

AWS Well-Architected Framework

Games Industry Lens



Games Industry Lens: AWS Well-Architected Framework

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

Amazon's trademarks and trade dress may not be used in connection with any product or service that is not Amazon's, in any manner that is likely to cause confusion among customers, or in any manner that disparages or discredits Amazon. All other trademarks not owned by Amazon are the property of their respective owners, who may or may not be affiliated with, connected to, or sponsored by Amazon.

Table of Contents

Introduction	i
Custom lens availability	1
Definitions	2
Gaming platform	3
Game server	4
Game client	6
Messaging	6
Live game operations (Live Ops)	7
General design principles	9
Scenarios	10
Game hosting for real-time synchronous gameplay	10
Game server processes	11
Session-based game server hosting with serverless backend	12
Multi-Region and hybrid architecture for low-latency games	14
Game backends	16
Container-based	17
Serverless-based	19
Game production in the cloud (GPIC)	22
CI/CD	23
Workstations	25
Game analytics pipeline	26
Well-Architected pillars	29
Operational excellence	29
Design principles	29
Definition	30
Best practices	30
Resources	44
Security	45
Design principles	46
Definition	46
Best practices	46
Resources	58
Reliability	59
Design principles	59

Definition	59
Best practices	59
Resources	69
Performance efficiency	70
Design principles	70
Definition	70
Best practices	71
Resources	81
Cost optimization	82
Design principles	82
Definition	83
Best practices	84
Resources	94
Conclusion	96
Contributors	97
Document history	98
Notices	99
AWS Glossary	100

Introduction

Publication date: **November 19, 2021** ([Document history](#))

The [AWS Well-Architected Framework](#) helps cloud architects build secure, high-performing, resilient, and efficient infrastructure for their applications and workloads. Based on five pillars — operational excellence, security, reliability, performance efficiency, and cost optimization — Well-Architected provides a consistent approach for customers and AWS Partners to evaluate architectures, remediate risks, and implement designs that deliver business value.

In this Lens we focus on how to design, architect, and deploy your games workloads in the AWS Cloud. We define components, explore common workload scenarios, and outline design principles that help you to apply the Well-Architected Framework. We recommend that you begin designing your architecture by considering the best practices and questions from the [AWS Well-Architected Framework whitepaper](#). This document provides supplemental best practices for games industry customers.

This Lens specifies best practices that are intended to address the unique characteristics of building and operating games in the cloud based on our experience working with games industry developers and publishers around the world. It provides guidance on how to design and operate your environment so that it is cost optimized and scalable for fluctuations in global player demand. This Lens also provides guidance for securing your game infrastructure and tuning performance in order to deliver a positive player experience.

This document is intended for those in technology roles, such as chief technology officers (CTOs), game studio technical directors, architects, developers, and operations team members. After reading this document, you will understand AWS best practices and strategies to use when designing architectures for games.

Custom lens availability

Custom lenses extend the best practice guidance provided by AWS Well-Architected Tool. AWS WA Tool allows you to create your own [custom lenses](#), or to use lenses created by others that have been shared with you.

To determine if a custom lens is available for the lens described in this whitepaper, reach out to your Technical Account Manager (TAM), Solutions Architect (SA), or AWS Support.

Definitions

The AWS Well-Architected Framework is based on five pillars: operational excellence, security, reliability, performance efficiency, and cost optimization. AWS provides multiple core components that allow you to design state-of-the-art architectures for your game workloads. In this section, we will present an overview of key definitions.

For the purposes of this paper, a game architecture encompasses the backend technical infrastructure required to build and operate a game. Some games might not have social, multiplayer, or other online features and might not require the use of certain aspects of backend technical infrastructure that are described in this paper. Refer to [Scenarios](#) for a detailed discussion of the different types of workloads that are frequently deployed to support a game architecture.

The AWS Cloud is built around *AWS Regions* and *Availability Zones*. A Region is a physical location in the world where we have multiple Availability Zones. Availability Zones consist of one or more discrete data centers, each with redundant power, networking, and connectivity, housed in separate facilities. Depending on the characteristics of your game, you might want to deploy certain components of your game architecture into multiple Regions to improve performance for players, or to deliver customized experiences for players based on their location.

There are many different types of games and the backend technical infrastructure that is required to support a game will differ depending on the type of game being developed. Popular types of games include first person shooter (FPS), role playing (RPG), multiplayer online (MMO), Battle Royale (BR), sports games, and puzzle games. There are also different game interaction modes that influence the architecture of the game, such as turn-based or simultaneous play, with very different performance characteristics.

Games are developed to be played on one or more gaming platforms including desktop, web, mobile, and consoles. They also might feature newer interaction modes, such as augmented reality (AR), virtual reality (VR), and game streaming platforms. It is becoming a common trend for games to support cross-platform gameplay, which means that players can save their game progression and resume gameplay on other platforms. They can also initiate gameplay sessions with players on other platforms. Video game monetization enables game publishers to generate revenue using different strategies such as by advertising, digital and retail-based game purchases, in-game purchases of downloadable content (DLC), known as microtransactions, and through required paid subscriptions to play the game. Some of the most common key performance indicators (KPIs) in the games industry include daily active users (DAU), monthly active users (MAU), concurrent users

(CCU), session duration, cost per install (CPI), player lifetime value (LTV), and variations of average revenue per user (ARPU).

Topics

- [Gaming platform](#)
- [Game server](#)
- [Game client](#)
- [Messaging](#)
- [Live game operations \(Live Ops\)](#)

Gaming platform

Video games are developed to be played on a gaming platform that provides a player experience, which is usually comprised of client input controls, graphics, client software (known as the game client) and hardware, and in some cases platform-exclusive features to support gameplay.

Gaming platforms are generally delineated into the following categories:

- **Consoles** – Purpose-built entertainment systems that are designed for playing games, which include popular examples such as Sony PlayStation, Microsoft Xbox, and Nintendo Switch. Consoles provide the ability to play games by installing physical or digitally distributed game content onto the console hardware that is manufactured by the gaming platform provider. Using this definition, a console might be handheld, such as the Nintendo Switch, or stationary and intended to be used in a home entertainment scenario, such as an Xbox or PlayStation.
- **Personal computer (PC) games** – Games that are played using computer software that is installed onto a client machine that can be customized by the player. For this reason, PC gaming is popular among players because of the flexibility and control that it provides.
- **Web games** – Games that are designed to be played using a web browser, and which usually provide the benefit of enabling a player to access the game across platforms (cross-platform) by default as web browsers can be installed on many different types of devices.
- **Mobile games** – Games that are developed to be played on a mobile phone, usually with a smart phone operating system, such as iOS or Android. Mobile games are usually downloaded from a digital app store and installed onto the phone.

In addition to the above platforms, there are also nascent platforms that are still relatively new and growing, but have much smaller market share compared to the more predominant platforms. Examples of gaming platforms in this category include augmented reality (AR), virtual reality (VR), and game streaming, which is sometimes referred to as cloud gaming. *Game streaming* involves rendering the gameplay in the cloud and streaming to a thin client, typically a web browser. Game streaming allows a player to play a game that is entirely hosted remotely, typically in the cloud by a game streaming service provider. In game streaming, the player connects to a cloud-based game through a web browser or a thin client provided by the cloud gaming service provider (gaming platform).

Game server

Game servers represent one of the most important aspects of the compute infrastructure for your game. *Game servers*, sometimes referred to as dedicated game servers, are used when developing a multiplayer game or when server authoritative processing of gameplay events is required. The game server is at the center of the game architecture, serving as the location where the core logic executes, which includes managing player and game state as well as managing the interactions between the connected game clients and the game server. The game server is usually one of the most performance-sensitive aspects of a game architecture because it is responsible for processing the inputs from a player's game client and properly distributing it to any other connected players in real-time. A badly performing game server can impact the overall performance of the game experience. Therefore, it is important to ensure that the game server performance is optimized and has sufficient capacity, especially when the game is launched and during peak gameplay periods.

For the purposes of this document, a game server (or game server instance) refers to the compute resources, such as a virtual machine (VM), that hosts one or more game server processes. A game server process represents a single instance of your game server build hosting a game session, which is an instance of your running game that players can connect to via a player session. For this reason, we refer to *game server process* and *game session* interchangeably in this document, due to the implied one-to-one relationship between a game session and the game server process that is hosting it. In AWS, there are multiple options for compute resources to host game servers, all of which provide access to scalable cloud-based capacity through elastic provisioning of resources.

[Amazon EC2](#) provides cloud-based virtual servers, known as instances, with support for multiple versions of Linux and Windows. You can create instances and manage them directly like any other server or virtual machine (VM). Typically, multiple game server processes are deployed to an instance in order to improve efficiency and reduce costs. Amazon EC2 is a good choice for game servers if you want the most control over the compute infrastructure.

