization_memory
- memory_total
- memory_used
- memory_free

事前設定済みの部分的な GPU メトリクス向けの systemd サービスは以下の場所にあります。

```
/opt/aws/amazon-cloudwatch-agent/etc/dlami-amazon-cloudwatch-agent-partial.json
```

以下のコマンドで systemd サービスを有効にして起動します。

```
sudo systemctl enable dlami-cloudwatch-agent@partial
sudo systemctl start dlami-cloudwatch-agent@partial
```

使用可能なすべての GPU メトリクスを設定する

dlami-cloudwatch-agent@all systemd サービスを使用して使用可能なすべての GPU メトリクスを設定します。このサービスは以下のメトリクスを設定します。

- utilization_gpu
- utilization_memory
- memory_total

- memory_used
- memory_free
- temperature_gpu
- power_draw
- fan_speed
- pcie_link_gen_current
- pcie_link_width_current
- encoder_stats_session_count
- encoder_stats_average_fps
- encoder_stats_average_latency
- clocks_current_graphics
- clocks_current_sm
- clocks_current_memory
- clocks_current_video

事前設定済みの使用可能なすべての GPU メトリクス向けの systemd サービスは以下の場所にあります。

```
/opt/aws/amazon-cloudwatch-agent/etc/dlami-amazon-cloudwatch-agent-all.json
```

以下のコマンドで systemd サービスを有効にして起動します。

```
sudo systemctl enable dlami-cloudwatch-agent@all
sudo systemctl start dlami-cloudwatch-agent@all
```

カスタム GPU メトリクスを設定する

事前設定されたメトリクスが要件を満たさない場合は、カスタム CloudWatch エージェント設定ファイルを作成できます。

カスタム設定を作成する

カスタム設定ファイルを作成するには、[CloudWatch 「エージェント設定ファイルを手動で作成または編集する」](#)の詳細なステップを参照してください。

この例では、スキーマ定義が `/opt/aws/amazon-cloudwatch-agent/etc/amazon-cloudwatch-agent.json` にあると仮定します。

カスタムファイルを使用してメトリクスを設定する

次のコマンドを実行して、カスタムファイルに従って CloudWatch エージェントを設定します。

```
sudo /opt/aws/amazon-cloudwatch-agent/bin/amazon-cloudwatch-agent-ctl \
-a fetch-config -m ec2 -s -c \
file:/opt/aws/amazon-cloudwatch-agent/etc/amazon-cloudwatch-agent.json
```

AWS CloudWatch エージェントのセキュリティパッチ適用

新しくリリースされた DLAMIs には、利用可能な最新の AWS CloudWatch エージェントセキュリティパッチが設定されています。以下のセクションを参照して、お使いのオペレーティングシステムに応じて、現在の DLAMI を最新のセキュリティパッチで更新してください。

Amazon Linux 2

`yum` を使用して、Amazon Linux 2 DLAMI の最新の AWS CloudWatch エージェントセキュリティパッチを取得します。

```
sudo yum update
```

Ubuntu

Ubuntu で DLAMI の最新の AWS CloudWatch セキュリティパッチを取得するには、Amazon S3 ダウンロードリンクを使用して AWS CloudWatch エージェントを再インストールする必要があります。

```
wget https://s3.region.amazonaws.com/amazoncloudwatch-agent-region/ubuntu/arm64/latest/
amazon-cloudwatch-agent.deb
```

Amazon S3 ダウンロードリンクを使用して AWS CloudWatch エージェントをインストールする方法の詳細については、[「サーバーでの CloudWatch エージェントのインストールと実行」](#)を参照してください。

プリインストールされた `gpumon.py` スクリプトを使用してメトリクスを設定する

`gpumon.py` というユーティリティは、DLAMI にプリインストールされています。GPU メモリ CloudWatch、GPU 温度、GPU Power など、GPU ごとの使用状況のモニタリングと統合され、サ

ポートされます。スクリプトは、モニタリング対象データを定期的に送信します CloudWatch。スクリプトのいくつかの設定を変更 CloudWatch することで、に送信されるデータの粒度を設定できます。ただし、スクリプトを開始する前に、メトリクスを受信する CloudWatch ように を設定する必要があります。

で GPU モニタリングをセットアップして実行する方法 CloudWatch

1. IAM ユーザーを作成するか、既存のユーザーを変更して、メトリクスをに発行するためのポリシーを設定します CloudWatch。新しいユーザーを作成する場合、認証情報をメモしてください。次のステップで必要になります。

検索する IAM ポリシーは「cloudwatch:PutMetricData」です。追加されるポリシーは次のようになります。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Action": [
        "cloudwatch:PutMetricData"
      ],
      "Effect": "Allow",
      "Resource": "*"
    }
  ]
}
```

Tip

IAM ユーザーの作成と のポリシーの追加の詳細については CloudWatch、「」の [CloudWatch ドキュメント](#) を参照してください。

2. DLAMI で、[AWS 構成](#) を実行し、IAM ユーザー認証情報を指定します。

```
$ aws configure
```

3. 実行する前に、gpumon ユーティリティにいくつかの変更が必要になる場合があります。gpumon ユーティリティと README は次のコードブロックに定義された場所にあります。gpumon.py スクリプトの詳細については、[スクリプトの Amazon S3 の場所](#) を参照してください。

```
Folder: ~/tools/GPUCloudWatchMonitor
Files:  ~/tools/GPUCloudWatchMonitor/gpumon.py
        ~/tools/GPUCloudWatchMonitor/README
```

オプション:

- インスタンスが us-east-1 でない場合、gpumon.py でリージョンを変更します。
 - namespace やレポート期間などの他のパラメータを で変更します
CloudWatchstore_reso。
4. 現在、スクリプトでは Python 3 のみがサポートされています。希望するフレームワークの Python 3 環境を有効化するか、DLAMI の一般的な Python 3 環境を有効化します。

```
$ source activate python3
```

5. gpumon ユーティリティをバックグラウンドで実行します。

```
(python3)$ python gpumon.py &
```

6. ブラウザを開いて <https://console.aws.amazon.com/cloudwatch/> にアクセスし、メトリクスを選択します。名前空間 DeepLearning 「トレーニング」があります。

 Tip

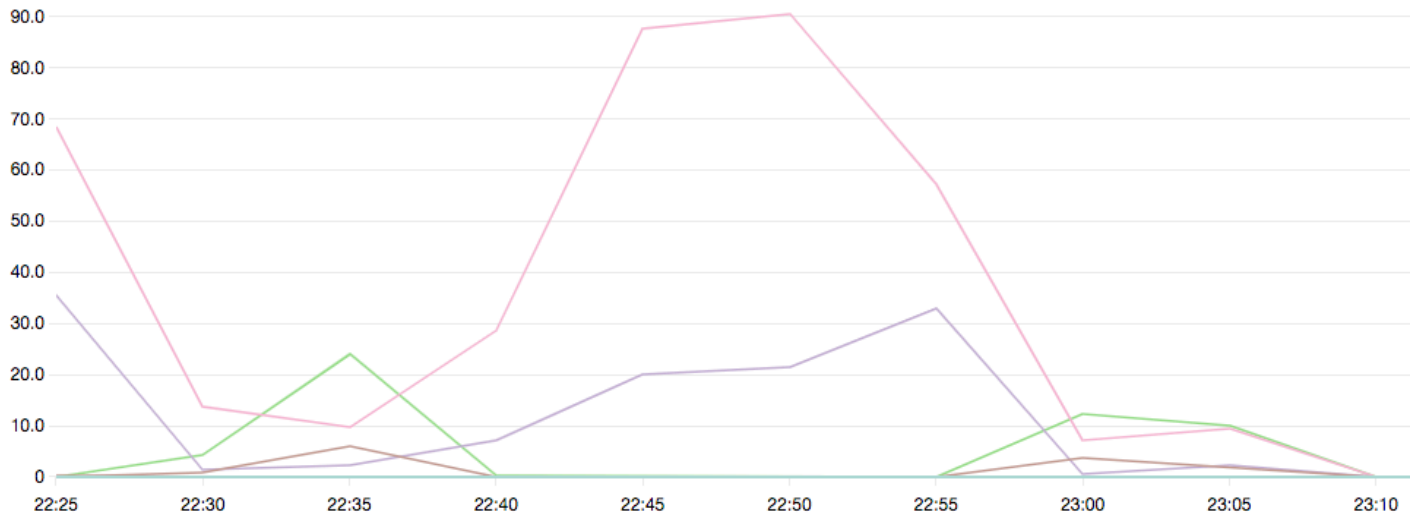
gpumon.py を変更することで名前空間を変更できます。store_reso を調整することで、レポートの間隔を調整することもできます。

以下は、p2.8xlarge インスタンスでトレーニングジョブをモニタリングする gpumon.py の実行に関する CloudWatch グラフレポートの例です。

GPU usage, Memory usage 

1h 3h 12h 1d 3d 1w custom

Various units



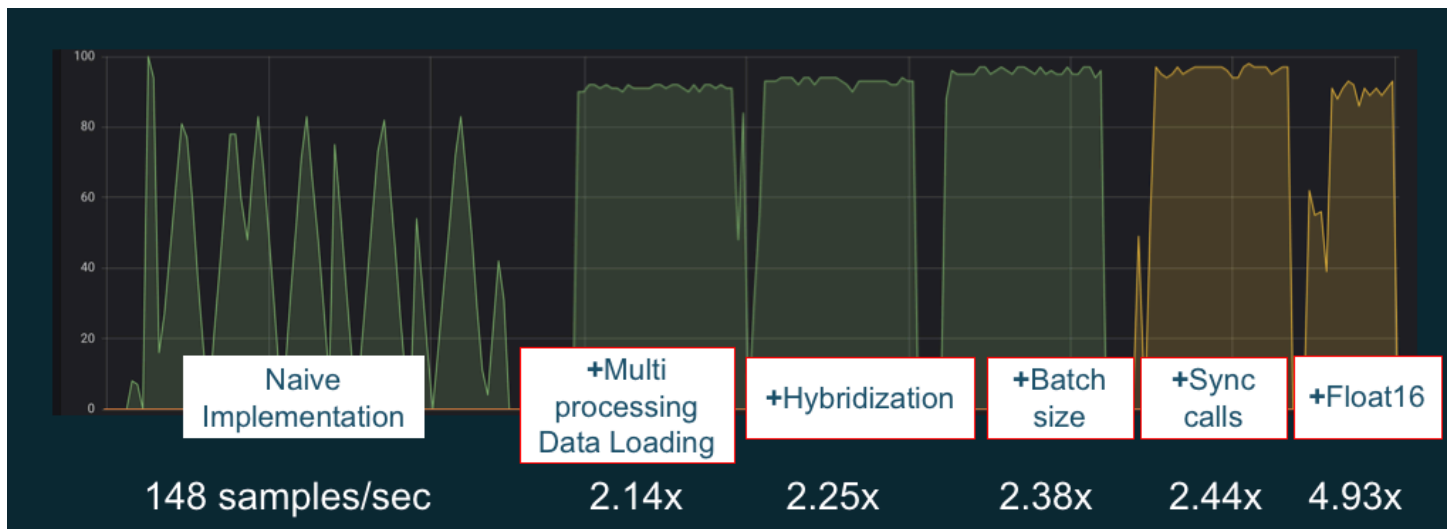
必要に応じて、GPU モニタリングおよび最適化に関する他のトピックも参照できます。

- [モニタリング](#)
 - [で GPU モニタリングする CloudWatch](#)
- [最適化](#)
 - [前処理](#)
 - [トレーニング](#)

最適化

GPU を最大限に活用するため、データパイプラインを最適化し、深層学習ネットワークをチューニングできます。次の図が示すように、ニューラルネットワークのネイティブまたは基本実装では、GPU が一貫性なく使用され、潜在能力が十分に引き出されない可能性があります。事前処理とデータのロードを最適化すると、CPU から GPU へのボトルネックを減らすことができます。ハイブリダイゼーションを使用し (フレームワークでサポートされている場合)、バッチサイズを調整してコールを同期することで、ニューラルネットワーク自体を調整することができます。ほとんどのフレームワークでは、多精度 (float16 または int8) トレーニングを使用することもできます。スループットの向上に劇的な効果が及ぶ可能性があります。

以下の図は、さまざまな最適化を適用した場合の累積的なパフォーマンス向上を示しています。結果は、処理するデータと最適化するネットワークによって異なります。



GPU パフォーマンスの最適化の例。グラフの出典: [MXNet と Gluon のパフォーマンスのヒント](#)

次のガイドでは、DLAMI を使用し、GPU パフォーマンスの向上に役立つオプションについて紹介します。

トピック

- [前処理](#)
- [トレーニング](#)

前処理

トランスフォーメーションやオーグメンテーションによるデータの事前処理は、多くの場合、CPU バウンド処理であり、パイプライン全体でボトルネックになる可能性があります。フレームワークには、画像処理用の組み込み演算子がありますが、DALI (データオーグメンテーションライブラリ) はフレームワークの組み込みオプションよりパフォーマンスが上回ることを実証しています。

- NVIDIA データオーグメンテーションライブラリ (DALI): DALI は、データオーグメンテーションを GPU にオフロードします。DLAMI にはプリインストールされていませんが、DLAMI や他の Amazon Elastic Compute Cloud インスタンスに DALI をインストールするか、サポートされているフレームワークコンテナをロードすることにより、DALI にアクセスできます。詳細については、NVIDIA ウェブサイトの [DALI プロジェクトページ](#) を参照してください。ユースケースの例とコードサンプルのダウンロードについては、[SageMaker 「前処理トレーニングパフォーマンスサンプル」](#) を参照してください。
- nvJPEG: C プログラマー向け GPU 加速 JPEG デコーダーライブラリ。1 つのイメージまたはバッチのデコードに加えて、深層学習に共通する後続のトランスフォーメーション操作もサポートさ

れています。nvJPEG には DALI が組み込まれています。または、[NVIDIA ウェブサイトの nvjpeg ページ](#)からダウンロードし、別個に使用することもできます。

必要に応じて、GPU モニタリングおよび最適化に関する他のトピックも参照できます。

- [モニタリング](#)
 - [で GPU モニタリングする CloudWatch](#)
- [最適化](#)
 - [前処理](#)
 - [トレーニング](#)

トレーニング

混合精度トレーニングでは、メモリの量が同じでより大規模なネットワークをデプロイしたり、単精度または倍精度ネットワークと比較してメモリの使用量を減らしたりすることができます。これにより、コンピューティングパフォーマンスが向上します。また、複数ノードに分散したトレーニングでは重要な要素である、少量かつ高速なデータ転送というメリットも得られます。混合精度トレーニングを利用するには、データキャストと損失スケールを調整する必要があります。混合精度をサポートするフレームワークでこれを行う方法について説明しているガイドを以下に示します。

- [NVIDIA Deep Learning SDK](#) - MXNet、PyTorch、および TensorFlow の混合精度実装について説明する NVIDIA ウェブサイトのドキュメント TensorFlow。

Tip

必ず、ウェブサイトで選択したフレームワークを確認し、「混合精度」または「fp16」を検索して最新の最適化手法を参照してください。以下の混合精度ガイドが役に立つ可能性があります。

- [TensorFlow \(ビデオ\) による混合精度トレーニング](#) - NVIDIA ブログサイト。
- [Mixed-precision training using float16 with MXNet](#) - MXNet ウェブサイト上のよくある質問記事。
- [NVIDIA Apex: NVIDIA ウェブサイトのブログ記事、で簡単に混合精度トレーニングを行うためのツール PyTorch。](#)

必要に応じて、GPU モニタリングおよび最適化に関する他のトピックも参照できます。

- [モニタリング](#)
 - [で GPU モニタリングする CloudWatch](#)
- [最適化](#)
 - [前処理](#)
 - [トレーニング](#)

DLAMI を備えた AWS Inferentia チップ

AWS Inferentia は、[によって設計されたカスタム機械学習チップ AWS](#) で、高性能な推論予測に使用できます。このチップを使用するには、Amazon Elastic Compute Cloud インスタンスをセットアップし、AWS Neuron ソフトウェア開発キット (SDK) を使用して Inferentia チップを呼び出します。お客様に最高の Inferentia エクスペリエンスを提供するために、Neuron が AWS Deep Learning AMI (DLAMI) に組み込まれています。

以下のトピックでは、Inferentia と DLAMI の使用を開始する方法について説明します。

内容

- [AWS Neuron を使用した DLAMI インスタンスの起動](#)
- [AWS Neuron での DLAMI の使用](#)

AWS Neuron を使用した DLAMI インスタンスの起動

最新の DLAMI は Inferentia AWS で使用できるようになり、AWS Neuron API パッケージが付属しています。DLAMI インスタンスを起動するには、[DLAMI の起動と設定](#)を参照してください。DLAMI を作成したら、以下のステップを使用して、AWS Inferentia チップと AWS Neuron リソースがアクティブであることを確認します。

内容

- [インスタンスの確認](#)
- [AWS Inferentia デバイスの識別](#)
- [リソースの使用状況の表示](#)
- [Neuron Monitor \(neuron-monitor\) の使用](#)
- [Neuron ソフトウェアのアップグレード](#)

インスタンスの確認

インスタンスを使用する前に、インスタンスが正しくセットアップされ、Neuron で設定されていることを確認します。

AWS Inferentia デバイスの識別

インスタンス上の Inferentia デバイスの数を確認するには、次のコマンドを使用します。