[Amazon Amazon GameLift](#) provides a fully-managed solution for dedicated game server hosting in the cloud as well as additional features such as matchmaking with Amazon GameLift FlexMatch. Amazon GameLift provides a layer of abstraction on top of Amazon EC2 to make game server management easier and is available in most AWS Regions so that you can host game servers close to players to reduce latency, achieve high availability, and significantly reduce costs by using Spot Instances. While Amazon GameLift can be integrated into existing game backends, it is especially useful for game developers who do not want to develop their own game server management and matchmaking solutions and want a solution that is managed by AWS and can scale as their game grows.

[Amazon Elastic Container Service \(Amazon ECS\)](#) is a fully managed container orchestration service that enables you to run Docker-based containers, and [Amazon Elastic Kubernetes Service \(Amazon EKS\)](#) enables you to run Docker-based containers that are built using Kubernetes. Using container technologies such as those provided by Amazon ECS and Amazon EKS can help you to improve the utilization of your compute infrastructure by allowing you to efficiently pack many game server processes or other game application instances into an EC2 instance. The use of containers can also improve developer productivity by enabling applications to be hosted using the same Docker image operating runtime that is used by developers on their local machines during development. You can further reduce operational overhead by using [AWS Fargate](#), which is a serverless compute platform for running containers and is compatible with both Amazon EKS and Amazon ECS. Fargate is best suited for use cases where you want to run game servers in containers without responsibility for operating the underlying instances that the containers run on.

[AWS Outposts](#) allows you to run AWS services in any data center or on-premises facility, which can enable games to run in on-premises environments and AWS using the same services to support a hybrid cloud adoption strategy. [AWS Local Zones](#) serve as extensions of AWS Regions to allow you to run your game servers and other latency-sensitive workloads closer to your players or development teams. Additionally, to reduce global network latency for your game servers, [AWS Global Accelerator](#) can be used to improve performance for player traffic to your game servers.

[AWS Lambda](#) is a serverless compute service that allows you to run code without provisioning or managing servers and is useful for asynchronous game server use cases, such as turn-based games. Lambda is also well suited for games with lightweight compute requirements, a small codebase, and where gameplay functionality can be designed using a stateless microservices architecture. Lambda functions run on an event-driven per-request basis, rather than running as part of a long-running game server process. Lambda provides the most runtime abstraction of the options described in this paper because the underlying application is provided out-of-the-box for developers to choose from to host their code.

When selecting your approach for game server hosting, you need to consider various requirements including operational overhead, legacy codebases, performance requirements, and scale. Amazon EC2 instances and containers are good options for legacy codebases, as they require the smallest amount of change to move to the cloud. EC2 instances also allow you to dedicate all the resources of a compute instance, while containers can make management and high utilization easier to achieve. Serverless functions offer the highest level of abstraction and allow you to define code that only runs in response to events, which can reduce costs.

Game client

The *game client* represents the software, hardware device, or both, that the player uses for playing a game. The game client provides the software for translating the player's inputs into messages that are sent to a server for processing. It's also responsible for handling the incoming responses from the server and rendering any outputs, such as graphics, to the player. In real-time networked multiplayer games, the game client usually maintains a persistent network connection to a game server for the duration of a gameplay session to reduce network latency and minimize processing time. However, the game client might also interact via representational state transfer (REST) with a game server or backend services.

Messaging

There are typically three primary categories of messages in games:

- Player engagement messaging targeted at a specific user or cohort of users, such as game invites or push notifications
- Group messaging between players, such as in-game chat
- Service-to-service messaging, such as JSON messages used to integrate two or more applications

A common strategy for sending and receiving these types messages is to use publisher-subscriber (pub-sub) and asynchronous processing architecture patterns. AWS provides several services that can help you to implement messaging in your game.

[Amazon Simple Notification Service \(Amazon SNS\)](#) provides a managed service for delivering messages between publishers and subscribers using a pub-sub architecture pattern. Publishers send messages via API to Amazon SNS, which delivers the messages asynchronously to subscribing applications, and can deliver push notifications directly to mobile clients or desktops with support

for some of the most widely used push notification services out of the box. Amazon SNS can be used for push notifications to clients as well as service-to-service messaging use cases.

[Amazon Simple Queue Service \(Amazon SQS\)](#) is a fully managed queue service that makes it easy to integrate game servers and your game regardless of the programming language used in each. Many games tasks can be decoupled and handled in the background such as updating a leaderboard or playtime values in a database. This approach is a very effective way to decouple various parts of your game and independently scale the player-facing features from backend processing.

[Amazon Managed Streaming for Apache Kafka \(Amazon MSK\)](#) is a fully managed service that makes it easy to build data streaming and producer/consumer applications using Apache Kafka, a popular open-source platform. Kafka is typically used for ingesting and processing real-time streaming data and can be used for service-to-service messaging

[Amazon ElastiCache \(Redis OSS\)](#) provides a fully managed in-memory data store which includes support for the popular pub/sub feature of Redis that is commonly used for developing chat room applications and high-performance service-to-service messaging. Redis also supports rich data types, such as lists and sets, that enable developers to use Redis for high-performance queuing.

[Amazon Pinpoint](#) provides user engagement messaging through email, SMS, voice, and push notifications. For example, Amazon Pinpoint can be used to deliver user engagement messages to players to invite them back to the game, and can be used for transactional use cases such as supporting multi-factor authentication tokens, and order confirmation and password reset emails.

Live game operations (Live Ops)

Live Operations (Live Ops) is a style of game management and operations that treats a game as a live service and continually delivers new features, updates, promotions, in-game events, and improvements to the launched game to improve the experience for the player community.

Traditionally, games were delivered as products rather than services, and new content and features were frequently incorporated into subsequent releases or sequels rather than into the launched product. With a Live Ops approach to game management, a game operations team can launch a game and maintain an engaged player community through experimentation, promotions, in-game events, and new innovations to keep players entertained. Although this approach has the benefit of unlocking new player engagement strategies and delivering recurring revenue streams, it requires more operational expertise. For example, to implement a successful Live Ops strategy, a developer

might need to integrate with cloud services or operate their own backend technical infrastructure, and require an effective way to identify and respond to issues that arise in the game, or within the player community, that can negatively impact the player experience.

General design principles

The AWS Well-Architected Framework identifies the following general design principles to facilitate good design in the cloud for games workloads:

Understand player behavior and usage patterns to evolve the game and protect players: To drive improvement continually to your game and effectively manage the player experience, it is important to have visibility into how players interact with game itself and with the rest of the players in the ecosystem. This helps you understand how to improve the game, manage costs, and monitor and react to unauthorized usage activity that poses risks to the player experience.

Use technologies that simplify game operations and increase development speed: Prioritize the adoption of technologies that can help you to improve speed and reduce the operational overhead of delivering new features and improvements to players. Games are hits-driven and there are many choices for players to consider, so moving quickly and adapting to change is critical for the success of a game. Consider whether you are comfortable operating your own software, or if you would prefer adopting a comparable managed service from AWS, AWS Partners, or both.

Optimize your architecture to improve metrics that reflect actual player experience: As you adapt and evolve your architecture over time, think about how these improvements and changes will impact player experience. Games workloads should be able to withstand and minimize the impact of failures to prevent widespread disruptions to gameplay. Game features and systems that aren't critically dependent on each other should be de-coupled to reduce the blast radius of failures and isolate player impacting issues.

Design infrastructure to meet peak player concurrency and dynamically scale as needed: Infrastructure should be designed to scale to meet player demand. Metrics, such as player session concurrency and number of logins, can be used to preemptively scale before systems become overloaded. Reactive system utilization metrics, such as CPU and memory consumption, can be used to scale after systems become overloaded. By dynamically scaling your infrastructure, you can reduce the costs of operating your game.

Implement runbooks to improve game operations: Operational runbooks should be used to help you consistently manage recurring game operations tasks. Runbooks should exist for common game operations workflows such as investigating and responding to player reports, managing infrastructure pre-scaling activities to prepare for anticipated large-scale events like new season launches and game content releases, and to address typical game maintenance activities.

Scenarios

In this section, we cover several scenarios that are common in a game architecture. Each scenario includes the common characteristics that drive the design and an example reference architecture diagram.

Scenarios

- [Game hosting for real-time synchronous gameplay](#)
- [Game backends](#)
- [Game production in the cloud \(GPIC\)](#)
- [Game analytics pipeline](#)

Game hosting for real-time synchronous gameplay

Real-time synchronous gameplay allows for two or more players to participate and interact in a game simultaneously where the state of the gameplay is shared between the connected players to create a real-time experience. Examples of synchronous games include first-person shooters, massively multiplayer online games (MMOG), sports and action games, or any online game where two or more players must be connected in order to share the play experience in real time.