```
neuron-ls
```

インスタンスに Inferentia デバイスがアタッチされている場合、出力は次のようになります。

```
+-----+-----+-----+-----+-----+
| NEURON | NEURON | NEURON | CONNECTED | PCI      |
| DEVICE | CORES  | MEMORY | DEVICES   | BDF     |
+-----+-----+-----+-----+-----+
| 0      | 4      | 8 GB   | 1         | 0000:00:1c.0 |
| 1      | 4      | 8 GB   | 2, 0      | 0000:00:1d.0 |
| 2      | 4      | 8 GB   | 3, 1      | 0000:00:1e.0 |
| 3      | 4      | 8 GB   | 2         | 0000:00:1f.0 |
+-----+-----+-----+-----+-----+
```

示されている出力は Inf1.6xLarge インスタンスから取得されたものであり、次の列が含まれています。

- NEURON DEVICE: に割り当てられた論理 ID NeuronDevice。この ID は、異なる を使用するよう
に複数のランタイムを設定するときに使用されます NeuronDevices。
- NEURON CORES: NeuronCores に存在する の数 NeuronDevice。
- NEURON MEMORY: 内の DRAM メモリの量 NeuronDevice。
- コネクテッドデバイス: NeuronDevices に接続されたその他 NeuronDevice。
- PCI BDF: の PCI Bus Device Function (BDF) ID NeuronDevice。

リソースの使用状況の表示

neuron-top コマンドを使用して、NeuronCore および vCPU 使用率、メモリ使用率、ロードされたモデル、および Neuron アプリケーションに関する有用な情報を表示します。引数 neuron-top なしで起動すると、 を利用するすべての機械学習アプリケーションのデータが表示されます NeuronCores。

neuron-top

アプリケーションが 4 つの を使用している場合 NeuronCores、出力は次の図のようになります。



Neuron ベースの推論アプリケーションを監視および最適化するためのリソースの詳細については、[Neuron Tools](#) を参照してください。

Neuron Monitor (neuron-monitor) の使用

Neuron Monitor は、システムで実行されている Neuron ランタイムからメトリクスを収集し、収集したデータを JSON 形式で標準出力にストリーミングします。これらのメトリクスは、設定ファイルを指定して設定するメトリクスグループに編成されます。Neuron Monitor の詳細については、[Neuron Monitor User Guide](#) を参照してください。

Neuron ソフトウェアのアップグレード

DLAMI 内で Neuron SDK ソフトウェアを更新する方法については、AWS 「Neuron [セットアップガイド](#)」を参照してください。

次のステップ

[AWS Neuron での DLAMI の使用](#)

AWS Neuron での DLAMI の使用

AWS Neuron SDK の一般的なワークフローは、以前にトレーニングした機械学習モデルをコンパイルサーバーでコンパイルすることです。その後、実行のためにアーティファクトを Inf1 インスタンスに配布します。AWS Deep Learning AMI (DLAMI) には、Inferentia を使用する Inf1 インスタンスで推論をコンパイルして実行するために必要なものがすべてプリインストールされています。

以下のセクションでは、Inferentia とともに DLAMI を使用する方法について説明します。

内容

- [TensorFlow-Neuron と AWS Neuron Compiler の使用](#)
- [AWS Neuron TensorFlow Serving の使用](#)
- [MXNet -Neuron と AWS Neuron Compiler の使用](#)
- [MXNet-Neuron モデルサービングの使用](#)
- [PyTorch-Neuron と AWS Neuron Compiler の使用](#)

TensorFlow-Neuron と AWS Neuron Compiler の使用

このチュートリアルでは、AWS Neuron コンパイラを使用して Keras ResNet-50 モデルをコンパイルし、保存されたモデルとして SavedModel 形式でエクスポートする方法を示します。この形式は、一般的な TensorFlow モデル交換可能な形式です。また、入力例を使用して Inf1 インスタンスで推論を実行する方法も学習します。

Neuron SDK の詳細については、[AWS Neuron SDK のドキュメント](#)を参照してください。

内容

- [前提条件](#)
- [Conda 環境のアクティブ化](#)
- [Resnet50 コンパイル](#)
- [ResNet50 推論](#)

前提条件

このチュートリアルを使用する前に、[AWS Neuron を使用した DLAMI インスタンスの起動](#) の設定ステップを完了しておく必要があります。また、深層学習および DLAMI の使用にも精通している必要があります。

Conda 環境のアクティブ化

次のコマンドを使用して TensorFlow-Neuron conda 環境をアクティブ化します。

```
source activate aws_neuron_tensorflow_p36
```

現在の Conda 環境を終了するには、次のコマンドを実行します。

```
source deactivate
```

Resnet50 コンパイル

次の内容を含む **tensorflow_compile_resnet50.py** という Python スクリプトを作成します。この Python スクリプトは Keras ResNet50 モデルをコンパイルし、保存されたモデルとしてエクスポートします。

```
import os
import time
import shutil
import tensorflow as tf
import tensorflow.neuron as tfn
import tensorflow.compat.v1.keras as keras
from tensorflow.keras.applications.resnet50 import ResNet50
from tensorflow.keras.applications.resnet50 import preprocess_input

# Create a workspace
WORKSPACE = './ws_resnet50'
os.makedirs(WORKSPACE, exist_ok=True)

# Prepare export directory (old one removed)
model_dir = os.path.join(WORKSPACE, 'resnet50')
compiled_model_dir = os.path.join(WORKSPACE, 'resnet50_neuron')
shutil.rmtree(model_dir, ignore_errors=True)
```

```
shutil.rmtree(compiled_model_dir, ignore_errors=True)

# Instantiate Keras ResNet50 model
keras.backend.set_learning_phase(0)
model = ResNet50(weights='imagenet')

# Export SavedModel
tf.saved_model.simple_save(
    session          = keras.backend.get_session(),
    export_dir       = model_dir,
    inputs           = {'input': model.inputs[0]},
    outputs          = {'output': model.outputs[0]})

# Compile using Neuron
tfn.saved_model.compile(model_dir, compiled_model_dir)

# Prepare SavedModel for uploading to Inf1 instance
shutil.make_archive(compiled_model_dir, 'zip', WORKSPACE, 'resnet50_neuron')
```

次のコマンドを使用してモデルをコンパイルします。

```
python tensorflow_compile_resnet50.py
```

コンパイル処理には数分かかります。完了すると、出力は次のようになります。

```
...
INFO:tensorflow:fusing subgraph neuron_op_d6f098c01c780733 with neuron-cc
INFO:tensorflow:Number of operations in TensorFlow session: 4638
INFO:tensorflow:Number of operations after tf.neuron optimizations: 556
INFO:tensorflow:Number of operations placed on Neuron runtime: 554
INFO:tensorflow:Successfully converted ./ws_resnet50/resnet50 to ./ws_resnet50/
resnet50_neuron
...
```

コンパイル後、保存されたモデルは **ws_resnet50/resnet50_neuron.zip** に圧縮されます。以下のコマンドを使用して、モデルを解凍し、推論用のサンプルイメージをダウンロードします。

```
unzip ws_resnet50/resnet50_neuron.zip -d .
```

```
curl -O https://raw.githubusercontent.com/awslabs/mxnet-model-server/master/docs/images/kitten_small.jpg
```

ResNet50 推論

次の内容を含む **tensorflow_infer_resnet50.py** という Python スクリプトを作成します。このスクリプトは、以前にコンパイルされた推論モデルを使用して、ダウンロードしたモデルに対して推論を実行します。

```
import os
import numpy as np
import tensorflow as tf
from tensorflow.keras.preprocessing import image
from tensorflow.keras.applications import resnet50

# Create input from image
img_sgl = image.load_img('kitten_small.jpg', target_size=(224, 224))
img_arr = image.img_to_array(img_sgl)
img_arr2 = np.expand_dims(img_arr, axis=0)
img_arr3 = resnet50.preprocess_input(img_arr2)
# Load model
COMPILED_MODEL_DIR = './ws_resnet50/resnet50_neuron/'
predictor_inferentia = tf.contrib.predictor.from_saved_model(COMPILED_MODEL_DIR)
# Run inference
model_feed_dict={'input': img_arr3}
infa_rslts = predictor_inferentia(model_feed_dict);
# Display results
print(resnet50.decode_predictions(infa_rslts["output"], top=5)[0])
```

次のコマンドを使用して、モデルに対して推論を実行します。

```
python tensorflow_infer_resnet50.py
```

出力は次のようになります。

```
...
[('n02123045', 'tabby', 0.6918919), ('n02127052', 'lynx', 0.12770271), ('n02123159',
'tiger_cat', 0.08277027), ('n02124075', 'Egyptian_cat', 0.06418919), ('n02128757',
'snow_leopard', 0.009290541)]
```


次のステップ

[AWS Neuron TensorFlow Serving の使用](#)

AWS Neuron TensorFlow Serving の使用

このチュートリアルでは、保存されたモデルを TensorFlow Serving で使用するようにエクスポートする前に、グラフを作成し、AWS Neuron コンパイルステップを追加する方法を示します。TensorFlow Serving は、ネットワーク全体で推論をスケールアップできるサービスシステムです。Neuron TensorFlow Serving は、通常の TensorFlow Serving と同じ API を使用します。唯一の違いは、保存されたモデルを Inferentia AWS 用にコンパイルする必要があり、エントリポイントはという名前の別のバイナリであることです `tensorflow_model_server_neuron`。バイナリは、`/usr/local/bin/tensorflow_model_server_neuron` にあり、DLAMI にあらかじめインストールされています。

Neuron SDK の詳細については、[AWS Neuron SDK のドキュメント](#)を参照してください。

内容

- [前提条件](#)
- [Conda 環境のアクティブ化](#)
- [保存したモデルのコンパイルとエクスポート](#)
- [保存したモデルの提供](#)
- [モデルサーバーへの推論リクエストを生成する](#)

前提条件

このチュートリアルを使用する前に、[AWS Neuron を使用した DLAMI インスタンスの起動](#) の設定ステップを完了しておく必要があります。また、深層学習および DLAMI の使用にも精通している必要があります。

Conda 環境のアクティブ化

次のコマンドを使用して TensorFlow-Neuron conda 環境をアクティブ化します。

```
source activate aws_neuron_tensorflow_p36
```

現在の Conda 環境を終了する必要がある場合は、次のコマンドを実行します。

```
source deactivate
```

保存したモデルのコンパイルとエクスポート

以下の内容が含まれた Python スクリプト `tensorflow-model-server-compile.py` を作成します。このスクリプトは、Neuron を使用してグラフを作成し、コンパイルします。次に、コンパイルされたグラフを保存されたモデルとしてエクスポートします。

```
import tensorflow as tf
import tensorflow.neuron
import os

tf.keras.backend.set_learning_phase(0)
model = tf.keras.applications.ResNet50(weights='imagenet')
sess = tf.keras.backend.get_session()
inputs = {'input': model.inputs[0]}
outputs = {'output': model.outputs[0]}

# save the model using tf.saved_model.simple_save
modeldir = "./resnet50/1"
tf.saved_model.simple_save(sess, modeldir, inputs, outputs)

# compile the model for Inferentia
neuron_modeldir = os.path.join(os.path.expanduser('~'), 'resnet50_inf1', '1')
tf.neuron.saved_model.compile(modeldir, neuron_modeldir, batch_size=1)
```

次のコマンドを使用してモデルをコンパイルします。

```
python tensorflow-model-server-compile.py
```

出力は次のようになります。

```
...
INFO:tensorflow:fusing subgraph neuron_op_d6f098c01c780733 with neuron-cc
INFO:tensorflow:Number of operations in TensorFlow session: 4638
INFO:tensorflow:Number of operations after tf.neuron optimizations: 556
INFO:tensorflow:Number of operations placed on Neuron runtime: 554
```

```
INFO:tensorflow:Successfully converted ./resnet50/1 to /home/ubuntu/resnet50_inf1/1
```

保存したモデルの提供

モデルをコンパイルしたら、次のコマンドを使用して、保存したモデルに `tensorflow_model_server_neuron` バイナリを提供できます。

```
tensorflow_model_server_neuron --model_name=resnet50_inf1 \  
  --model_base_path=$HOME/resnet50_inf1/ --port=8500 &
```

出力は次のようになります。コンパイルされたモデルは、推論の準備のためにサーバーによって Inferentia デバイスの DRAM にステージングされます。