Characteristics of real-time synchronous gameplay architectures include:

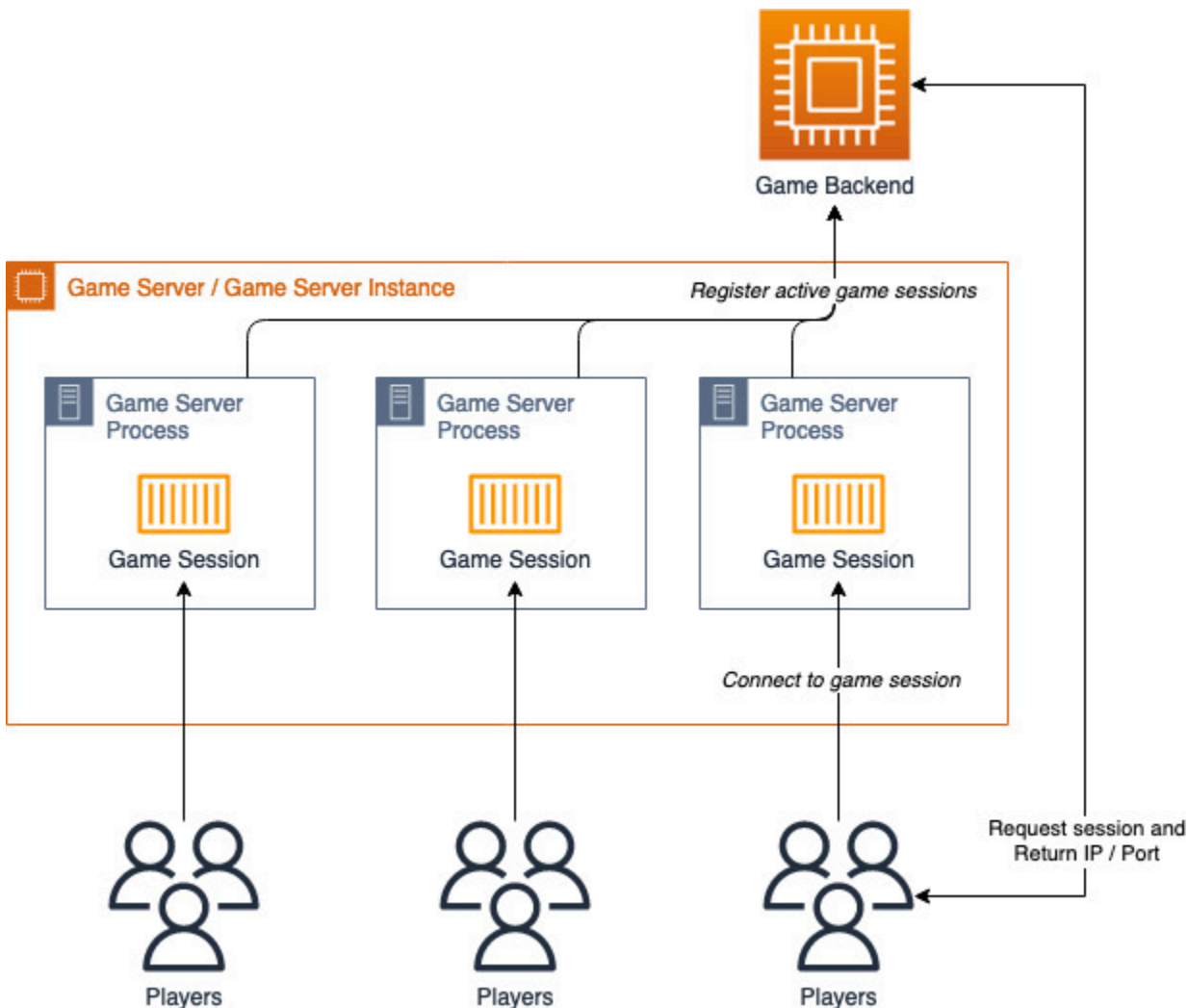
- Games are hosted as game sessions through game server processes that run on game servers. Game servers are hosted in multiple data centers and AWS Regions globally.
- Game clients can join a game session either by requesting a match from a centralized matchmaking service hosted in the game backend platform, or by selecting a match from a predefined list of available game servers. The game client is provided with an IP address and port to connect to.
- Many synchronous games are latency sensitive, for example, first-person shooters and massively multiplayer online games. These games typically have a pre-defined latency tolerance that is carefully measured and optimized to reduce the lag experience that can sometimes occur for players in high-latency situations. This latency information is determined by instrumenting game clients to ping the available game server locations to capture metrics such as latency, network jitter, and other important metrics for the gameplay experience. These metrics are sent to a central metrics collection service in the game backend platform so that live operations teams can

monitor game health. During the matchmaking process, game clients can provide their current latency data as one of the request parameters when requesting a match, and the matchmaking service uses can use that latency data as one of the variables when selecting a game server to host the player.

- Typically, the gameplay is conducted over a mix of protocols — game servers using faster UDP-based messaging, and matchmaking, authentication, and other client-server traffic using HTTPS.
- Game servers are frequent targets for malicious activities and should be protected with a DDoS protection solution like AWS Shield Advanced.

Game server processes

The following diagram illustrates a typical game server architecture. It describes the logical relationship between a game server instance and the game server processes that host game sessions.



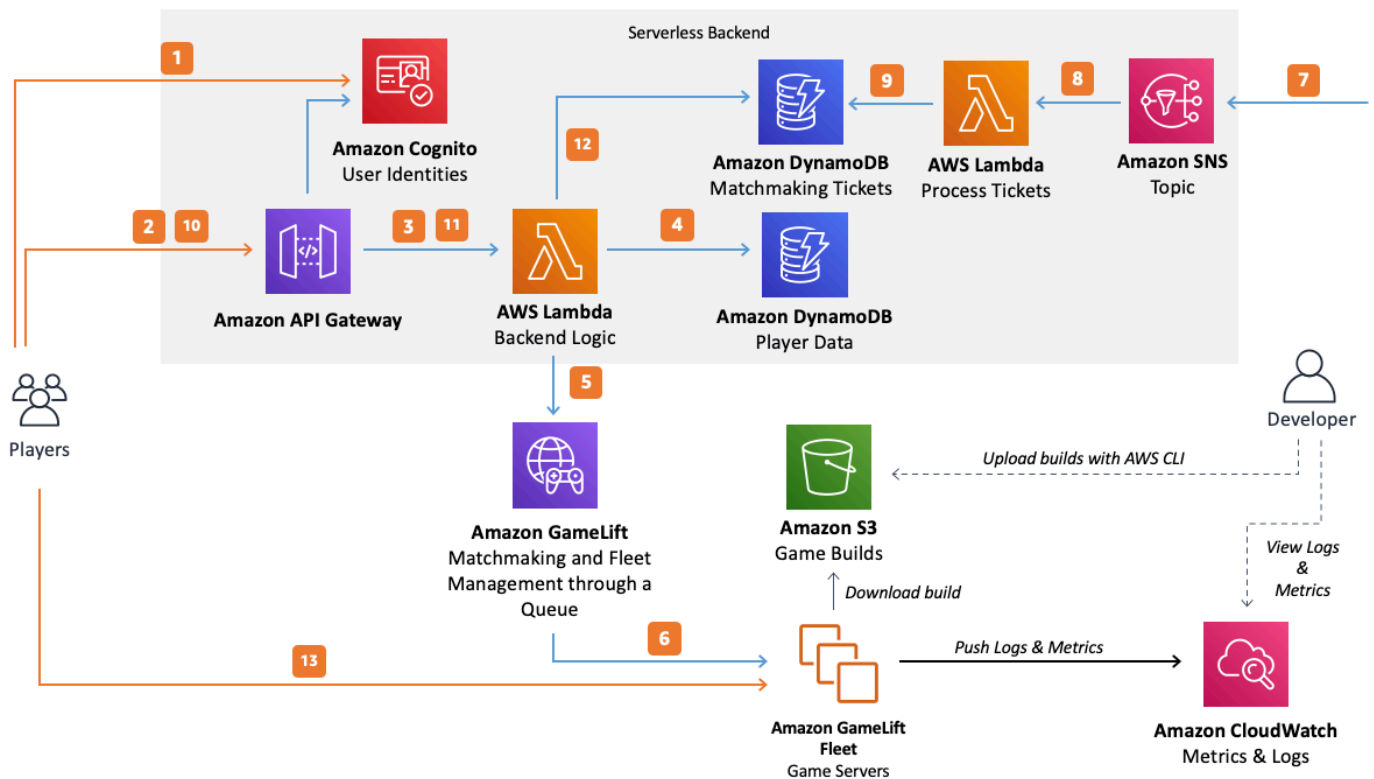
Logical game server architecture

- An Amazon EC2 instance is used as a game server, also known as a game server instance. The game server hosts one or more game server processes and each is running a copy of your game server build. Typically, multiple game server processes are running on a game server instance to utilize compute resources efficiently and reduce costs. When a game session is active and ready to host player sessions, its status is updated with the game backend (usually a matchmaking service) so that it can begin to be used to host players.
- The game backend can provide the player's game client with the IP address and server port that is hosting a game session so that they can connect to play.