```
...  
2019-11-22 01:20:32.075856: I external/org_tensorflow/tensorflow/cc/saved_model/  
loader.cc:311] SavedModel load for tags { serve }; Status: success. Took 40764  
microseconds.  
2019-11-22 01:20:32.075888: I tensorflow_serving/servables/tensorflow/  
saved_model_warmup.cc:105] No warmup data file found at /home/ubuntu/resnet50_inf1/1/  
assets.extra/tf_serving_warmup_requests  
2019-11-22 01:20:32.075950: I tensorflow_serving/core/loader_harness.cc:87]  
Successfully loaded servable version {name: resnet50_inf1 version: 1}  
2019-11-22 01:20:32.077859: I tensorflow_serving/model_servers/  
server.cc:353] Running gRPC ModelServer at 0.0.0.0:8500 ...
```

モデルサーバーへの推論リクエストを生成する

次の内容で `tensorflow-model-server-infer.py` という Python スクリプトを作成します。このスクリプトは、サービスフレームワークである gRPC を介して推論を実行します。

```
import numpy as np  
import grpc  
import tensorflow as tf  
from tensorflow.keras.preprocessing import image  
from tensorflow.keras.applications.resnet50 import preprocess_input  
from tensorflow_serving.apis import predict_pb2  
from tensorflow_serving.apis import prediction_service_pb2_grpc  
from tensorflow.keras.applications.resnet50 import decode_predictions
```

```
if __name__ == '__main__':
    channel = grpc.insecure_channel('localhost:8500')
    stub = prediction_service_pb2_grpc.PredictionServiceStub(channel)
    img_file = tf.keras.utils.get_file(
        "./kitten_small.jpg",
        "https://raw.githubusercontent.com/awslabs/mxnet-model-server/master/docs/
images/kitten_small.jpg")
    img = image.load_img(img_file, target_size=(224, 224))
    img_array = preprocess_input(image.img_to_array(img)[None, ...])
    request = predict_pb2.PredictRequest()
    request.model_spec.name = 'resnet50_inf1'
    request.inputs['input'].CopyFrom(
        tf.contrib.util.make_tensor_proto(img_array, shape=img_array.shape))
    result = stub.Predict(request)
    prediction = tf.make_ndarray(result.outputs['output'])
    print(decode_predictions(prediction))
```

次のコマンドで gRPC を使用して、モデルの推論を実行します。

```
python tensorflow-model-server-infer.py
```

出力は次のようになります。

```
[(['n02123045', 'tabby', 0.6918919), ('n02127052', 'lynx', 0.12770271), ('n02123159',
'tiger_cat', 0.08277027), ('n02124075', 'Egyptian_cat', 0.06418919), ('n02128757',
'snow_leopard', 0.009290541)]]
```

MXNet-Neuron と AWS Neuron Compiler の使用

MXNet-Neuron コンパイル API は、AWS Inferentia デバイスで実行できるモデルグラフをコンパイルする方法を提供します。

この例では、API を使用して ResNet-50 モデルをコンパイルし、それを使用して推論を実行します。

Neuron SDK の詳細については、[AWS Neuron SDK のドキュメント](#)を参照してください。

内容

- [前提条件](#)

- [Conda 環境のアクティブ化](#)
- [Resnet50 コンパイル](#)
- [ResNet50 推論](#)

前提条件

このチュートリアルを使用する前に、[AWS Neuron を使用した DLAMI インスタンスの起動](#) の設定ステップを完了しておく必要があります。また、深層学習および DLAMI の使用にも精通している必要があります。

Conda 環境のアクティブ化

次のコマンドを使用して、MXNet-Neuron Conda 環境をアクティブにします。

```
source activate aws_neuron_mxnet_p36
```

現在の Conda 環境を終了するには、次のコマンドを実行します。

```
source deactivate
```

Resnet50 コンパイル

次の内容で `mxnet_compile_resnet50.py` という Python スクリプトを作成します。このスクリプトは、MXNet-Neuron コンパイル Python API を使用して ResNet-50 モデルをコンパイルします。

```
import mxnet as mx
import numpy as np

print("downloading...")
path='http://data.mxnet.io/models/imagenet/'
mx.test_utils.download(path+'resnet/50-layers/resnet-50-0000.params')
mx.test_utils.download(path+'resnet/50-layers/resnet-50-symbol.json')
print("download finished.")

sym, args, aux = mx.model.load_checkpoint('resnet-50', 0)
```

```
print("compile for inferentia using neuron... this will take a few minutes...")
inputs = { "data" : mx.nd.ones([1,3,224,224], name='data', dtype='float32') }

sym, args, aux = mx.contrib.neuron.compile(sym, args, aux, inputs)

print("save compiled model...")
mx.model.save_checkpoint("compiled_resnet50", 0, sym, args, aux)
```

次のコマンドを使用してモデルをコンパイルします。

```
python mxnet_compile_resnet50.py
```

コンパイルには数分かかります。コンパイルが終了すると、次のファイルが現在のディレクトリに表示されます。

```
resnet-50-0000.params
resnet-50-symbol.json
compiled_resnet50-0000.params
compiled_resnet50-symbol.json
```

ResNet50 推論

次の内容で **mxnet_infer_resnet50.py** という Python スクリプトを作成します。このスクリプトは、サンプルイメージをダウンロードし、それを使用して、コンパイルされたモデルを持つ推論を実行します。

```
import mxnet as mx
import numpy as np

path='http://data.mxnet.io/models/imagenet/'
mx.test_utils.download(path+'synset.txt')

fname = mx.test_utils.download('https://raw.githubusercontent.com/aws-labs/mxnet-model-server/master/docs/images/kitten_small.jpg')
img = mx.image.imread(fname)

# convert into format (batch, RGB, width, height)
img = mx.image.imresize(img, 224, 224)
# resize
```

```
img = img.transpose((2, 0, 1))
# Channel first
img = img.expand_dims(axis=0)
# batchify
img = img.astype(dtype='float32')

sym, args, aux = mx.model.load_checkpoint('compiled_resnet50', 0)
softmax = mx.nd.random_normal(shape=(1,))
args['softmax_label'] = softmax
args['data'] = img
# Inferentia context
ctx = mx.neuron()

exe = sym.bind(ctx=ctx, args=args, aux_states=aux, grad_req='null')
with open('synset.txt', 'r') as f:
    labels = [l.rstrip() for l in f]

exe.forward(data=img)
prob = exe.outputs[0].asnumpy()
# print the top-5
prob = np.squeeze(prob)
a = np.argsort(prob)[::-1]
for i in a[0:5]:
    print('probability=%f, class=%s' %(prob[i], labels[i]))
```

次のコマンドを使用して、コンパイルされたモデルで推論を実行します。

```
python mxnet_infer_resnet50.py
```

出力は次のようになります。

```
probability=0.642454, class=n02123045 tabby, tabby cat
probability=0.189407, class=n02123159 tiger cat
probability=0.100798, class=n02124075 Egyptian cat
probability=0.030649, class=n02127052 lynx, catamount
probability=0.016278, class=n02129604 tiger, Panthera tigris
```

次のステップ

[MXNet-Neuron モデルサービングの使用](#)

MXNet-Neuron モデルサービングの使用

このチュートリアルでは、事前にトレーニングされた MXNet モデルを使用して、マルチモデルサーバー (MMS) でリアルタイムのイメージ分類を実行する方法を学習します。MMS は、機械学習または深層学習フレームワークを使用してトレーニングされた深層学習モデルを提供するための柔軟性と easy-to-use ツールです。このチュートリアルには、AWS Neuron を使用したコンパイルステップと、MXNet を使用した MMS の実装が含まれています。

Neuron SDK の詳細については、[AWS Neuron SDK のドキュメント](#)を参照してください。

内容

- [前提条件](#)
- [Conda 環境のアクティブ化](#)
- [コード例のダウンロード](#)
- [モデルのコンパイル](#)
- [推論の実行](#)

前提条件

このチュートリアルを使用する前に、[AWS Neuron を使用した DLAMI インスタンスの起動](#) の設定ステップを完了しておく必要があります。また、深層学習および DLAMI の使用にも精通している必要があります。

Conda 環境のアクティブ化

次のコマンドを使用して、MXNet-Neuron Conda 環境をアクティブにします。

```
source activate aws_neuron_mxnet_p36
```

現在の Conda 環境を終了するには、次のコマンドを実行します。

```
source deactivate
```

コード例のダウンロード

この例を実行するには、次のコマンドを使用してコード例をダウンロードします。

```
git clone https://github.com/aws-labs/multi-model-server
cd multi-model-server/examples/mxnet_vision
```


モデルのコンパイル

次の内容で `multi-model-server-compile.py` という Python スクリプトを作成します。このスクリプトは、ResNet50 個のモデルを Inferentia デバイスタarget にコンパイルします。

```
import mxnet as mx
from mxnet.contrib import neuron
import numpy as np

path='http://data.mxnet.io/models/imagenet/'
mx.test_utils.download(path+'resnet/50-layers/resnet-50-0000.params')
mx.test_utils.download(path+'resnet/50-layers/resnet-50-symbol.json')
mx.test_utils.download(path+'synset.txt')

nn_name = "resnet-50"

#Load a model
sym, args, auxs = mx.model.load_checkpoint(nn_name, 0)

#Define compilation parameters# - input shape and dtype
inputs = {'data' : mx.nd.zeros([1,3,224,224], dtype='float32')}

# compile graph to inferentia target
csym, cargs, cauxs = neuron.compile(sym, args, auxs, inputs)

# save compiled model
mx.model.save_checkpoint(nn_name + "_compiled", 0, csym, cargs, cauxs)
```

モデルをコンパイルするには、次のコマンドを使用します。

```
python multi-model-server-compile.py
```

出力は次のようになります。

```
...
[21:18:40] src/nnvm/legacy_json_util.cc:209: Loading symbol saved by previous version
v0.8.0. Attempting to upgrade...
[21:18:40] src/nnvm/legacy_json_util.cc:217: Symbol successfully upgraded!
[21:19:00] src/operator/subgraph/build_subgraph.cc:698: start to execute partition
graph.
[21:19:00] src/nnvm/legacy_json_util.cc:209: Loading symbol saved by previous version
v0.8.0. Attempting to upgrade...
```

```
[21:19:00] src/nnvm/legacy_json_util.cc:217: Symbol successfully upgraded!
```

次の内容を含むと `signature.json` という名前のファイルを作成し、入力名と形状を設定します。

```
{
  "inputs": [
    {
      "data_name": "data",
      "data_shape": [
        1,
        3,
        224,
        224
      ]
    }
  ]
}
```

次のコマンドを使用して `synset.txt` ファイルを呼び出すことができます。このファイルは、ImageNet 予測クラスの名前のリストです。

```
curl -O https://s3.amazonaws.com/model-server/model_archive_1.0/examples/
squeeze_net_v1.1/synset.txt
```

`model_server_template` フォルダ内のテンプレートに続くカスタムサービスクラスを作成します。次のコマンドを使用して、テンプレートを現在の作業ディレクトリにコピーします。

```
cp -r ../model_service_template/* .
```

次のように `mxnet_model_service.py` モジュールを編集して、`mx.cpu()` コンテキストを `mx.neuron()` コンテキストに置き換えます。MXNet-Neuron は `NDArray` および `Glueon API` をサポートしていないため、`model_input` の不要なデータコピーをコメントアウトする必要があります。

```
...
self.mxnet_ctx = mx.neuron() if gpu_id is None else mx.gpu(gpu_id)
...
#model_input = [item.as_in_context(self.mxnet_ctx) for item in model_input]
```

次のコマンドを使用して、モデルアーカイバでモデルをパッケージ化します。

```
cd ~/multi-model-server/examples
model-archiver --force --model-name resnet-50_compiled --model-path mxnet_vision --
handler mxnet_vision_service:handle
```

推論の実行

マルチモデルサーバーを起動し、次のコマンドを使用して RESTful API を使用するモデルをロードします。neuron-rtd がデフォルト設定で実行されていることを確認します。

```
cd ~/multi-model-server/
multi-model-server --start --model-store examples > /dev/null # Pipe to log file if you
  want to keep a log of MMS
curl -v -X POST "http://localhost:8081/models?
  initial_workers=1&max_workers=4&synchronous=true&url=resnet-50_compiled.mar"
sleep 10 # allow sufficient time to load model
```

以下のコマンドでサンプルイメージを使用して推論を実行します。

```
curl -O https://raw.githubusercontent.com/awslabs/multi-model-server/master/docs/
  images/kitten_small.jpg
curl -X POST http://127.0.0.1:8080/predictions/resnet-50_compiled -T kitten_small.jpg
```

出力は次のようになります。

```
[
  {
    "probability": 0.6388034820556641,
    "class": "n02123045 tabby, tabby cat"
  },
  {
    "probability": 0.16900072991847992,
    "class": "n02123159 tiger cat"
  },
  {
    "probability": 0.12221276015043259,
    "class": "n02124075 Egyptian cat"
  },
  {
    "probability": 0.028706775978207588,
    "class": "n02127052 lynx, catamount"
  },
  {
```

```
"probability": 0.01915954425930977,  
"class": "n02129604 tiger, Panthera tigris"  
}  
]
```

テスト後にクリーンアップするには、RESTful API を使用して削除コマンドを発行し、次のコマンドを使用してモデルサーバーを停止します。

```
curl -X DELETE http://127.0.0.1:8081/models/resnet-50_compiled  
  
multi-model-server --stop
```

以下の出力が表示されます。

```
{  
  "status": "Model \"resnet-50_compiled\" unregistered"  
}  
Model server stopped.  
Found 1 models and 1 NCGs.  
Unloading 10001 (MODEL_STATUS_STARTED) :: success  
Destroying NCG 1 :: success
```

PyTorch-Neuron と AWS Neuron Compiler の使用

PyTorch-Neuron コンパイル API は、AWS Inferentia デバイスで実行できるモデルグラフをコンパイルする方法を提供します。

トレーニング済みモデルは、Inf1 インスタンスにデプロイする前に、Inferentia ターゲットにコンパイルする必要があります。次のチュートリアルでは、torchvision ResNet50 モデルをコンパイルし、保存された TorchScript モジュールとしてエクスポートします。このモデルは、推論を実行するために使用されます。

便宜上、このチュートリアルではコンパイルと推論の両方に Inf1 インスタンスを使用します。実際には、c5 インスタンスファミリーなどの別のインスタンスタイプを使用してモデルをコンパイルできます。その後、コンパイルされたモデルを Inf1 推論サーバーにデプロイする必要があります。詳細については、[AWS 「Neuron PyTorch SDK ドキュメント」](#) を参照してください。

内容

- [前提条件](#)
- [Conda 環境のアクティブ化](#)

- [Resnet50 コンパイル](#)
- [ResNet50 推論](#)

前提条件

このチュートリアルを使用する前に、[AWS Neuron を使用した DLAMI インスタンスの起動](#) の設定ステップを完了しておく必要があります。また、深層学習および DLAMI の使用にも精通している必要があります。

Conda 環境のアクティブ化

次のコマンドを使用して PyTorch-Neuron conda 環境をアクティブ化します。

```
source activate aws_neuron_pytorch_p36
```

現在の Conda 環境を終了するには、次のコマンドを実行します。

```
source deactivate
```

Resnet50 コンパイル

次の内容で **pytorch_trace_resnet50.py** という Python スクリプトを作成します。このスクリプトは、PyTorch-Neuron コンパイル Python API を使用して ResNet-50 モデルをコンパイルします。

Note

torchvision と torch パッケージのバージョン間には、torchvision モデルのコンパイル時に注意すべき依存関係があります。これらの依存関係ルールは、pip を介して管理できます。torchvision==0.6.1 は torch==1.5.1 リリースと対応し、torchvision==0.8.2 は torch==1.7.1 リリースに対応しています。

```
import torch
import numpy as np
import os
import torch_neuron
from torchvision import models
```

```
image = torch.zeros([1, 3, 224, 224], dtype=torch.float32)

## Load a pretrained ResNet50 model
model = models.resnet50(pretrained=True)

## Tell the model we are using it for evaluation (not training)
model.eval()
model_neuron = torch.neuron.trace(model, example_inputs=[image])

## Export to saved model
model_neuron.save("resnet50_neuron.pt")
```

コンパイルスクリプトを実行します。

```
python pytorch_trace_resnet50.py
```

コンパイルには数分かかります。コンパイルが完了すると、コンパイルされたモデルは `resnet50_neuron.pt` としてローカルディレクトリに保存されます。

ResNet50 推論

次の内容で **pytorch_infer_resnet50.py** という Python スクリプトを作成します。このスクリプトは、サンプルイメージをダウンロードし、それを使用して、コンパイルされたモデルを持つ推論を実行します。

```
import os
import time
import torch
import torch_neuron
import json
import numpy as np

from urllib import request

from torchvision import models, transforms, datasets

## Create an image directory containing a small kitten
os.makedirs("./torch_neuron_test/images", exist_ok=True)
request.urlretrieve("https://raw.githubusercontent.com/aws-labs/mxnet-model-server/master/docs/images/kitten_small.jpg",
                   "./torch_neuron_test/images/kitten_small.jpg")
```

```
## Fetch labels to output the top classifications
request.urlretrieve("https://s3.amazonaws.com/deep-learning-models/image-models/
imagenet_class_index.json","imagenet_class_index.json")
idx2label = []

with open("imagenet_class_index.json", "r") as read_file:
    class_idx = json.load(read_file)
    idx2label = [class_idx[str(k)][1] for k in range(len(class_idx))]

## Import a sample image and normalize it into a tensor
normalize = transforms.Normalize(
    mean=[0.485, 0.456, 0.406],
    std=[0.229, 0.224, 0.225])

eval_dataset = datasets.ImageFolder(
    os.path.dirname("./torch_neuron_test/"),
    transforms.Compose([
        transforms.Resize([224, 224]),
        transforms.ToTensor(),
        normalize,
    ])
)

image, _ = eval_dataset[0]
image = torch.tensor(image.numpy()[np.newaxis, ...])

## Load model
model_neuron = torch.jit.load( 'resnet50_neuron.pt' )

## Predict
results = model_neuron( image )

# Get the top 5 results
top5_idx = results[0].sort()[1][-5:]

# Lookup and print the top 5 labels
top5_labels = [idx2label[idx] for idx in top5_idx]

print("Top 5 labels:\n {}".format(top5_labels) )
```

次のコマンドを使用して、コンパイルされたモデルで推論を実行します。

```
python pytorch_infer_resnet50.py
```

出力は次のようになります。

```
Top 5 labels:  
['tiger', 'lynx', 'tiger_cat', 'Egyptian_cat', 'tabby']
```

ARM64 DLAMI

AWS ARM64 GPU DLAMIsは、深層学習ワークロードに高いパフォーマンスとコスト効率を提供するように設計されています。具体的には、G5g インスタンスタイプには Arm64-based [AWS Graviton2 プロセッサ](#) が搭載されています。これは、によってゼロから構築 AWS され、お客様がクラウドでワークロードを実行する方法に最適化されています。AWS ARM64 GPU DLAMIs は、Docker、NVIDIA Docker、NVIDIA Driver、CUDA、CuDNN、NCCL、TensorFlow および や などの一般的な機械学習フレームワークで事前設定されています PyTorch。

G5g インスタンスタイプでは、Graviton2 の価格およびパフォーマンスの利点を活用することで、GPU アクセラレーション機能を備えた x86 ベースのインスタンスと比較して、GPU アクセラレーションを利用した深層学習モデルを大幅に低いコストでデプロイできます。

ARM64 DLAMI を選択する

選択した ARM64 DLAMI を使用して [G5g インスタンス](#) を起動します。

DLAMI を起動する step-by-step 手順については、[「Launching and Configuring a DLAMI」を参照してください。](#)

最新の ARM64 DLAMIs [「DLAMI のリリースノート」](#) を参照してください。

使用を開始する

以下のトピックでは、ARM64 DLAMI の使用を開始する方法について説明します。

内容

- [ARM64 GPU PyTorch DLAMI の使用](#)

ARM64 GPU PyTorch DLAMI の使用

AWS Deep Learning AMI は Arm64 プロセッサベースの GPUs で使用できるようになり、用に最適化されています PyTorch。ARM64 GPU PyTorch DLAMI には [PyTorch](#)、深層学習のトレーニングと推論のユースケース [TorchServe](#) 用に、[TorchVision](#)、および で事前設定された Python 環境が含まれています。

内容

- [PyTorch Python 環境の検証](#)
- [でトレーニングサンプルを実行する PyTorch](#)
- [で推論サンプルを実行する PyTorch](#)

PyTorch Python 環境の検証

G5g インスタンスに接続し、次のコマンドを使用して Base Conda 環境を有効化します。

```
source activate base
```

コマンドプロンプトには、PyTorch TorchVision、およびその他のライブラリを含む基本 Conda 環境で作業していることを示す必要があります。

```
(base) $
```

PyTorch 環境のデフォルトのツールパスを確認します。

```
(base) $ which python
(base) $ which pip
(base) $ which conda
(base) $ which mamba
>>> import torch, torchvision
>>> torch.__version__
>>> torchvision.__version__
>>> v = torch.autograd.Variable(torch.randn(10, 3, 224, 224))
>>> v = torch.autograd.Variable(torch.randn(10, 3, 224, 224)).cuda()
>>> assert isinstance(v, torch.Tensor)
```

でトレーニングサンプルを実行する PyTorch

サンプル MNIST トレーニングジョブを実行します。

```
git clone https://github.com/pytorch/examples.git
cd examples/mnist
python main.py
```

出力は以下のようになります。

```
...
Train Epoch: 14 [56320/60000 (94%)]    Loss: 0.021424
Train Epoch: 14 [56960/60000 (95%)]    Loss: 0.023695
Train Epoch: 14 [57600/60000 (96%)]    Loss: 0.001973
Train Epoch: 14 [58240/60000 (97%)]    Loss: 0.007121
Train Epoch: 14 [58880/60000 (98%)]    Loss: 0.003717
Train Epoch: 14 [59520/60000 (99%)]    Loss: 0.001729
Test set: Average loss: 0.0275, Accuracy: 9916/10000 (99%)
```

で推論サンプルを実行する PyTorch

次のコマンドを使用して、事前トレーニング済みの densenet161 モデルをダウンロードし、を使用して推論を実行します TorchServe。

```
# Set up TorchServe
cd $HOME
git clone https://github.com/pytorch/serve.git
mkdir -p serve/model_store
cd serve

# Download a pre-trained densenet161 model
wget https://download.pytorch.org/models/densenet161-8d451a50.pth >/dev/null

# Save the model using torch-model-archiver
torch-model-archiver --model-name densenet161 \
  --version 1.0 \
  --model-file examples/image_classifier/densenet_161/model.py \
  --serialized-file densenet161-8d451a50.pth \
  --handler image_classifier \
  --extra-files examples/image_classifier/index_to_name.json \
  --export-path model_store

# Start the model server
torchserve --start --no-config-snapshots \
  --model-store model_store \
  --models densenet161=densenet161.mar &> torchserve.log
```

```
# Wait for the model server to start
sleep 30

# Run a prediction request
curl http://127.0.0.1:8080/predictions/densenet161 -T examples/image_classifier/
kitten.jpg
```

出力は以下のようになります。

```
{
  "tiger_cat": 0.4693363308906555,
  "tabby": 0.4633873701095581,
  "Egyptian_cat": 0.06456123292446136,
  "lynx": 0.0012828150065615773,
  "plastic_bag": 0.00023322898778133094
}
```

次のコマンドを使用して densenet161 モデルの登録を解除し、サーバーを停止します。

```
curl -X DELETE http://localhost:8081/models/densenet161/1.0
torchserve --stop
```

出力は以下のようになります。

```
{
  "status": "Model \"densenet161\" unregistered"
}
TorchServe has stopped.
```

推論

このセクションでは、DLAMI のフレームワークとツールを使用して推論を実行する方法についてチュートリアルを提供します。

推論ツール

- [TensorFlow サービス](#)

モデル提供

Deep Learning AMI with Conda にインストールされているモデル提供のオプションを以下に示します。いずれかのオプションをクリックして、使用方法をご覧ください。

トピック

- [TensorFlow サービス](#)
- [TorchServe](#)

TensorFlow サービス

[TensorFlow Serving](#) は、機械学習モデル用の柔軟で高性能なサービスシステムです。

tensorflow-serving-api は Deep Learning AMI with Conda にあらかじめインストールされています。MNIST モデルをトレーニング、エクスポート、処理するサンプルスクリプトが、~/examples/tensorflow-serving/ にあります。

これらの例を実行するには、まず Deep Learning AMI with Conda に接続し、TensorFlow 環境をアクティブ化します。

```
$ source activate tensorflow2_p310
```

ここで、処理するサンプルスクリプトのフォルダにディレクトリを変更します。

```
$ cd ~/examples/tensorflow-serving/
```

事前トレーニング済み Inception モデルを提供する

以下に示しているのは、Inception のような異なるモデルを提供するために試すことができる例です。原則として、提供可能なモデルとクライアントスクリプトが DLAMI にすでにダウンロードされている必要があります。

インセプションモデルによる推論の提供とテスト

1. モデルをダウンロードします。

```
$ curl -O https://s3-us-west-2.amazonaws.com/tf-test-models/INCEPTION.zip
```

2. モデルを展開します。

```
$ unzip INCEPTION.zip
```

3. ハスキーの画像をダウンロードします。

```
$ curl -O https://upload.wikimedia.org/wikipedia/commons/b/b5/Siberian_Husky_bi-eyed_Flickr.jpg
```

4. サーバーを起動します。Amazon Linux の場合、`model_base_path` に使用されるディレクトリを `/home/ubuntu` から `/home/ec2-user` に変更する必要があります。

```
$ tensorflow_model_server --model_name=INCEPTION --model_base_path=/home/ubuntu/examples/tensorflow-serving/INCEPTION/INCEPTION --port=9000
```

5. サーバーをフォアグラウンドで実行している場合、続行するには別のターミナルセッションを開始する必要があります。新しいターミナルを開き、TensorFlow を有効にします `source activate tensorflow2_p310`。次に、任意のテキストエディタを使用して、以下の内容のスクリプトを作成します。このスクリプトに `inception_client.py` という名前を付けます。このスクリプトは画像のファイル名をパラメータとして受け取り、事前トレーニング済みモデルから予測結果を取得します。

```
from __future__ import print_function

import grpc
import tensorflow as tf
import argparse

from tensorflow_serving.apis import predict_pb2
from tensorflow_serving.apis import prediction_service_pb2_grpc

parser = argparse.ArgumentParser(
    description='TF Serving Test',
    formatter_class=argparse.ArgumentDefaultsHelpFormatter
)
parser.add_argument('--server_address', default='localhost:9000',
                    help='Tenforflow Model Server Address')
parser.add_argument('--image', default='Siberian_Husky_bi-eyed_Flickr.jpg',
                    help='Path to the image')
args = parser.parse_args()

def main():
```

```
channel = grpc.insecure_channel(args.server_address)
stub = prediction_service_pb2_grpc.PredictionServiceStub(channel)
# Send request
with open(args.image, 'rb') as f:
    # See prediction_service.proto for gRPC request/response details.
    request = predict_pb2.PredictRequest()
    request.model_spec.name = 'INCEPTION'
    request.model_spec.signature_name = 'predict_images'

    input_name = 'images'
    input_shape = [1]
    input_data = f.read()
    request.inputs[input_name].CopyFrom(
        tf.make_tensor_proto(input_data, shape=input_shape))

    result = stub.Predict(request, 10.0) # 10 secs timeout
    print(result)

print("Inception Client Passed")

if __name__ == '__main__':
    main()
```

6. サーバーの場所、ポート、ハスキーの写真のファイル名をパラメータとして渡してスクリプトを実行します。

```
$ python3 inception_client.py --server=localhost:9000 --image Siberian_Husky_bi-
eyed_Flickr.jpg
```

MNIST モデルをトレーニングして提供する

このチュートリアルでは、モデルをエクスポートして、`tensorflow_model_server` アプリケーションで処理します。最終的に、サンプルのクライアントスクリプトを使用してモデルサーバーをテストできます。

MNIST モデルをトレーニングおよびエクスポートするスクリプトを実行します。スクリプトの唯一の引数として、モデルを保存するフォルダのロケーションを指定する必要があります。ここでは、`mnist_model` に入力するだけです。スクリプトによってフォルダが作成されます。

```
$ python mnist_saved_model.py /tmp/mnist_model
```

このスクリプトは出力に時間がかかることがあるため、少し待ちます。トレーニングが完了して最後にモデルがエクスポートされると、次の内容が表示されます。

```
Done training!  
Exporting trained model to mnist_model/1  
Done exporting!
```

次のステップでは、`tensorflow_model_server` を実行してエクスポートされたモデルを処理します。

```
$ tensorflow_model_server --port=9000 --model_name=mnist --model_base_path=/tmp/  
mnist_model
```

サーバーをテストするクライアントスクリプトが用意されています。

これ テストするには、新しいターミナルウィンドウを開く必要があります。

```
$ python mnist_client.py --num_tests=1000 --server=localhost:9000
```

その他の機能と例

TensorFlow Serving の詳細については、[TensorFlow のウェブサイト](#)を参照してください。

Amazon Elastic Inference で TensorFlow Serving を使用することもできます。<https://docs.aws.amazon.com/elastic-inference/latest/developerguide/what-is-ei.html> 詳細については、[TensorFlow Serving で Elastic Inference を使用する](#)方法に関するガイドを参照してください。

TorchServe

TorchServe は、Deep Learning AMI with Conda がプリインストールされた PyTorch. TorchServe comes からエクスポートされた深層学習モデルを提供するための柔軟なツールです。

の使用の詳細については TorchServe、[「ドキュメントのモデルサーバー PyTorch」](#)を参照してください。

トピック

で画像分類モデルを提供する TorchServe

このチュートリアルでは、を使用して画像分類モデルを提供する方法を示します TorchServe。が提供する DenseNet-161 モデルを使用します PyTorch。サーバーが実行されると、予測リクエストを

リッスンします。イメージをアップロードすると (この例では猫の画像)、サーバーはモデルがトレーニングされたクラスから最適な 5 つの一致するクラスの予測を返します。

でサンプル画像分類モデルを提供するには TorchServe

1. Deep Learning AMI with Conda v34 以降の Amazon Elastic Compute Cloud (Amazon EC2) インスタンスに接続します。
2. `pytorch_p310` 環境をアクティブ化します。