Session-based game server hosting with serverless backend

When developing an architecture for your game, you should consider the features and capabilities you need, and the level of operational management overhead that you are prepared to own. To provide the best balance between ease of operations and flexibility, you can build your game using managed services from cloud providers that give you the control to develop and customize your own custom game features, while also reducing your burden to deploy and manage infrastructure.

Hosting a session-based multiplayer game requires having server infrastructure to host the game server processes as well as a scalable backend for matchmaking and session management. The following reference architecture shows how Amazon GameLift managed hosting and a serverless backend can be used to manage your session-based games.



Amazon GameLift managed hosting for session-based games

The diagram describes the process of getting players into games running on Amazon GameLift managed game hosting. It includes the following steps:

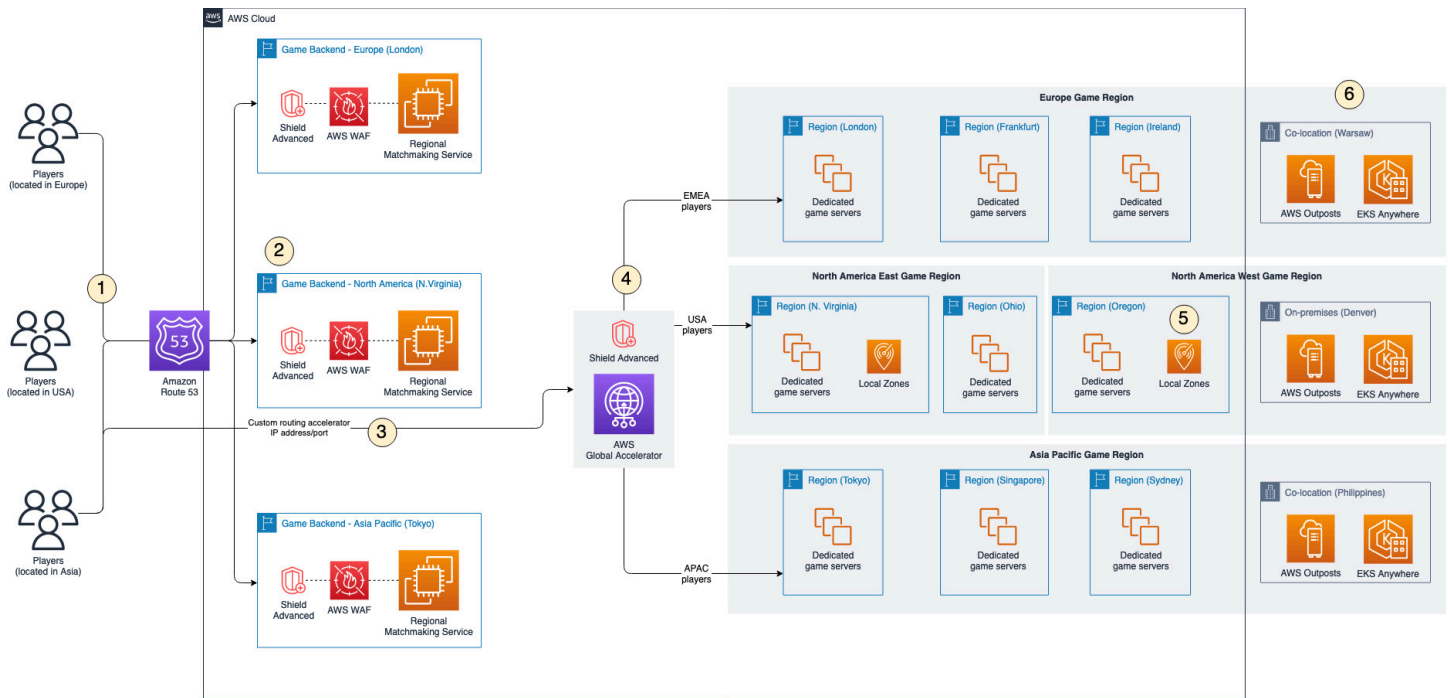
1. Game client requests an **Amazon Cognito** identity from Amazon Cognito. This can optionally be connected to external identity providers.
2. The game client receives temporary access credentials and requests a game session through an API hosted with **API Gateway** by using the credentials from Amazon Cognito.
3. API Gateway invokes an **AWS Lambda** function.
4. The Lambda function requests player data from a **DynamoDB** table. The Amazon Cognito identity is used to securely request the correct player data because the authenticated identity is provided in the request.
5. Using the correct player data for any additional information (like player skill level), the Lambda function requests a match through **Amazon GameLift FlexMatch** matchmaking. You can define a FlexMatch matchmaking configuration with JSON-based configuration documents. The game client can generate latency metrics by pinging server endpoints in various Regions, and the latency data can be used to support latency-based matchmaking.

6. After FlexMatch matches a group of players with suitable latency to a Region, it requests a game session placement through a Amazon GameLift queue. The queue contains fleets with one or more registered Region locations.
7. When the session is placed on one of the fleet's locations, an event notification is sent to an **Amazon SNS** topic.
8. A Lambda function will receive the Amazon SNS event and process it.
9. If the Amazon SNS message is a MatchmakingSucceeded event, the Lambda function writes the result to DynamoDB with the server port and IP address. A time-to-live (TTL) value is used to make sure that matchmaking tickets are deleted from DynamoDB when they no longer needed.
10. The game client makes a signed request to API Gateway to check the status of the matchmaking ticket on a specific interval.
11. API Gateway invokes a Lambda function that checks the matchmaking ticket status.
12. The Lambda function checks DynamoDB to determine whether the ticket has succeeded. If it has succeeded, the Lambda function sends the IP address, port, and the player session ID back to the client. If the ticket hasn't succeeded yet, the Lambda function sends a response declaring that the match is not ready.
13. The game client connects to the game server using the port and IP address provided by the backend. It sends the player session ID to the game server and the game server validates it using the Amazon GameLift Server SDK.

Alternatively, you can modify the preceding architecture to use API Gateway WebSockets with Amazon GameLift. In this approach, communication between the game client and your game backend service occur using a [websockets-based implementation](#). This implementation can be used so that the game backend Lambda function initiates a server-side message to the game client over a WebSocket rather than implementing a polling model.

Multi-Region and hybrid architecture for low-latency games

This section describes a multi-Region and hybrid architecture for low-latency games.



Reducing latency with network acceleration and game servers deployed globally

1. Players in a globally available game can originate from anywhere. When a player requests a game session or match, their game client sends a request to the game backend service registered with **Amazon Route 53**.
2. The game backend is deployed in multiple AWS Regions closest to player populations. Each game backend includes a Regional matchmaking service that finds a game session from across the game Regions. Although a player's matchmaking request is processed by a Regional matchmaking service near them, the matchmaking service is capable of routing a player to a game session in any game Region, if necessary. This action improves resiliency and performance. Additionally, each game backend service uses **AWS WAF** and **Shield Advanced** to provide layer-7 web filtering and bot control, and protections against distribution denial of service (DDoS). You have many options for how you build your game backend service, such as serverless, containers, EC2 instances, or hosting the game backend service in your own data centers.
3. To improve the experience for players by reducing network latency and jitter, a custom routing accelerator is deployed using **AWS Global Accelerator**, which automatically optimizes the routing of traffic from the game client to the game server. You configure the [custom routing accelerator](#) to map Global Accelerator listener ports to your game server's EC2 instance port. The game client connects to the Global Accelerator IP and port that deterministically routes the players to the correct game server IP and port that is hosting the game session.

4. Your game includes player-friendly logical game Regions that represent a collection of game server hosting locations that are geographically close to each other, such as *North America*, or *Asia Pacific*. You might also choose to create more granular regions, such as *North America East* and *North America West*. To reduce latency and increase geographic coverage, a combination of different game server hosting solutions can be used to improve the player experience. Prioritize the use of AWS Regions where possible, as these locations are fully featured and contain the largest capacity footprint.
5. Use **Local Zones** to host game servers in geographic locations of underserved players where you do not have existing hosting facilities or an AWS Region is not available.
6. Deploy **Outposts** into your existing on-premises data centers and co-location providers to create a seamless control plane and management experience across each deployment location using fully-managed racks and rack-mounted servers. Outposts are also beneficial if you don't have existing server capacity in your on-premises environment. However, if you want a hybrid implementation with software running on your own existing server infrastructure, you can use **EKS Anywhere**, which allows you to create and run clusters of containers on your own infrastructure with connectivity back to Amazon EKS, which provides a consistent console view of your Kubernetes clusters in AWS and on premises.