```
source activate pytorch_p310
```

3. TorchServe リポジトリをクローンし、モデルを保存するディレクトリを作成します。

```
git clone https://github.com/pytorch/serve.git
mkdir model_store
```

4. モデルアーカイバを使用してモデルをアーカイブします。 `extra-files` パラメータは TorchServe リポジトリのファイルを使用するため、必要に応じてパスを更新します。モデルアーカイブの詳細については、[「のトークモデルアーカイブ」](#)を参照してください [TorchServe](#)。

```
wget https://download.pytorch.org/models/densenet161-8d451a50.pth
torch-model-archiver --model-name densenet161 --version 1.0 --model-file ./
serve/examples/image_classifier/densenet_161/model.py --serialized-file
densenet161-8d451a50.pth --export-path model_store --extra-files ./serve/examples/
image_classifier/index_to_name.json --handler image_classifier
```

5. TorchServe を実行してエンドポイントを開始します。 `> /dev/null` を追加すると、ログ出力が抑止されます。

```
torchserve --start --ncs --model-store model_store --models densenet161.mar > /dev/
null
```

6. 子猫のイメージをダウンロードし、予測エンドポイントに送信します TorchServe。

```
curl -O https://s3.amazonaws.com/model-server/inputs/kitten.jpg
curl http://127.0.0.1:8080/predictions/densenet161 -T kitten.jpg
```


予測エンドポイントは次のような上位 5 つの予測に類似する予測を JSON で返します。ここでは、エジプシャンマウが含まれている可能性が 47%、続いてトラネコが含まれている可能性が 46% となっています。

```
{
  "tiger_cat": 0.46933576464653015,
  "tabby": 0.463387668132782,
  "Egyptian_cat": 0.0645613968372345,
  "lynx": 0.0012828196631744504,
  "plastic_bag": 0.00023323058849200606
}
```

7. テストが完了したら、以下のようにサーバーを停止します。

```
torchserve --stop
```

その他の例

TorchServe には、DLAMI インスタンスで実行できるさまざまな例があります。プロジェクト [TorchServe リポジトリの例のページ](#) で表示できます。

詳細情報

Docker のセットアップ方法 TorchServe や最新の TorchServe 機能など、詳細な TorchServe ドキュメントについては、「」の「[TorchServe プロジェクトページ](#)」を参照してください GitHub。

DLAMI のアップグレード

ここでは、DLAMI のアップグレードについての情報および DLAMI 上のソフトウェアを更新するためのヒントを示します。

トピック

- [DLAMI の新しいバージョンへのアップグレード](#)
- [ソフトウェア更新のヒント](#)
- [新しい更新の通知を受け取る](#)

DLAMI の新しいバージョンへのアップグレード

DLAMI のシステムイメージは、新しい深層学習フレームワークのリリース、CUDA やその他のソフトウェア更新、およびパフォーマンスチューニングを活用するため、定期的に更新されています。DLAMI をしばらく使用していて、更新を活用したい場合は、新しいインスタンスを起動する必要があります。また、データセット、チェックポイント、またはその他の重要なデータを手動で移行する必要があります。代わりに、Amazon EBS を使用してデータを保持し、新しい DLAMI にアタッチできます。このようにして頻繁にアップグレードし、データの移行にかかる時間を最小限に抑えることができます。

Note

DLAMI 間で Amazon EBS ボリュームをアタッチおよび移行する場合は、DLAMI と新しいボリュームの両方が同じアベイラビリティゾーンにある必要があります。

1. Amazon EC2 コンソールを使用して、新しい Amazon EBS ボリュームを作成します。詳細な手順については、[Amazon EBS ボリュームの作成](#)を参照してください。
2. 新しく作成した Amazon EBS ボリュームを既存の DLAMI にアタッチします。詳細な手順については、[Amazon EBS ボリュームのアタッチ](#)を参照してください。
3. データセット、チェックポイント、設定ファイルなどのデータを転送します。
4. DLAMI を起動します。詳細な手順については、[DLAMI を起動および設定する](#)を参照してください。
5. 古い DLAMI から Amazon EBS ボリュームをデタッチします。詳細な手順については、[Amazon EBS ボリュームのデタッチ](#)を参照してください。

6. Amazon EBS ボリュームを新しい DLAMI にアタッチします。ステップ 2 から手順に従ってボリュームをアタッチします。
7. データが新しい DLAMI で利用できることを確認したら、古い DLAMI を停止して削除します。詳細なクリーンアップ手順については、「[クリーンアップ](#)」を参照してください。

ソフトウェア更新のヒント

時として、DLAMI のソフトウェアを手動で更新する必要がある場合があります。一般的には、pip を使用して Python パッケージを更新することをお勧めします。また、Deep Learning AMI with Conda の Conda 環境内のパッケージの更新にも pip を使用してください。特定のフレームワークまたはソフトウェアのアップグレードとインストールの手順については、それぞれのウェブサイト参照してください。

Note

パッケージが正常に更新されることは保証できません。互換性のない依存関係を持つ環境にパッケージを更新しようとする、更新が失敗する可能性があります。そのような場合は、ライブラリのメンテナンス担当者に連絡して、パッケージの依存関係を更新できるかどうかを確認してください。また、更新を許可するように環境を変更することもできます。ただし、このような変更を行う場合は、既存のパッケージを削除または更新する必要があり、この環境の安定性が保証されなくなります。

AWS Deep Learning AMI には、多くの Conda 環境と多くのパッケージがプリインストールされています。プリインストールされているパッケージの数が多いため、互換性が保証されているパッケージのセットを見つけることは困難です。「環境が矛盾しています。パッケージプランをよく確認してください」という警告が表示されることがあります。DLAMI では、DLAMI が提供する環境が正しいことを保証しますが、ユーザーがインストールしたパッケージが正しく機能することは保証できません。

新しい更新の通知を受け取る

Note

AWS Deep Learning AMIs には、セキュリティパッチのリリース頻度が毎週あります。これらの増分セキュリティパッチについては、リリース通知が送信されますが、公式リリースノートには含まれない場合があります。

新しい DLAMI がリリースされるたびに通知を受け取ることができます。通知は、以下のトピックを使用して [Amazon SNS](#) で公開されます。

```
arn:aws:sns:us-west-2:767397762724:dlami-updates
```

メッセージは、新しい DLAMI が公開されたときにここに投稿されます。AMI のバージョン、メタデータ、およびリージョン AMI ID がメッセージに含まれます。

これらのメッセージは複数の方法で受信できます。次の方法を使用することをお勧めします。

1. [\[Amazon SNS コンソール\]](#) を開きます。
2. ナビゲーションバーで、必要に応じて AWS リージョンを米国西部 (オレゴン) に変更します。サブスクライブする SNS 通知が作成されたリージョンを選択する必要があります。
3. ナビゲーションペインで、サブスクリプション、サブスクリプションの作成 を選択します。
4. [サブスクリプションの作成] ダイアログボックスで、次の操作を行います。
 - a. トピック ARN の場合、次の Amazon リソースネーム (ARN) をコピーして貼り付けます。
arn:aws:sns:us-west-2:767397762724:dlami-updates
 - b. プロトコルで、[Amazon SQS , AWS Lambda, Email, Email-JSON] から 1 つ選択します。
 - c. エンドポイントには、通知の受信に使用するリソースの E メールアドレスまたは Amazon リソースネーム (ARN) を入力します。
 - d. [サブスクリプションを作成] を選択します。
5. AWS 「通知 - サブスクリプションの確認」という件名の確認メールが届きます。メールを開いて [サブスクリプションを確定] を選択して受信登録を完了します。

のセキュリティ AWS Deep Learning AMI

のクラウドセキュリティが最優先事項 AWS です。お客様は AWS、セキュリティを最も重視する組織の要件を満たすように構築されたデータセンターとネットワークアーキテクチャからメリットを得られます。

セキュリティは、AWS とユーザーの間で共有される責任です。[責任共有モデル](#)では、これをクラウドのセキュリティおよびクラウド内のセキュリティと説明しています。

- クラウドのセキュリティ — クラウドで AWS サービスを実行するインフラストラクチャを保護する責任 AWS は AWS にあります。AWS また、では、安全に使用できるサービスも提供しています。コンプライアンス [AWS プログラム](#) コンプライアンスプログラム の一環として、サードパーティーの監査者は定期的にセキュリティの有効性をテストおよび検証。DLAMI に適用されるコンプライアンスプログラムの詳細については、「[コンプライアンスプログラム AWS による対象範囲内のサービスコンプライアンスプログラム](#)」を参照してください。
- クラウドのセキュリティ — お客様の責任は、使用する AWS サービスによって決まります。また、お客様は、データの機密性、会社の要件、適用される法律や規制など、その他の要因についても責任を負います。

このドキュメントは、DLAMI を使用する際に責任共有モデルを適用する方法を理解するのに役立ちます 以下のトピックでは、セキュリティおよびコンプライアンスの目的を達成するために DLAMI を設定する方法を示します。また、DLAMI リソースのモニタリングや保護に役立つ他の AWS のサービスの使用方法についても説明します。

詳細については、「[Amazon EC2 におけるセキュリティ](#)」を参照してください。

トピック

- [でのデータ保護 AWS Deep Learning AMI](#)
- [での Identity and Access Management AWS Deep Learning AMI](#)
- [でのログ記録とモニタリング AWS Deep Learning AMI](#)
- [のコンプライアンス検証 AWS Deep Learning AMI](#)
- [の耐障害性 AWS Deep Learning AMI](#)
- [のインフラストラクチャセキュリティ AWS Deep Learning AMI](#)

でのデータ保護 AWS Deep Learning AMI

責任 AWS [共有モデル](#)、Deep Learning AWS AMI でのデータ保護に適用されます。このモデルで説明されているように、AWS はすべての を実行するグローバルインフラストラクチャを保護する責任があります AWS クラウド。お客様は、このインフラストラクチャでホストされているコンテンツに対する管理を維持する責任があります。また、使用する AWS のサービスのセキュリティ設定と管理タスクもユーザーの責任となります。データプライバシーの詳細については、「[データプライバシーのよくある質問](#)」を参照してください。欧州でのデータ保護の詳細については、AWS セキュリティブログに投稿された「[AWS 責任共有モデルおよび GDPR](#)」のブログ記事を参照してください。

データ保護の目的で、認証情報を保護し AWS アカウント、AWS IAM Identity Center または AWS Identity and Access Management (IAM) を使用して個々のユーザーを設定することをお勧めします。この方法により、それぞれのジョブを遂行するために必要な権限のみが各ユーザーに付与されます。また、次の方法でデータを保護することもお勧めします:

- 各アカウントで多要素認証 (MFA) を使用します。
- SSL/TLS を使用して AWS リソースと通信します。TLS 1.2 は必須であり TLS 1.3 がお勧めです。
- を使用して API とユーザーアクティビティのログ記録を設定します AWS CloudTrail。
- AWS 暗号化ソリューションと、内のすべてのデフォルトのセキュリティコントロールを使用します AWS のサービス。
- Amazon Macie などの高度なマネージドセキュリティサービスを使用します。これらは、Amazon S3 に保存されている機密データの検出と保護を支援します。
- コマンドラインインターフェイスまたは API AWS を介して にアクセスするときに FIPS 140-2 検証済みの暗号化モジュールが必要な場合は、FIPS エンドポイントを使用します。利用可能な FIPS エンドポイントの詳細については、「[連邦情報処理規格 \(FIPS\) 140-2](#)」を参照してください。

お客様の E メールアドレスなどの極秘または機密情報は、タグ、または名前フィールドなどの自由形式のテキストフィールドに配置しないことを強くお勧めします。これは、コンソール、API、または AWS のサービス SDK を使用して DLAMI AWS CLI または他の を使用する場合も同様です。AWS SDKs 名前に使用する自由記述のテキストフィールドやタグに入力したデータは、課金や診断ログに使用される場合があります。外部サーバーへの URL を提供する場合は、そのサーバーへのリクエストを検証するための認証情報を URL に含めないように強くお勧めします。

での Identity and Access Management AWS Deep Learning AMI

AWS Identity and Access Management (IAM) は、管理者が AWS リソースへのアクセスを安全に制御 AWS のサービス するのに役立つです。IAM 管理者は、誰が認証 (サインイン) され、DLAMI リソースを使用する認可を受ける (許可がある) ことができるかを制御します。IAM は、追加料金なしで AWS のサービス 使用できる です。

Identity and Access Management の詳細については、「[Amazon EC2 の Identity and Access Management](#)」を参照してください。

トピック

- [アイデンティティを使用した認証](#)
- [ポリシーを使用したアクセスの管理](#)
- [Amazon EMR での IAM](#)

アイデンティティを使用した認証

認証とは、ID 認証情報 AWS を使用して にサインインする方法です。として、IAM ユーザーとして AWS アカウントのルートユーザー、または IAM ロールを引き受けて認証 (にサインイン AWS) される必要があります。

ID ソースを介して提供された認証情報を使用して、フェデレーテッド ID AWS として にサインインできます。AWS IAM Identity Center (IAM Identity Center) ユーザー、会社のシングルサインオン認証、Google または Facebook の認証情報は、フェデレーテッド ID の例です。フェデレーテッド ID としてサインインする場合、IAM ロールを使用して、前もって管理者により ID フェデレーションが設定されています。フェデレーション AWS を使用して にアクセスすると、間接的にロールを引き受けることとなります。

ユーザーのタイプに応じて、AWS Management Console または AWS アクセスポータルにサインインできます。へのサインインの詳細については AWS、「ユーザーガイド」の「[へのサインイン AWS アカウント](#)方法AWS サインイン」を参照してください。

AWS プログラムで にアクセスする場合、 は Software Development Kit (SDK) とコマンドラインインターフェイス (CLI) AWS を提供し、認証情報を使用してリクエストに暗号で署名します。AWS ツールを使用しない場合は、リクエストに自分で署名する必要があります。推奨される方法を使用してリクエストを自分で署名する方法の詳細については、IAM [ユーザーガイドの API AWS リクエスト](#)の署名を参照してください。

使用する認証方法を問わず、追加セキュリティ情報の提供をリクエストされる場合もあります。例えば、AWS では、アカウントのセキュリティを強化するために多要素認証 (MFA) を使用することをお勧めします。詳細については、「AWS IAM Identity Center ユーザーガイド」の「[多要素認証](#)」および「IAM ユーザーガイド」の「[AWSでの多要素認証 \(MFA\) の使用](#)」を参照してください。

AWS アカウント ルートユーザー

を作成するときは AWS アカウント、アカウント内のすべての AWS のサービス およびリソースへの完全なアクセス権を持つ 1 つのサインインアイデンティティから始めます。この ID は AWS アカウント ルートユーザーと呼ばれ、アカウントの作成に使用した E メールアドレスとパスワードでサインインすることでアクセスできます。日常的なタスクには、ルートユーザーを使用しないことを強くお勧めします。ルートユーザーの認証情報は保護し、ルートユーザーでしか実行できないタスクを実行するときに使用します。ルートユーザーとしてサインインする必要があるタスクの完全なリストについては、IAM ユーザーガイドの「[ルートユーザー認証情報が必要なタスク](#)」を参照してください。

IAM ユーザーとグループ

[IAM ユーザー](#)は、単一のユーザーまたはアプリケーションに対して特定のアクセス許可 AWS アカウントを持つ内のアイデンティティです。可能であれば、パスワードやアクセスキーなどの長期的な認証情報を保有する IAM ユーザーを作成する代わりに、一時認証情報を使用することをお勧めします。ただし、IAM ユーザーでの長期的な認証情報が必要な特定のユースケースがある場合は、アクセスキーをローテーションすることをお勧めします。詳細については、IAM ユーザーガイドの「[長期的な認証情報を必要とするユースケースのためにアクセスキーを定期的にローテーションする](#)」を参照してください。

[IAM グループ](#)は、IAM ユーザーの集団を指定するアイデンティティです。グループとしてサインインすることはできません。グループを使用して、複数のユーザーに対して一度に権限を指定できます。多数のユーザーグループがある場合、グループを使用することで権限の管理が容易になります。例えば、IAMAdmins という名前のグループを設定して、そのグループに IAM リソースを管理する許可を与えることができます。

ユーザーは、ロールとは異なります。ユーザーは 1 人の人または 1 つのアプリケーションに一意に関連付けられますが、ロールはそれを必要とする任意の人が引き受けるようになっています。ユーザーには永続的な長期の認証情報がありますが、ロールでは一時的な認証情報が提供されます。詳細については、「IAM ユーザーガイド」の「[IAM ユーザー \(ロールではなく\) の作成が適している場合](#)」を参照してください。

IAM ロール

[IAM ロール](#)は、特定のアクセス許可 AWS アカウント を持つ 内のアイデンティティです。これは IAM ユーザーに似ていますが、特定のユーザーには関連付けられていません。ロール を切り替える AWS Management Console ことで、[IAM ロール](#)を一時的に引き受けることができます。ロール を引き受けるには、または AWS API AWS CLI オペレーションを呼び出すか、カスタム URL を使用します。ロールを使用する方法の詳細については、「IAM ユーザーガイド」の「[IAM ロールの使用](#)」を参照してください。

IAM ロールと一時的な認証情報は、次の状況で役立ちます:

- フェデレーションユーザーアクセス - フェデレーティッド ID に許可を割り当てるには、ロールを作成してそのロールの許可を定義します。フェデレーティッド ID が認証されると、その ID はロールに関連付けられ、ロールで定義されている許可が付与されます。フェデレーションの詳細については、「IAM ユーザーガイド」の「[Creating a role for a third-party Identity Provider](#)」(サードパーティーアイデンティティプロバイダー向けロールの作成)を参照してください。IAM Identity Center を使用する場合は、許可セットを設定します。アイデンティティが認証後にアクセスできるものを制御するため、IAM Identity Center は、権限セットを IAM のロールに関連付けます。アクセス許可セットの詳細については、「AWS IAM Identity Center ユーザーガイド」の「[アクセス許可セット](#)」を参照してください。
- 一時的な IAM ユーザー権限 - IAM ユーザーまたはロールは、特定のタスクに対して複数の異なる権限を一時的に IAM ロールで引き受けることができます。
- クロスアカウントアクセス - IAM ロールを使用して、自分のアカウントのリソースにアクセスすることを、別のアカウントの人物 (信頼済みプリンシパル) に許可できます。クロスアカウントアクセス権を付与する主な方法は、ロールを使用することです。ただし、一部の では AWS のサービス、(ロールをプロキシとして使用する代わりに) ポリシーをリソースに直接アタッチできます。クロスアカウントアクセスのロールとリソースベースのポリシーの違いについては、「[IAM ユーザーガイド](#)」の「[IAM でのクロスアカウントリソースアクセス](#)」を参照してください。
- クロスサービスアクセス — 一部の は、他の の機能 AWS のサービス を使用します AWS のサービス。例えば、あるサービスで呼び出しを行うと、通常そのサービスによって Amazon EC2 でアプリケーションが実行されたり、Amazon S3 にオブジェクトが保存されたりします。サービスでは、呼び出し元プリンシパルの許可、サービスロール、またはサービスリンクロールを使用してこれを行う場合があります。
- 転送アクセスセッション (FAS) — IAM ユーザーまたはロールを使用して でアクションを実行する場合 AWS、ユーザーはプリンシパルと見なされます。一部のサービスを使用する際に、アクションを実行することで、別のサービスの別のアクションがトリガーされることがあります。FAS は、 を呼び出すプリンシパルのアクセス許可を AWS のサービス、ダウンストリーム

サービス AWS のサービス へのリクエストのリクエストと組み合わせて使用します。FAS リクエストは、サービスが他の AWS のサービス またはリソースとのやり取りを完了する必要があるリクエストを受け取った場合にのみ行われます。この場合、両方のアクションを実行するためのアクセス許可が必要です。FAS リクエストを行う際のポリシーの詳細については、「[転送アクセスセッション](#)」を参照してください。

- サービスロール - サービスがユーザーに代わってアクションを実行するために引き受ける [IAM ロール](#)です。IAM 管理者は、IAM 内からサービスロールを作成、変更、削除できます。詳細については、「IAM ユーザーガイド」の「[AWS のサービスにアクセス許可を委任するロールの作成](#)」を参照してください。
- サービスにリンクされたロール - サービスにリンクされたロールは、にリンクされたサービスロールの一種です AWS のサービス。サービスは、ユーザーに代わってアクションを実行するロールを引き受けることができます。サービスにリンクされたロールは に表示され AWS アカウント、サービスによって所有されます。IAM 管理者は、サービスリンクロールの許可を表示できますが、編集することはできません。
- Amazon EC2 で実行されているアプリケーション - IAM ロールを使用して、EC2 インスタンスで実行され、AWS CLI または AWS API リクエストを行うアプリケーションの一時的な認証情報を管理できます。これは、EC2 インスタンス内でのアクセスキーの保存に推奨されます。AWS ロールを EC2 インスタンスに割り当て、そのすべてのアプリケーションで使用できるようにするには、インスタンスにアタッチされたインスタンスプロファイルを作成します。インスタンスプロファイルにはロールが含まれ、EC2 インスタンスで実行されるプログラムは一時的な認証情報を取得できます。詳細については、「IAM ユーザーガイド」の「[Amazon EC2 インスタンスで実行されるアプリケーションに IAM ロールを使用して許可を付与する](#)」を参照してください。

IAM ロールと IAM ユーザーのどちらを使用するかについては、「IAM ユーザーガイド」の「[\(IAM ユーザーではなく\) IAM ロールをいつ作成したら良いのか?](#)」を参照してください。

ポリシーを使用したアクセスの管理