Game backends

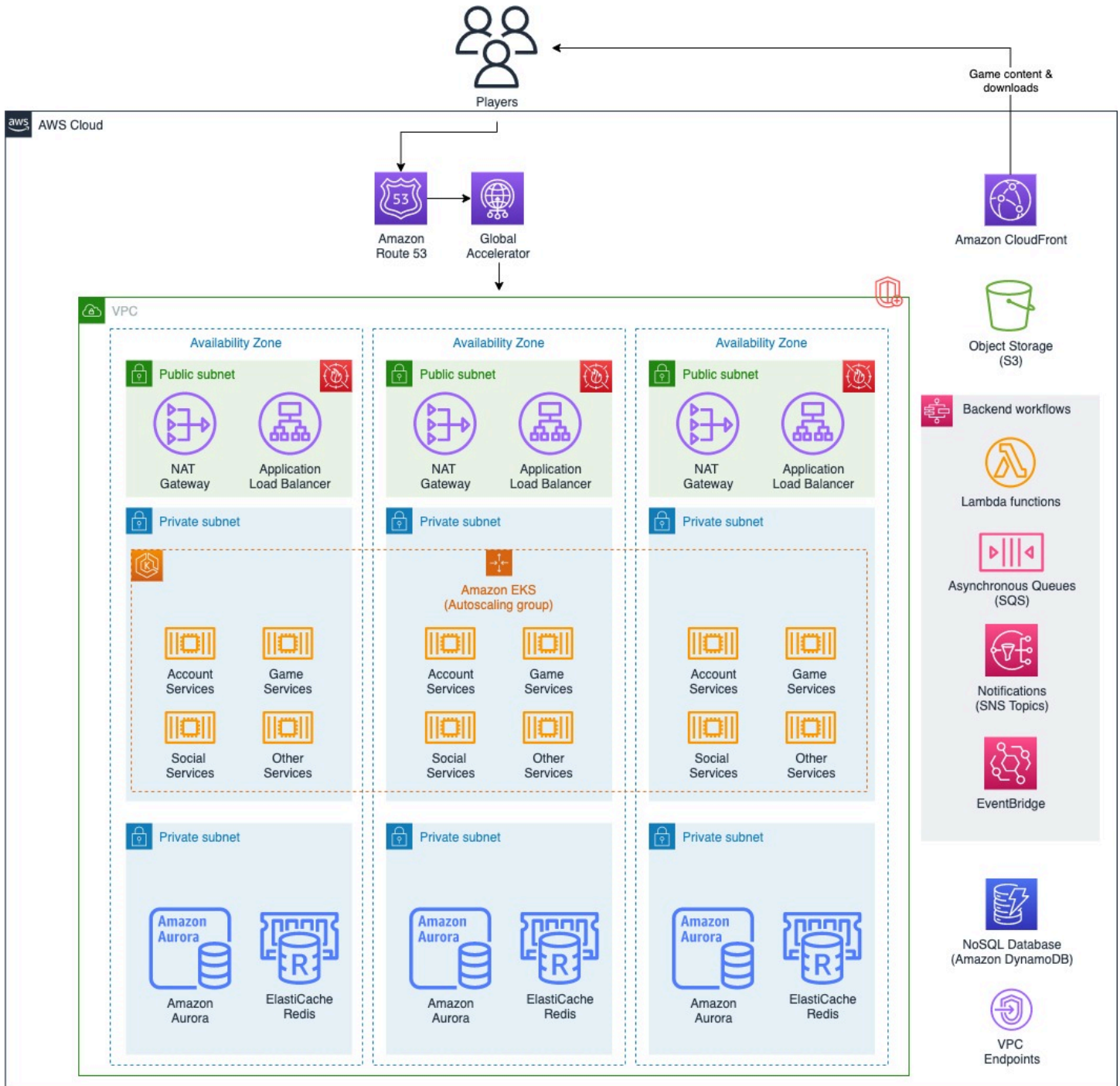
Game backends are used to manage game and player state, as well as integrate social- and platform-level features into the game that support the gaming experience. Player profile management, item and inventory storage, and stats and leaderboards are examples of services hosted in game backends. Game backends are typically built as REST APIs that are accessed by clients using HTTPS. However, other approaches are also common, such as WebSockets that provide bidirectional channels for use cases, such as client notifications for in-game chat and presence. Game backends can be deployed using a variety of different deployment architectures, including using instances, containers, or a serverless architecture.

Topics

- [Container-based game backend architecture](#)
- [Serverless-based game backend architecture](#)

Container-based game backend architecture

This section outlines a container-based game backend architecture



Hosting a game backend using containers

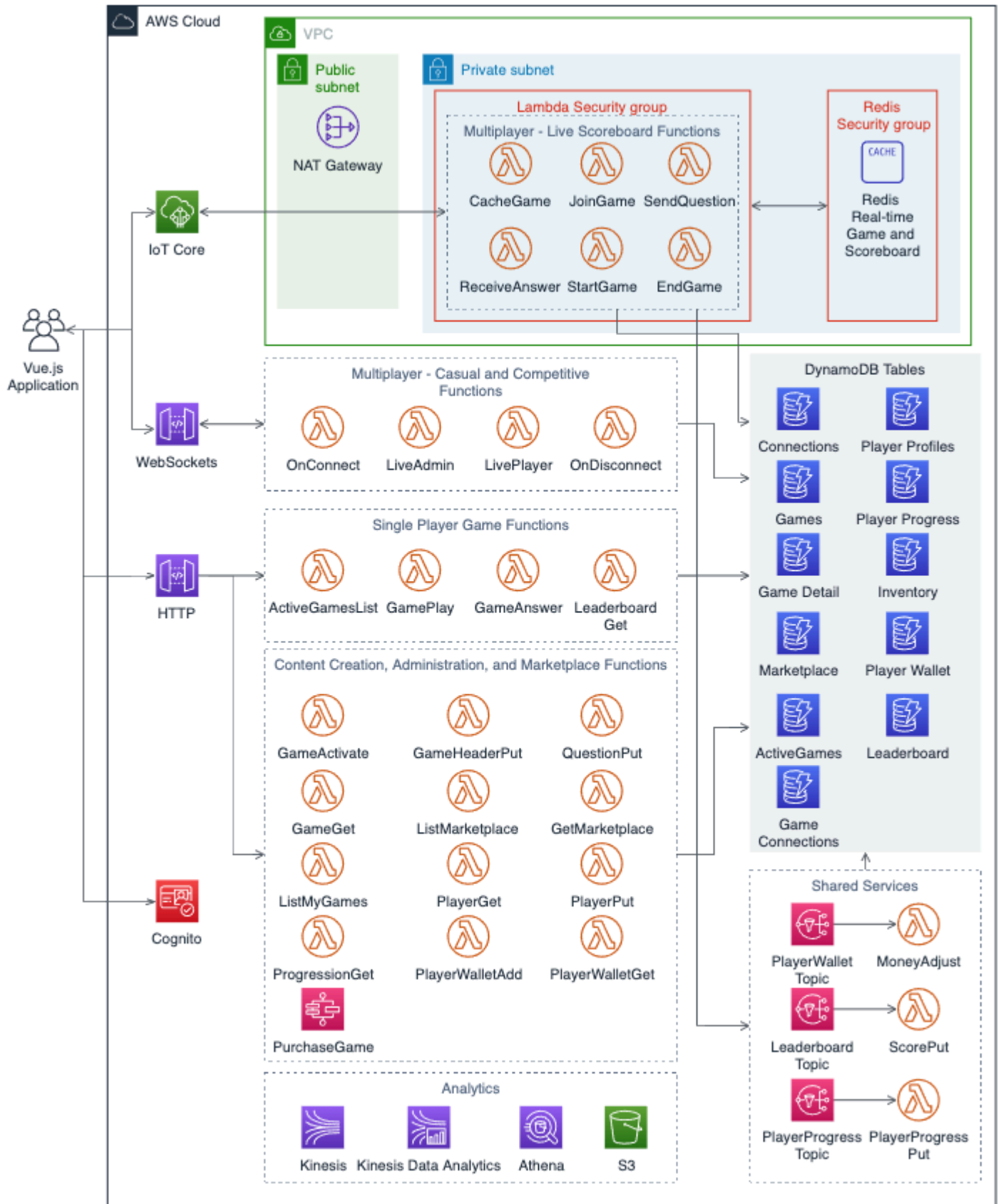
1. Players access the game using game client software, which can be distributed to them through a gaming platform, a digital storefront, or a direct download from a content delivery network