でアクセスを制御する AWS には、ポリシーを作成し、AWS ID またはリソースにアタッチします。ポリシーは AWS、アイデンティティまたはリソースに関連付けられているときにアクセス許可を定義する のオブジェクトです。は、プリンシパル(ユーザー、ルートユーザー、またはロールセッション) がリクエストを行うときに、これらのポリシー AWS を評価します。ポリシーでの権限により、リクエストが許可されるか拒否されるかが決まります。ほとんどのポリシーは JSON ドキュメント AWS として に保存されます。JSON ポリシードキュメントの構造と内容の詳細については、「IAM ユーザーガイド」の「[JSON ポリシー概要](#)」を参照してください。

管理者は AWS JSON ポリシーを使用して、誰が何にアクセスできるかを指定できます。つまり、どのプリンシパルがどんなリソースにどんな条件でアクションを実行できるかということです。

デフォルトでは、ユーザーやロールに権限はありません。IAM 管理者は、リソースで必要なアクションを実行するための権限をユーザーに付与する IAM ポリシーを作成できます。その後、管理者はロールに IAM ポリシーを追加し、ユーザーはロールを引き継ぐことができます。

IAM ポリシーは、オペレーションの実行方法を問わず、アクションの許可を定義します。例えば、iam:GetRole アクションを許可するポリシーがあるとします。そのポリシーを持つユーザーは、AWS Management Console、AWS CLI または AWS API からロール情報を取得できます。

アイデンティティベースのポリシー

アイデンティティベースポリシーは、IAM ユーザー、ユーザーのグループ、ロールなど、アイデンティティにアタッチできる JSON 権限ポリシードキュメントです。これらのポリシーは、ユーザーとロールが実行できるアクション、リソース、および条件をコントロールします。アイデンティティベースのポリシーを作成する方法については、「IAM ユーザーガイド」の「[IAM ポリシーの作成](#)」を参照してください。

アイデンティティベースのポリシーは、さらにインラインポリシーまたはマネージドポリシーに分類できます。インラインポリシーは、単一のユーザー、グループ、またはロールに直接埋め込まれています。管理ポリシーは、内の複数のユーザー、グループ、ロールにアタッチできるスタンドアロンポリシーです AWS アカウント。管理ポリシーには、AWS 管理ポリシーとカスタマー管理ポリシーが含まれます。マネージドポリシーまたはインラインポリシーのいずれかを選択する方法については、「IAM ユーザーガイド」の「[マネージドポリシーとインラインポリシーの比較](#)」を参照してください。

リソースベースのポリシー

リソースベースのポリシーは、リソースに添付する JSON ポリシードキュメントです。リソースベースのポリシーには例として、IAM ロールの信頼ポリシー や Amazon S3 バケットポリシー があげられます。リソースベースのポリシーをサポートするサービスでは、サービス管理者はポリシーを使用して特定のリソースへのアクセスを制御できます。ポリシーがアタッチされているリソースの場合、指定されたプリンシパルがそのリソースに対して実行できるアクションと条件は、ポリシーによって定義されます。リソースベースのポリシーでは、[プリンシパルを指定する](#)必要があります。プリンシパルには、アカウント、ユーザー、ロール、フェデレーティッドユーザー、またはを含めることができます AWS のサービス。

リソースベースのポリシーは、そのサービス内にあるインラインポリシーです。リソースベースのポリシーでは、IAM の AWS マネージドポリシーを使用できません。

アクセスコントロールリスト (ACL)

アクセスコントロールリスト (ACL) は、どのプリンシパル (アカウントメンバー、ユーザー、またはロール) がリソースにアクセスするための許可を持つかを制御します。ACL はリソースベースのポリシーに似ていますが、JSON ポリシードキュメント形式は使用しません。

Amazon S3、AWS WAF、および Amazon VPC は、ACLs。ACL の詳細については、『Amazon Simple Storage Service デベロッパーガイド』の「[アクセスコントロールリスト \(ACL\) の概要](#)」を参照してください。

その他のポリシータイプ

AWS は、一般的ではない追加のポリシータイプをサポートします。これらのポリシータイプでは、より一般的なポリシータイプで付与された最大の権限を設定できます。

- **アクセス許可の境界** - アクセス許可の境界は、アイデンティティベースのポリシーによって IAM エンティティ (IAM ユーザーまたはロール) に付与できる権限の上限を設定する高度な機能です。エンティティにアクセス許可の境界を設定できます。結果として得られる権限は、エンティティのアイデンティティベースポリシーとそのアクセス許可の境界の共通部分になります。Principal フィールドでユーザーまたはロールを指定するリソースベースのポリシーでは、アクセス許可の境界は制限されません。これらのポリシーのいずれかを明示的に拒否した場合、権限は無効になります。アクセス許可の境界の詳細については、「IAM ユーザーガイド」の「[IAM エンティティのアクセス許可の境界](#)」を参照してください。
- **サービスコントロールポリシー (SCPs)** - SCPs は、の組織または組織単位 (OU) に対する最大アクセス許可を指定する JSON ポリシーです AWS Organizations。AWS Organizations は、AWS アカウント ビジネスが所有する複数の をグループ化して一元管理するためのサービスです。組織内のすべての機能を有効にすると、サービスコントロールポリシー (SCP) を一部またはすべてのアカウントに適用できます。SCP は、各 を含むメンバーアカウントのエンティティのアクセス許可を制限します AWS アカウントのルートユーザー。Organizations と SCP の詳細については、AWS Organizations ユーザーガイドの「[SCP の仕組み](#)」を参照してください。
- **セッションポリシー** - セッションポリシーは、ロールまたはフェデレーションユーザーの一時的なセッションをプログラムで作成する際にパラメータとして渡す高度なポリシーです。結果としてセッションの権限は、ユーザーまたはロールのアイデンティティベースポリシーとセッションポリシーの共通部分になります。また、リソースベースのポリシーから権限が派生する場合があります。これらのポリシーのいずれかを明示的に拒否した場合、権限は無効になります。詳細については、「IAM ユーザーガイド」の「[セッションポリシー](#)」を参照してください。

複数のポリシータイプ

1つのリクエストに複数のタイプのポリシーが適用されると、結果として作成される権限を理解するのがさらに難しくなります。複数のポリシータイプが関与する場合にリクエストを許可するかどうか AWS を決定する方法については、IAM ユーザーガイドの「[ポリシー評価ロジック](#)」を参照してください。

Amazon EMR での IAM

Amazon EMR AWS Identity and Access Management を使用して、ユーザー、AWS リソース、グループ、ロール、ポリシーを定義できます。また、これらのユーザーとロールがアクセスできる AWS サービスを制御できます。

Amazon EMR での IAM の使用の詳細については、[Amazon EMR での AWS Identity and Access Management](#) を参照してください。

でのログ記録とモニタリング AWS Deep Learning AMI

AWS Deep Learning AMI インスタンスには、GPU 使用状況統計を Amazon に報告するユーティリティなど、いくつかの GPU モニタリングツールが付属しています CloudWatch。詳細については、「[GPU のモニタリングおよび最適化](#)」と「[Amazon EC2 のモニタリング](#)」を参照してください。

使用状況の追跡

次の AWS Deep Learning AMI オペレーティングシステムディストリビューションには、ガインスタンスタイプ、インスタンス ID、DLAMI タイプ、および OS 情報を収集 AWS できるようにするコードが含まれています。DLAMI 内で使用されているコマンドに関する情報は収集も保持もされません。DLAMI に関するその他の情報は収集も保持もされません。

- Ubuntu 16.04
- Ubuntu 18.04
- Ubuntu 20.04
- Amazon Linux 2

DLAMI の使用状況の追跡をオプトアウトするには、起動時に Amazon EC2 インスタンスにタグを追加します。タグでは、キー OPT_OUT_TRACKING を使用し、関連した値を true に設定する必要があります。詳細については、[Amazon EC2 リソースにタグを付ける](#) を参照してください。

のコンプライアンス検証 AWS Deep Learning AMI

サードパーティーの監査者は、複数のコンプライアンスプログラム AWS Deep Learning AMI の一環としてのセキュリティと AWS コンプライアンスを評価します。サポートされているコンプライアンスプログラムの詳細については、「[Amazon EC2 のコンプライアンス検証](#)」を参照してください。

特定のコンプライアンスプログラムの対象となる AWS サービスのリストについては、「[コンプライアンスプログラム AWS による対象範囲内のサービスコンプライアンスプログラム](#)」を参照してください。一般的な情報については、[AWS 「コンプライアンスプログラム」](#)を参照してください。

を使用して、サードパーティーの監査レポートをダウンロードできます AWS Artifact。詳細については、「[アーティファクトでの AWS レポートのダウンロード](#)」を参照してください。

DLAMI を使用する際のお客様のコンプライアンス責任は、お客様のデータの機密性、貴社のコンプライアンス目的、適用される法律および規制によって決まります。では、コンプライアンスに役立つ以下のリソース AWS を提供しています。

- 「[セキュリティ&コンプライアンスクイックリファレンスガイド](#)」 - これらのデプロイガイドには、アーキテクチャ上の考慮事項の説明と、AWSでセキュリティとコンプライアンスに重点を置いたベースライン環境をデプロイするための手順が記載されています。
- [AWS コンプライアンスリソース](#) - このワークブックとガイドのコレクションは、お客様の業界や地域に適用される場合があります。
- 「[デベロッパーガイド](#)」の「[ルールによるリソースの評価](#)」 - この AWS Config サービスは、リソース設定が社内プラクティス、業界ガイドライン、および規制にどの程度準拠しているかを評価します。AWS Config
- [AWS Security Hub](#) - この AWS サービスは、内のセキュリティ状態を包括的に把握 AWS し、セキュリティ業界標準とベストプラクティスへの準拠を確認するのに役立ちます。

の耐障害性 AWS Deep Learning AMI

AWS グローバルインフラストラクチャは、AWS リージョンとアベイラビリティゾーンを中心に構築されています。AWS リージョンは、低レイテンシー、高スループット、および高度に冗長なネットワークで接続された、物理的に分離された複数のアベイラビリティゾーンを提供します。アベイラビリティゾーンでは、ゾーン間で中断することなく自動的にフェイルオーバーするアプリケーションとデータベースを設計および運用することができます。アベイラビリティゾーンは、従

来の単一または複数のデータセンターインフラストラクチャよりも可用性が高く、フォールトトレラントで、スケーラブルです。

AWS リージョンとアベイラビリティゾーンの詳細については、[AWS 「グローバルインフラストラクチャ」](#)を参照してください。

データの耐障害性とバックアップのニーズをサポートするための機能については、「[Amazon EC2 の耐障害性](#)」を参照してください。

のインフラストラクチャセキュリティ AWS Deep Learning AMI

のインフラストラクチャセキュリティ AWS Deep Learning AMI は Amazon EC2 によってサポートされています。詳細については、「[Amazon EC2 におけるインフラストラクチャセキュリティ](#)」を参照してください。

DLAMI に関する重要な変更点

よくある質問

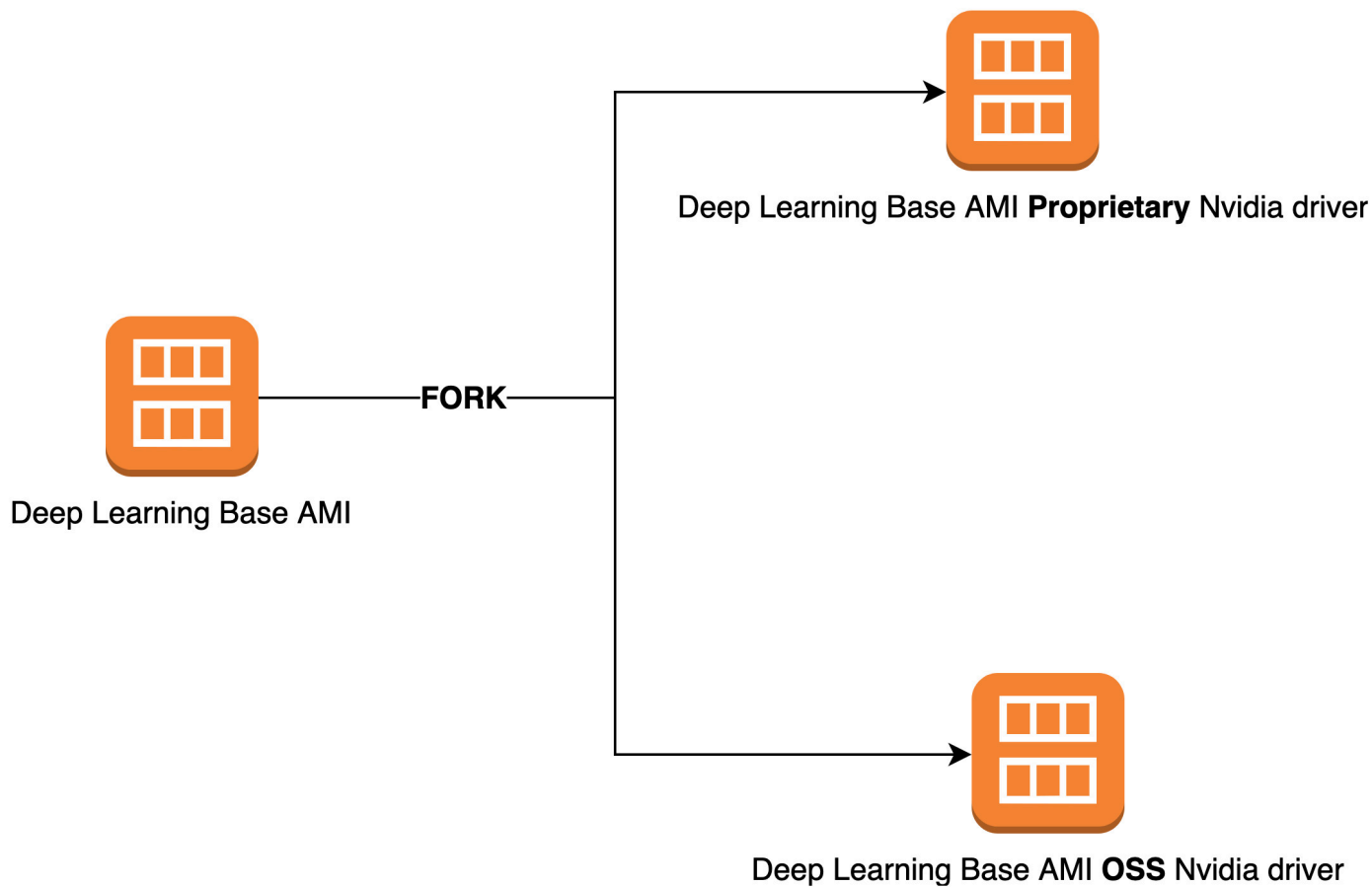
- [違いは何か](#)
- [なぜこの変更が必要なのか?](#)
- [この変更の影響を受けるのはどの DLAM ですか?](#)
- [これはあなたにとってどのような意味がありますか?](#)
- [新しい DLAMI をいつ使い始めるべきですか?](#)
- [新しい DLAMI では機能が低下することはありますか?](#)
- [DLC についてはどうですか?](#)

違いは何か

2023 年 11 月 15 日、AWS Deep Learning AMI (DLAMI) は次の 2 つのグループに分割されます。

- Nvidia 独自のドライバー (P3、P3dn、G3 をサポートするため) を使用する DLAM。
- Nvidia OSS ドライバー (G4dn、G5、P4、P5 をサポートするため) を使用する DLAM。

その結果、2 つのカテゴリのそれぞれについて、新しい名前と新しい AMI ID を使用して新しい DLAMI が作成されます。これらの DLAMI は交換できません。つまり、一方のグループの DLAMI は、もう一方のグループがサポートするインスタンスをサポートしません。たとえば、p5 をサポートする DLAMI は g3 をサポートせず、その逆も同様です。



なぜこの変更が必要なのか？

現在、NVIDIA GPU 用の DLAMI には NVIDIA 独自のカーネルドライバーが含まれています。しかし、最近、上流の Linux カーネルコミュニティは、NVIDIA GPU ドライバーなどのプロプライエタリなカーネルドライバーを他のカーネルドライバーと通信できないようにする変更を受け入れました。この変更により、GPU が分散トレーニングに EFA を効率的に使用できるようにするメカニズムである P4/P5 シリーズインスタンスの GPUDirect RDMA が無効になります。その結果、DLAMI は OpenRM ドライバー (NVIDIA オープンソースドライバー) を使用します。このドライバーは、G4dn、G5、P4、および P5 をサポートするオープンソースの EFA ドライバーとリンクされています。ただし、この OpenRM ドライバーは古いインスタンス (P3、G3 など) をサポートしません。したがって、両方のタイプのインスタンスをサポートする最新の高性能で安全な DLAMI を引き続き提供するために、DLAMI を 2 つのグループに分割します。1 つは OpenRM ドライバー (G4dn、G5、P4、P5 をサポート) で、もう 1 つは古い専用ドライバー (古いインスタンス P3、P3dn、G3 をサポート) です。

この変更の影響を受けるのはどの DLAM ですか？

すべての DLAMI がこの変更の影響を受けます。

これはあなたにとってどのような意味がありますか？

新しい DLAMI は、互換性のあるインスタンスタイプで実行されている限り、現在の DLAMI の機能、パフォーマンス、セキュリティを引き続き提供します。DLAMI を使用している場合は、各 DLAMI のリリースノート (こちらを参照) に記載されている互換性のあるインスタンスのいずれかで DLAMI が起動されていることを確認する必要があります。たとえば、次のような変更に対応する必要があります。

- 適切な CLI クエリで DLAMI を呼び出す (以下を参照)
- 互換性のあるインスタンスタイプでコンソールと CLI から DLAMIs を起動する

EC2 コンソールから DLAMI を起動する場合クイックスタート:各 DLAMI の説明には、EC2 コンソールでサポートされているインスタンスタイプが記載されています。DLAMI は互換性のあるインスタンスで起動する必要があります。

| | | |
|------------------------------|---|--------|
| ubuntu® Verified provider | Deep Learning Base GPU AMI (Ubuntu 20.04) 20231018 ami-05f9aedeafddcf112 (64-bit (x86)) Supported EC2 instances: P5*, P4*, P3*, G3*, G5*, G4dn. Release notes: https://aws.amazon.com/releases/notes/aws-deep-learning-base-gpu-ami-ubuntu-20-04/ | Select |
| ubuntu® Verified provider | Deep Learning AMI GPU PyTorch 2.0.1 (Ubuntu 20.04) 20231003 ami-005656037407fcf99 (64-bit (x86)) Supported EC2 instances: P5, P4d, P4de, P3, P3dn, G5, G4dn, G3. Release notes: https://docs.aws.amazon.com/dlami/latest/devguide/appendix-ami-release-notes.html | Select |
| ubuntu® Verified provider | Deep Learning AMI Neuron PyTorch 1.13 (Ubuntu 20.04) 20231003 ami-0f337e1c69255b2b6 (64-bit (x86)) Supported EC2 instances: Inf1, Trn1, Trn1n, Inf2. Release notes: https://docs.aws.amazon.com/dlami/latest/devguide/appendix-ami-release-notes.html | Select |

CLI を使用して DLAMIs を起動する場合は、クエリを変更する必要があります。例:

現在、すべてのインスタンス [P3、P3dn、G3、G4dn、G5、P4、P5] をサポートするベース DLAMI には、次の CLI クエリが使用されています。