- (CDN), such as **Amazon CloudFront**. CDNs provide caching at edge locations to accelerate the performance for users downloading content. For example, CloudFront can be used to distribute the game client software to your players as well as the game assets and other content.
2. **Global Accelerator** provides traffic acceleration and customizable controls for routing traffic from player game clients to your load balancers as well as routing traffic across Regions for multi-Region and failover purposes. The custom domain names for your game backend REST APIs are configured in Route 53 to route traffic to Global Accelerator endpoints. **Shield Advanced** provides DDoS mitigation for your accelerator and game backend.
 3. NAT Gateway and Application Load Balancer are deployed to public subnets in each of the Availability Zones used by the game backend to help ensure high availability within the Region. **Web Application Firewall** is deployed on the Application Load Balancer to provide layer-7 web traffic filtering.
 4. The game backend is hosted as a collection of individual container-based microservices deployed into an **Amazon EKS** cluster in private subnets spread across Availability Zones for resiliency. Automatic scaling dynamically adjusts the capacity of the services and cluster nodes based on resource utilization, which is typically correlated with player demand. Whereas **Cluster Autoscaler** automatically adjusts the number of nodes in the cluster, the **Horizontal Pod Autoscaler** automatically scales the pods deployed into the cluster.
 5. Game and player data is stored in backend databases and caches that are deployed into private subnets across Availability Zones with replication between primary and replica nodes. **Aurora** is a popular choice for use cases such as player profiles, entitlements, and in-game purchasing, which can have more complex querying requirements and may benefit from the relational data modeling capabilities of MySQL and PostgreSQL. **ElastiCache for Redis** is useful for building high-performance leaderboards, pub/sub messaging, and for caching frequently accessed data to reduce latency and load on databases. **DynamoDB** is a fully managed NoSQL data store that is ideal for unpredictable access patterns and the ability to scale to virtually unlimited throughput for use cases such as player and game state data, session data, inventory and item stores, or use cases where you want a global database with minimal overhead.
 6. Asynchronous processing workflows should be used to perform work that can be done in the background, such as updating leaderboards or sending friend requests. Configure your game backend to push this type of work into **Amazon SQS** queues to help you scale as your game grows, or consider using **Amazon SNS** topics to distribute the work between many consumer application queues for parallel processing. Use **Lambda functions** to perform the processing in an event-driven manner to reduce your compute infrastructure costs and management overhead. For workflows that are long-lived or require task coordination with multiple steps,

consider orchestrating the entire workflow using **Step Functions**. **Amazon EventBridge** can be used to initiate functions to respond to AWS service and custom application events.

Serverless-based game backend architecture

Many game developers do not want to manage infrastructure, and instead prefer to build their games using technologies that allow them to focus on software. A serverless architecture is recommended in this scenario because it allows you to build and release features more quickly, and with less operational overhead. Serverless architectures are designed using cloud services that can scale dynamically based on demand without needing to setup, manage, and scale servers. The following reference architecture illustrates how to build a game on using a serverless architecture.



Serverless-based game backend reference architecture

This reference architecture illustrates a web-based trivia game that provides single player and multiplayer features.

- **Player authentication:** Players authenticate using Amazon Cognito, which provides secure authentication with a user directory for player identity management.
- **Game logic as serverless functions:** All of the game's features and backend business logic run as Lambda functions that are initiated in response to events, which keeps costs down because you only pay when the function runs. Lambda gives you the flexibility to write each game feature as a separate microservice using a programming language of your choice. For example, you might choose to develop .NET Lambda functions if you have experience using C# to build Unity games, or you might choose to develop Node.js Lambda functions if you desire to program a frontend and backend for a web-based game both in JavaScript.
- **NoSQL Data Store for game and player data:** Use DynamoDB to store your player and game data, as it is purpose-built for storing large amounts of data from microservices. As illustrated in this architecture, it is a best practice to use separate data stores for each game feature's data storage needs, which makes it easier for you to monitor and manage the features independently. This also helps create boundaries of separation if feature or service ownership changes within your team. In this reference architecture, DynamoDB tables are used to store data such as connection state, game details, player progress, and leaderboard information.
- **Single player gameplay:** Single player features allow players to perform actions like selecting and playing a game and viewing the leaderboard. These features are implemented as RESTful backend services hosted with an Amazon API Gateway HTTP API that invokes the appropriate Lambda function to get and set data in DynamoDB tables. When gameplay completes, the backend also sends notifications to SNS topics that initiate Lambda functions asynchronously to store the player's progress and stats.
- **Multiplayer gameplay:** Multiplayer game features require the players to be able to interact with the game for point-to-point communication, as well as broadcast and receive updates from other connected players. A WebSockets implementation is suitable for point-to-point communication in a lightweight game, such as trivia. Players can establish a WebSockets connection to Amazon API Gateway WebSockets, which manages the connection and only invokes the Lambda functions when there are messages to send or receive for a player. For use cases where one-to-many communication is required between players, AWS IoT Core provides support for messaging using WebSockets over MQTT, which allows clients to subscribe to topics and act on the messages it receives. In this architecture, WebSockets over MQTT is used to support use cases such

as broadcasting live in-game updates and asking questions to all connected players. As an alternative to AWS IoT Core, you might choose Redis Pub/Sub for message delivery, or Redis Streams if you need message retention.

- **Use VPC-enabled Lambda functions to access resources in your private subnets:** Configure VPC-enabled Lambda functions to access resources in your VPC's private subnets, such as ElastiCache for Redis, which is used to reduce query times for low-latency datasets like live leaderboards.

For more information, refer to the [Simple trivia service code sample](#).

Game production in the cloud (GPIC)

Game production in the cloud (GPIC) refers to the infrastructure and tools that are required for the game development lifecycle to build, test, and develop a game. Game development is collaborative between users and the infrastructure requirements frequently change throughout the stages of development. Many game developers are embracing globally distributed and remote development teams, which requires technology that supports this type of development. Game developers can host all or part of these environments in AWS and use the global availability of AWS Regions to place resources closer to users, and manage their development environments more cost effectively by scaling compute and storage as needed.

The environments can vary depending on game developer needs, but they typically include developer workstations for artists, designers, engineers, QA testers, contractors, and other personnel to perform their work. These environments also typically include a build farm comprised of source code repositories for users to check-in their changes and the CI/CD infrastructure for building, packaging, and testing the developed artifacts.

These game production architectures have the following characteristics:

- Users should be able to access a virtual workstation through a web browser or local desktop client, such as [NICE DCV](#), that provides them with a low latency streaming session to access the same software and tools that they would have access to if they were working on a machine in an office or development studio. These virtual workstations, typically a cloud-based server, should allow a user to collaborate and work on their projects entirely in a cloud environment over a LAN or WAN. When users are not actively using the machines, their work should be backed up to durable cloud storage, for example a source control repository or file system such as [Amazon](#)

[Elastic File System \(EFS\)](#) and [Amazon FSx](#), and their machine should be shut down to reduce costs.

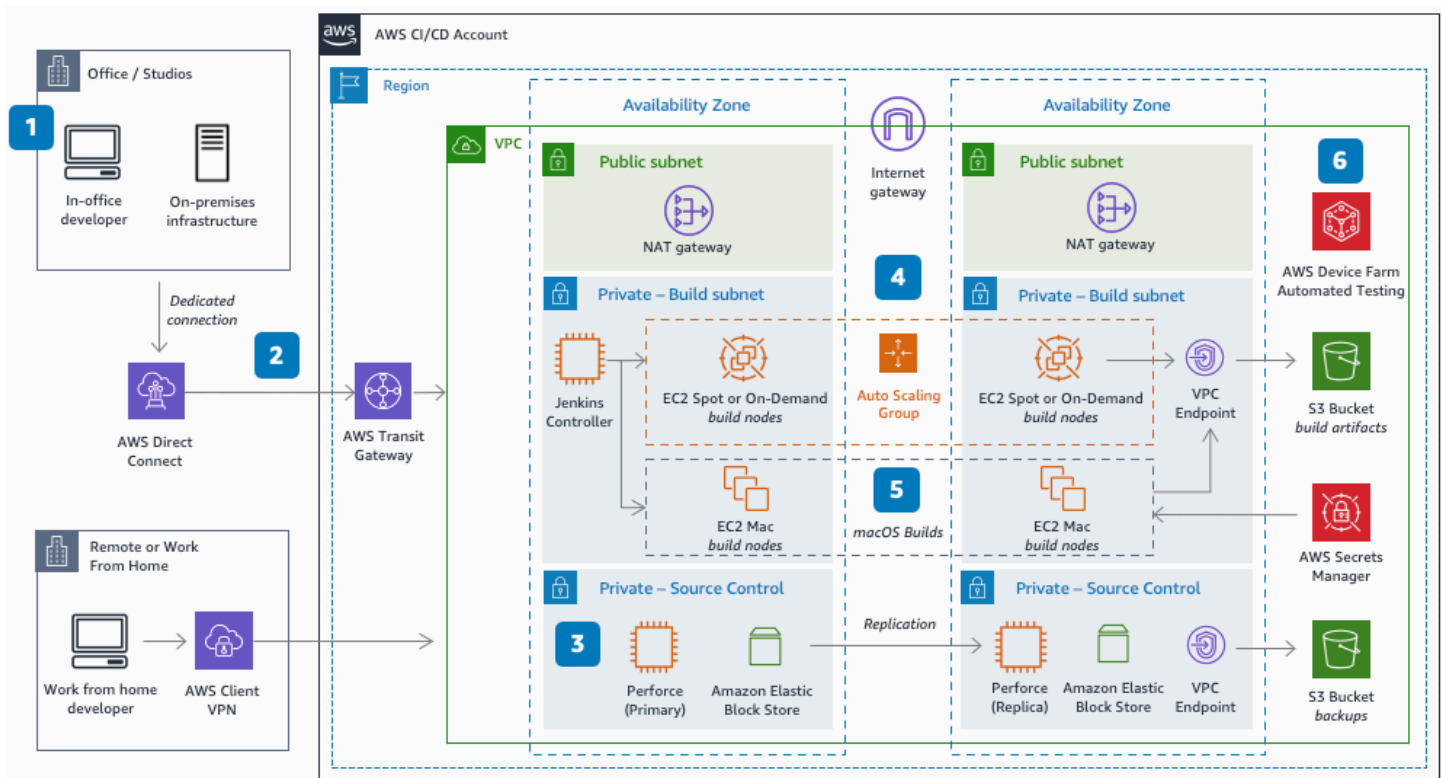
- Source control repositories, such as Perforce, should be designed with high availability using replication between Availability Zones or between on-premises environments with backups stored in cloud storage such as [Amazon S3](#). For example, a cloud-based Perforce server should include a primary commit server hosted in one Availability Zone with replication to a standby server hosted in another Availability Zone in the same Region.
- Game development build farm resources should be designed with automatic scaling so that compute resources are provisioned as they are required, and [EC2 Spot Instances](#) should be used to reduce the costs incurred when scaling out the number of servers required for builds.

Topics

- [Game production in the cloud – CI/CD](#)
- [Game production in the cloud – Workstations](#)

Game production in the cloud – CI/CD

You must have CI/CD infrastructure when developing games. A game development CI/CD pipeline is typically comprised of highly available source control servers and storage, compute resources to run your builds, and software to perform automated testing, along with the proper network connectivity from your development machines. The following reference architecture demonstrates how to offload game builds from remote or on-premises game development environments to the AWS Cloud to aid developers in migrating or building new build farms.



Offload game builds to the cloud

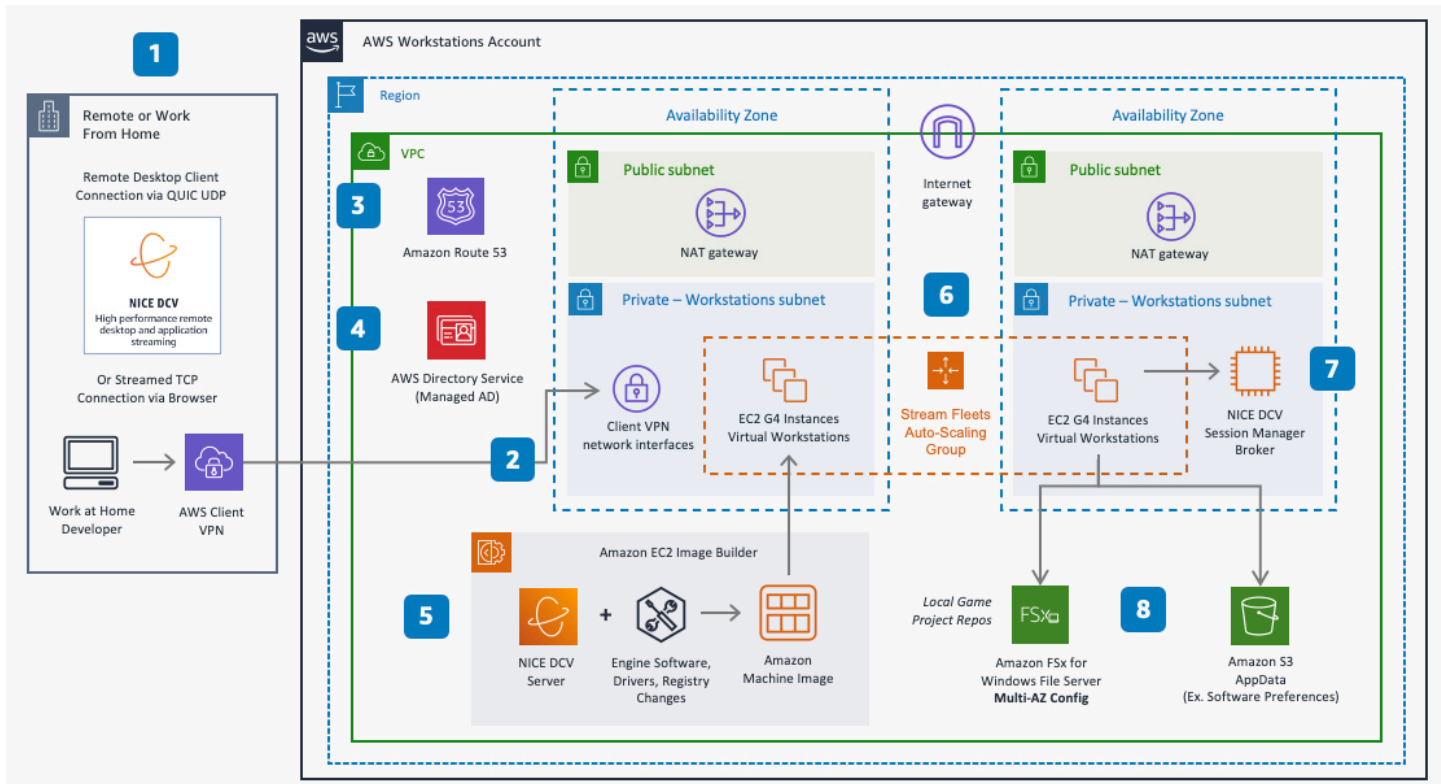
1. **AWS Direct Connect** provides a low latency, private dedicated connection to AWS for in-office developers. Remote developers use **AWS Client VPN**.
2. **AWS Transit Gateway** simplifies network management for connectivity between VPCs and from on-premises networks.
3. Perforce manages source and version control (CI) backed by **Amazon EBS** storage for quickly accessed, persistent data. Perforce Helix Core (P4D) is available in the **AWS Marketplace**.
4. Commits start a build (CD) in Jenkins when developers push changes to Perforce tied to a branch. Perforce triggers POST a JSON payload to Jenkins. The Jenkins controller calls engine “headless” CLI commands to run and parallelize the build process across ephemeral, Docker nodes (such as **Amazon EC2 Spot Instances**), or **Amazon EC2 On-Demand Instances**. Developers can increase availability by using two Jenkins controllers, one in each Availability Zone, behind a load balancer. For some engines, developers might need additional licensing infrastructure configured in additional subnets to vend licenses for the build context each time a concurrent build is run.
5. The Xcode portion of iOS builds is offloaded to **Amazon EC2 Mac instances** to sign, build, and export the .IPA file, splitting the process and reducing build times. **AWS Secrets Manager** holds provisioning profiles, private keys, and certificates.

6. Build artifacts are delivered to **Amazon S3**, which sends notifications of success or failure. **AWS Device Farm** enables automated testing.

Game production in the cloud – Workstations

The game development process can be highly distributed with teams working on various aspects of the game from many locations, and it is important that they have access to workstations and shared storage in order to be productive. In a simple scenario, a game can be developed by a single game development studio located in a specific office location. However, in practice, it is more common for developers to contribute to a game project from many geographic locations, such as a single game studio with many globally distributed teams in different cities or countries, or work-for-hire studios, contractors, and co-development studio partners who may support a project by bringing specific expertise such as art development or game testing, QA, and localization specialization to the project. Game development also must be able to support remote work so that developers can be productive from home or other locations. An increasingly popular trend in the games industry is for newly developed studios to be fully remote without the expectation of working in an office environment.

The following reference architecture demonstrates how to use AWS for hosting remote game development workstations using the NICE DCV protocol.