```
aws ec2 describe-images --region us-east-1 --owners amazon \
--filters 'Name=name,Values=Deep Learning Base AMI (Amazon Linux 2) ??????????'
'Name=state,Values=available' \
--query 'reverse(sort_by(Images, &CreationDate))[1].ImageId' --output text
```

新しい CLI クエリは次のようになります。

P3、P3dn、および G3 をサポートするベース DLAMI の場合:

```
aws ec2 describe-images --region us-east-1 --owners amazon \  
--filters 'Name=name,Values=Deep Learning Base Proprietary Nvidia Driver AMI (Amazon Linux 2) Version ??.' 'Name=state,Values=available' \  
--query 'reverse(sort_by(Images, &CreationDate))[:1].ImageId' --output text
```

G4dn、G5、P4、および P5 をサポートするベース DLAMI の場合:

```
aws ec2 describe-images --region us-east-1 --owners amazon \  
--filters 'Name=name,Values=Deep Learning Base OSS Nvidia Driver AMI (Amazon Linux 2) Version ??.' 'Name=state,Values=available' \  
--query 'reverse(sort_by(Images, &CreationDate))[:1].ImageId' --output text
```

[新しい AMI に関する最新のリリースノートをこちらで参照してください。EC2 インスタンスで AMI を起動する方法については、こちらの手順を参照してください。](#)

新しい DLAMI をいつ使い始めるべきですか？

最新のフレームワーク、依存関係、パッチ、機能に対応するため、新しい DLAMI をできるだけ早く使い始める必要があります。[オプションで、2023 年 11 月 8 日より前にリリースされた Amazon Linux 2 DLAMI を使用している場合は、2023 年 11 月 30 日までその DLAMI のライブパッチ \(手順はこちら\) を継続することもできます。](#)

新しい DLAMI では機能が低下することはありますか？

いいえ、新しい DLAMI でも機能が失われることはありません。分割後の新しい DLAMI は、互換性のあるインスタンスで実行されている限り、分割前の古い DLAMI のすべての機能、パフォーマンス、セキュリティを引き続き提供します。DLAMI を 2 つのグループに分け、最新かつ高性能で安全な DLAMI をさまざまなインスタンスで引き続き提供できるようにしています。

DLC についてはどうですか？

DLC には NVIDIA ドライバーが含まれていないため、この変更による影響はありません。ただし、DLC は基盤となるインスタンスと互換性のある AMI で実行されるようにする必要があります。

関連情報

トピック

- [フォーラム](#)
- [関連する ブログ記事](#)
- [よくある質問](#)

フォーラム

- [フォーラム: AWS Deep Learning AMI](#)

関連する ブログ記事

- [Deep Learning AMI に関連する記事の更新されたリスト](#)
- [AWS Deep Learning AMI を起動する \(10 分\)](#)
- [最適化した TensorFlow 1.6 を使用して、Amazon EC2 C5 および P3 インスタンスでのトレーニングを高速化する](#)
- [機械学習の実践者向けの新しい AWS Deep Learning AMI](#)
- [新しいトレーニングコースの追加: AWS での機械学習および深層学習の入門](#)
- [AWS での深層学習のジャーニー](#)

よくある質問

- Q: DLAMI に関する製品発表を追跡するにはどうすればよいですか？

ここでは 2 通りの方法をお勧めします。

- [Deep Learning AMI に関連する記事の更新されたリスト](#)にあるブログカテゴリ「AWS Deep Learning AMIs」をブックマークしてください。
- [フォーラム: AWS Deep Learning AMI](#) を「監視」してください。
- Q: NVIDIA ドライバと CUDA はインストールされますか？

はい。一部の DLAMI は、バージョンが異なります。[Deep Learning AMI with Conda](#) には、すべての DLAMI の最新バージョンが含まれています。詳細については、[CUDA のインストール環境](#)およ

[びフレームワークのバインド](#) を参照してください。特定の AMI のリリースノート参照して、インストールされる内容を確認することもできます。

- Q. cuDNN はインストールされますか？

はい。

- Q. GPU が検出されたことや、GPU の現在のステータスを確認するには、どうすればよいですか？

`nvidia-smi` を実行します。これにより、1 つ以上の GPU (個数はインスタンスタイプによって決まります) と現在のメモリ使用率が表示されます。

- Q. 仮想環境は、自分専用にセットアップされますか？

はい。ただし、[Deep Learning AMI with Conda](#) の場合に限られます。

- Q: どのバージョンの Python がインストールされますか？

各 DLAMI には、Python 2 および 3 の両方がインストールされます。[Deep Learning AMI with Conda](#) には、各フレームワークの両方のバージョン用の環境が含まれます。

- Q. Keras はインストールされますか？

AMI によって異なります。[Deep Learning AMI with Conda](#) では、各フレームワークのフロントエンドとして Keras が使用できます。Keras のバージョンは、フレームワークでの Keras のサポート状況によって異なります。

- Q. 無料で利用できますか？

DLAMI はすべて無料で利用できます。ただし、選択したインスタンスタイプによっては、インスタンスを無料で利用できない場合もあります。詳細については、[DLAMI の価格設定](#) を参照してください。

- Q. フレームワークから CUDA エラーや GPU 関連のメッセージが表示されます。何が問題なのでしょうか？

使用しているインスタンスタイプを確認してください。多くの例やチュートリアルが機能するには、GPU が必要です。`nvidia-smi` を実行しても GPU が表示されない場合は、1 つ以上の GPU を含むインスタンスを使用して別の DLAMI を起動する必要があります。詳細については、[DLAMI のインスタンスタイプを選択する](#) を参照してください。

- Q. Docker は使用できますか？

Deep Learning AMI with Conda のバージョン 14 から、Docker は事前にインストールされています。GPU を利用するには、GPU インスタンス上で [nvidia-docker](#) を使用する必要があります。

- Q: Linux DLAMI を利用できるリージョンを教えてください。

| リージョン | コード |
|--------------------|----------------|
| 米国東部 (オハイオ) | us-east-2 |
| 米国東部 (バージニア北部) | us-east-1 |
| GovCloud | us-gov-west-1 |
| 米国西部 (北カリフォルニア) | us-west-1 |
| 米国西部 (オレゴン) | us-west-2 |
| 北京 (中国) | cn-north-1 |
| 寧夏 (中国) | cn-northwest-1 |
| アジアパシフィック (ムンバイ) | ap-south-1 |
| アジアパシフィック (ソウル) | ap-northeast-2 |
| アジアパシフィック (シンガポール) | ap-southeast-1 |
| アジアパシフィック (シドニー) | ap-southeast-2 |
| アジアパシフィック (東京) | ap-northeast-1 |
| カナダ (中部) | ca-central-1 |
| 欧州 (フランクフルト) | eu-central-1 |
| 欧州 (アイルランド) | eu-west-1 |

| リージョン | コード |
|------------|-----------|
| 欧州 (ロンドン) | eu-west-2 |
| 欧州 (パリ) | eu-west-3 |
| SA (サンパウロ) | sa-east-1 |

- Q: Windows DLAMI を利用できるリージョンを教えてください。

| リージョン | コード |
|--------------------|----------------|
| 米国東部 (オハイオ) | us-east-2 |
| 米国東部 (バージニア北部) | us-east-1 |
| GovCloud | us-gov-west-1 |
| 米国西部 (北カリフォルニア) | us-west-1 |
| 米国西部 (オレゴン) | us-west-2 |
| 北京 (中国) | cn-north-1 |
| アジアパシフィック (ムンバイ) | ap-south-1 |
| アジアパシフィック (ソウル) | ap-northeast-2 |
| アジアパシフィック (シンガポール) | ap-southeast-1 |
| アジアパシフィック (シドニー) | ap-southeast-2 |
| アジアパシフィック (東京) | ap-northeast-1 |

| リージョン | コード |
|--------------|--------------|
| カナダ (中部) | ca-central-1 |
| 欧州 (フランクフルト) | eu-central-1 |
| 欧州 (アイルランド) | eu-west-1 |
| 欧州 (ロンドン) | eu-west-2 |
| 欧州 (パリ) | eu-west-3 |
| SA (サンパウロ) | sa-east-1 |

DLAMI のリリースノート

Note

AWS Deep Learning AMIには、セキュリティパッチの定期的な夜間リリースがあります。これらの増分セキュリティパッチは公式リリースノートには含まれていません。

サポートされていないフレームワークリリースノートについては、[DLAMI サポートポリシーページ](#)を参照してください。

Base DLAMI

GPU

- X86
 - [AWS Deep Learning Base AMI \(Amazon Linux 2\)](#)
 - [AWS Deep Learning Base AMI \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning Base AMI \(Ubuntu 20.04\)](#)
- ARM64
 - [AWS Deep Learning Base ARM64 AMI \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning Base ARM64 AMI \(Amazon Linux 2\)](#)

AWS Neuron

- X86
 - [AWS Deep Learning Base AMI Neuron \(Amazon Linux 2\)](#)
 - [AWS Deep Learning Base AMI Neuron \(Ubuntu 20.04\)](#)

クアルコム

- X86
 - [AWS Deep Learning Base Qualcomm AMI \(Amazon Linux 2\)](#)

シングルフレーム DLAMI

PyTorch固有の AMI

GPU

- X86
 - [AWS Deep Learning AMI GPU PyTorch 2.3 \(Ubuntu 20.04\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 2.3 \(Amazon Linux 2\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 2.2 \(Ubuntu 20.04\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 2.2 \(Amazon Linux 2\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 1.13 \(Amazon Linux 2\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 1.13 \(Ubuntu 20.04\)](#)
- ARM64
 - [AWS Deep Learning ARM64 AMI GPU PyTorch 2.3 \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning ARM64 AMI GPU PyTorch 2.2 \(Ubuntu 20.04\)](#)

AWS Neuron

- X86
 - [AWS Deep Learning AMI Neuron PyTorch 1.13 \(Amazon Linux 2\)](#)
 - [AWS Deep Learning AMI Neuron PyTorch 1.13 \(Ubuntu 20.04\)](#)

TensorFlow固有の AMI

GPU

- X86
 - [AWS Deep Learning AMI GPU TensorFlow 2.16 \(Amazon Linux 2\)](#)
 - [AWS Deep Learning AMI GPU TensorFlow 2.16 \(Ubuntu 20.04\)](#)
 - [AWS Deep Learning AMI GPU TensorFlow 2.15 \(Amazon Linux 2\)](#)
 - [AWS Deep Learning AMI GPU TensorFlow 2.15 \(Ubuntu 20.04\)](#)
 - [AWS Deep Learning AMI GPU TensorFlow 2.13 \(Amazon Linux 2\)](#)
 - [AWS Deep Learning AMI GPU TensorFlow 2.13 \(Ubuntu 20.04\)](#)

AWS Neuron

- X86
 - [AWS Deep Learning AMI Neuron TensorFlow 2.10 \(Amazon Linux 2\)](#)
 - [AWS Deep Learning AMI Neuron TensorFlow 2.10 \(Ubuntu 20.04\)](#)

マルチフレームワーク DLAMI

Note

機械学習フレームワークを1つだけ使用する場合は、[シングルフレーム DLAMI](#)

GPU

- X86
 - [AWS Deep Learning AMI \(Amazon Linux 2\)](#)

AWS Neuron

- X86
 - [AWS Deep Learning AMI Neuron \(Ubuntu 22.04\)](#)

DLAMI の廃止に関する通知

次の表に、AWS Deep Learning AMI での機能の廃止に関する情報を示します。

| 廃止される機能 | 廃止日 | 廃止の通知 |
|--------------|-----------------|--|
| Ubuntu 16.04 | 2021 年 10 月 7 日 | Ubuntu Linux 16.04 LTS は、2021 年 4 月 30 日に 5 年間の LTS ウィンドウが終了し、ベンダーによってサポートされなくなりました。2021 年 10 月から、新規リリースでの Deep Learning Base AMI (Ubuntu 16.04) に対する更新はなくなりました。以前のリリースは、引き続き利用可能です。 |
| Amazon Linux | 2021 年 10 月 7 日 | Amazon Linux は 2020 年 12 月から サポート終了 です。2021 年 10 月から、新規リリースでの Deep Learning AMI (Amazon Linux) に対する更新はなくなりました。Deep Learning AMI (Amazon Linux) の以前のリリースは引き続き利用できます。 |
| Chainer | 2020 年 7 月 1 日 | Chainer は、2019 年 12 月時点での メジャーリリースのサポート終了 を公表しました。そのため、2020 年 7 月以降は、DLAMI に Chainer Conda 環境は含まれなくなります。これらの |

| 廃止される機能 | 廃止日 | 廃止の通知 |
|------------|-----------------|--|
| | | <p>環境を含む DLAMI の以前のリリースは引き続き利用できます。これらのフレームワークのオープンソースコミュニティによってセキュリティ修正が公開されている場合にのみ、これらの環境の更新を提供します。</p> |
| Python 3.6 | 2020 年 6 月 15 日 | <p>お客様の要望により、新しい TF/MX/PT リリースに対しては Python 3.7 に移行中です。</p> |
| Python 2 | 2020 年 1 月 1 日 | <p>Python オープンソースコミュニティは、Python 2 のサポートを正式に終了しました。</p> <p>TensorFlow、PyTorch、MXNet コミュニティもまた、TensorFlow 1.15、TensorFlow 2.1、PyTorch 1.4、MxNet 1.6.0の各リリースが、Python 2 をサポートする最後のリリースとなると発表しています。</p> |

AWS Deep Learning AMI デベロッパーガイドのドキュメント履歴

| 変更 | 説明 | 日付 |
|--|--|------------------|
| ARM64 DLAMI | は、Arm64 プロセッサベースの GPUs でイメージをサポートする AWS Deep Learning AMI ようになりました。 | 2021 年 11 月 29 日 |
| TensorFlow 2 | Deep Learning AMI with Conda に CUDA 10 を搭載した TensorFlow 2 が付属するようになりました。 | 2019 年 12 月 3 日 |
| AWS 推論 | Deep Learning AMI で、Inferentia AWS ハードウェアと AWS Neuron SDK がサポートされるようになりました。 | 2019 年 12 月 3 日 |
| 受信モデルでの TensorFlow Serving の使用 | Inception モデルで推論を使用する例が、Elastic Inference の有無にかかわらず、TensorFlow Serving に追加されました。 | 2018 年 11 月 28 日 |
| Elastic 推論 | Elastic 推論の前提条件と関連情報がセットアップガイドに追加されました。 | 2018 年 11 月 28 日 |
| ナイトリービルド PyTorch からのインストール | をアンインストールし PyTorch、Deep Learning AMI with Conda PyTorch に の夜間ビルドをインストールする方 | 2018 年 9 月 25 日 |

法を説明するチュートリアルが追加されました。

[Docker が DLAMI に事前にインストールされている](#)

Deep Learning AMI with Conda の v14 から、Docker および NVIDIA の GPU 対応 Docker バージョンが事前にインストールされています。

2018 年 9 月 25 日

[Conda チュートリアル](#)

MOTD の例では、最新リリースを反映するように更新されました。

2018 年 7 月 23 日

以前の更新:

次の表は、2018 年 7 月 AWS Deep Learning AMI 以前の の各リリースにおける重要な変更点を示しています。

| 変更 | 説明 | 日付 |
|---------------------------|---|-----------------|
| TensorFlow Horovod で | TensorFlow と Horovod ImageNet でトレーニングするためのチュートリアルを追加しました。 | 2018 年 6 月 6 日 |
| アップグレードガイド | アップグレードガイドを追加しました。 | 2018 年 5 月 15 日 |
| 新しいリージョンと新しい 10 分間チュートリアル | 追加された新しいリージョン: 米国西部 (北カリフォルニア)、南米、カナダ (中部)、欧州 (ロンドン)、および欧州 (パリ)。また、「Deep Learning AMI の使用開始」という 10 分間チュートリアルの最初のリリースが追加されました。 | 2018 年 4 月 26 日 |

| 変更 | 説明 | 日付 |
|--|---|------------------|
| Chainer チュートリアル | マルチ GPU、単一の GPU および CPU モデルで Chainer を使用するためのチュートリアルが追加されました。CUDA 統合は、複数のフレームワークで CUDA 8 から CUDA 9 にアップグレードされました。 | 2018 年 2 月 28 日 |
| Linux AMIs v3.0、MXNet Model Server、TensorFlow Serving、およびの紹介 TensorBoard | MXNet Model Server v0.1.5、TensorFlow Serving v1.4.0、および v0.4.0 を使用した新しいモデルおよびビジュアライゼーション提供機能を持つ Conda AMIs のチュートリアルを追加しました TensorBoard。AMI およびフレームワーク CUDA 機能については、Conda および CUDA の概要で説明されています。最新のリリースノートは https://aws.amazon.com/releasenotes/ に移動しました。 | 2018 年 1 月 25 日 |
| Linux AMI v2.0 | Base、Source、および Conda AMI が NCCL 2.1 で更新されました。ソース AMI と Conda AMIs MXNet v1.0、PyTorch 0.3.0、Keras 2.0.9 で更新されました。 | 2017 年 12 月 11 日 |
| 2 つの Windows AMI オプションを追加 | Windows 2012 R2 および 2016 AMI をリリース: AMI セレクションガイドおよびリリースノートに追加。 | 2017 年 11 月 30 日 |

| 変更 | 説明 | 日付 |
|---------------|----------------------------------|-------------|
| 初回のドキュメントリリース | 変更されたトピック/セクションへのリンクを含む変更の詳細な説明。 | 2017年11月15日 |

AWS 用語集

AWS の最新の用語については、「AWS の用語集リファレンス」の「[AWS 用語集](#)」を参照してください。

翻訳は機械翻訳により提供されています。提供された翻訳内容と英語版の間で齟齬、不一致または矛盾がある場合、英語版が優先します。