Stream game development from anywhere with NICE DCV

1. **NICE DCV** is a streaming protocol that supports 4K, 60-FPS streaming. Developers using a browser connect via TCP connections whereas desktop clients can use QUIC UDP over port 8443 for increased performance.
2. Developers use **AWS Client VPN** for a secure connection to network interfaces in the workstation subnets with source network address translation (SNAT).
3. **Amazon Route 53** provides private DNS for the resources in the VPC as well as inbound and outbound DNS forwarding.
4. **AWS Directory Service** provides managed Microsoft Active Directory to enable local game project storage mapped to individual users.
5. Workstations are created using an **Amazon Machine Image (AMI)** built with **Image Builder**. Images include NICE DCV Server, developer software, registry changes, and drivers such as GPU gaming drivers or peripheral drivers. **AWS Marketplace** includes common AMIs used for workstations.
6. Fleets of workstations use graphics **Amazon EC2** instance types that provide GPUs and are scaled using EC2 Auto Scaling groups.
7. A Session Manager Broker enables management of NICE DCV sessions.
8. Local file storage of projects is hosted in **Amazon FSx for Windows File Server**. Developers commit to a separate CI/CD pipeline by pushing from local storage to source control.

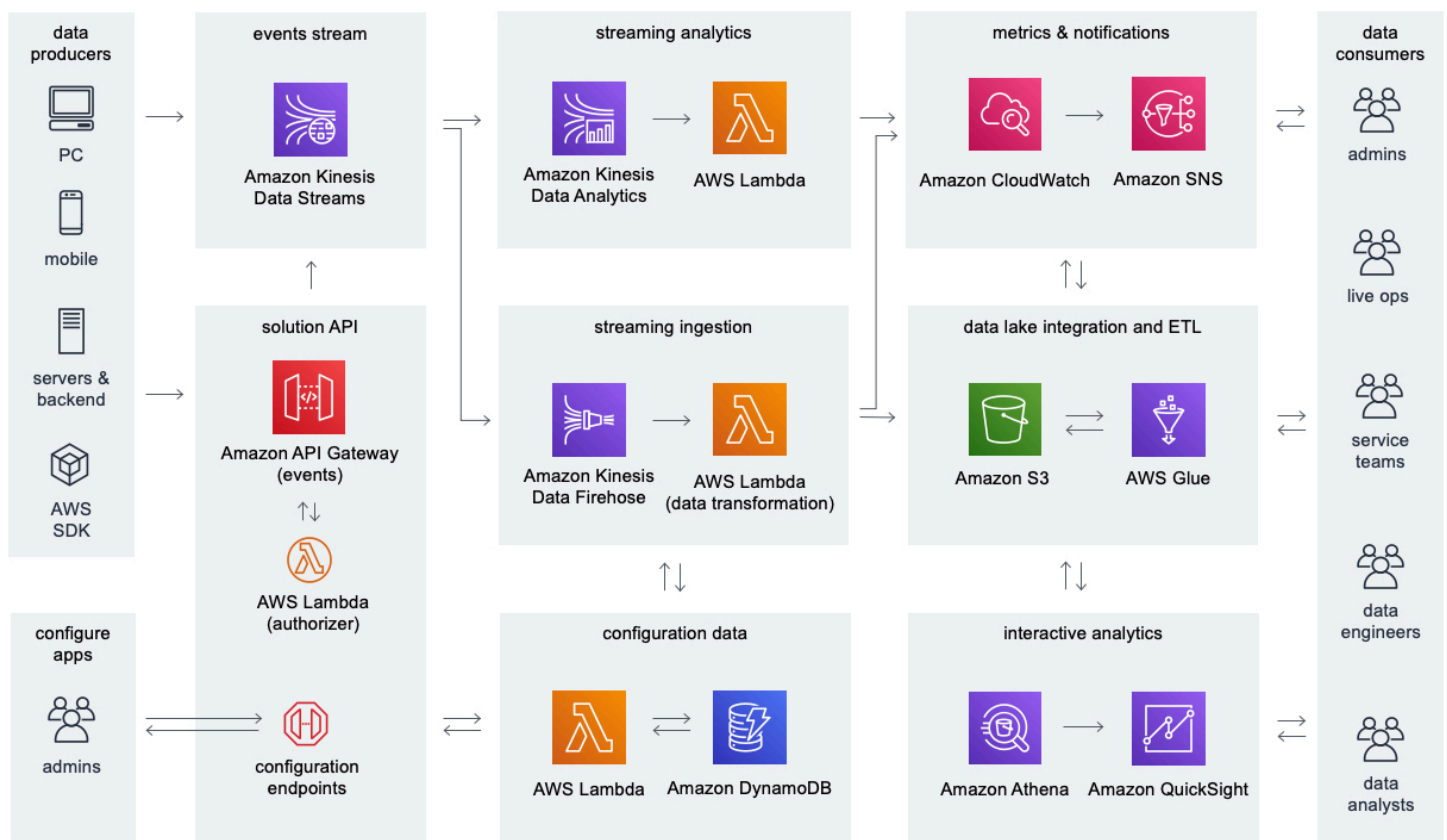
Game analytics pipeline

Game developers are increasingly looking for ways to better understand player behavior so that they can improve the gameplay experience to retain and grow their player base. Game analytics represents the technical infrastructure and processes that is required to understand and analyze all of the data that is generated from the game and related services. This typically requires the use of an analytics pipeline architecture that can support this end-to-end process, such as the [Game Analytics Pipeline](#) solution implementation.

Game Analytics architectures have the following characteristics:

- Data sources send data in a common format such as JSON and typically include game servers and game backend services, as well as game clients including PC, mobile devices, and game consoles.

- A game analytics pipeline automates the entire workflow of ingesting and storing the raw data, and processing it into usable output formats so that it can be analyzed efficiently and cost effectively by data consumers, such as end users and analytics applications.
- Game analytics pipelines provide support for ingesting and processing high volumes of real-time data in order to scale as a game grows.
- Provide support for both real-time and batch reporting use cases. For example, real-time dashboards and alerts are typically used by Live Ops teams to monitor game infrastructure and player behavior to detect issues. Data analyst teams typically rely on ad-hoc and batch reporting to understand trends over time.



Serverless game analytics pipeline for gameplay telemetry

Game data is ingested from game clients, game servers, and other applications. The streaming data is ingested into Amazon S3 for data lake integration and interactive analytics. Streaming analytics processes real-time events and generates metrics. Data consumers analyze metrics data in Amazon CloudWatch and raw events in Amazon S3.

1. **Solution API and configuration data:** Use Amazon API Gateway to provide a REST API for administering your game analytics pipeline and storing configuration data in Amazon

DynamoDB using Lambda functions. You can build an internal portal on top of this API or a custom command line interface for administration. REST API also provides server-authentication for ingesting gameplay data from data sources and forwarding the telemetry data to Amazon Kinesis Data Streams for real-time processing and ingestion into storage.

2. **Events stream:** Kinesis Data Streams captures streaming data from your game for data processing and storage.
3. **Streaming analytics:** Managed Service for Apache Flink analyzes streaming event data from the Kinesis Data Streams and can generate custom metrics and alerts that are published to CloudWatch using Lambda functions.
4. **Metrics and notifications:** Use Amazon CloudWatch to monitor your solution's metrics, logs, and alarms. Use Amazon SNS for sending notifications to on-call engineers and other data consumers.
5. **Streaming ingestion:** Use Kinesis Data Firehose to easily consume your streaming data from Kinesis Data Streams and deliver it to your data lake in Amazon S3 for long-term storage, transformation, and integration with other data.
6. **Data lake integration and ETL:** Use Glue for ETL processing workflows and to organize your metadata in the Glue Data Catalog, which provides the basis for a data lake for integration with flexible analytics tools.
7. **Interactive analytics:** End users can use Amazon Athena to perform ad-hoc interactive queries on the datasets stored in Amazon S3, and QuickSight can be used to build dashboards.

Refer to the [Game Analytics Pipeline](#) for an automated reference implementation of an analytics pipeline that can be deployed into your account using CloudFormation.

Well-Architected pillars

This section describes the Games Industry Lens in the context of the five pillars of the Well-Architected Framework. Each pillar discusses design principles, definitions, best practices, evaluation questions, considerations, key services, and useful links.

Pillars

- [Operational excellence](#)
- [Security](#)
- [Reliability](#)
- [Performance efficiency](#)
- [Cost optimization](#)

Operational excellence

The operational excellence pillar focuses on best practices for deploying and operating cloud-based games at any scale. It is important to focus on operational excellence to maintain a positive player experience and take preventative measures to prepare for and recover from issues that impact their experience.

Design principles

In addition to the design principles from the Well-Architected Framework whitepaper, the following design principles can help you to achieve operational excellence.

Define measurable and achievable objectives for game operations teams, adapt as necessary:

Due to the hits-driven nature of games, it is difficult to determine ahead of time how many players will show up to play your game when it launches, or what expectations players will have for your ongoing game operations. It is important to set ambitious but achievable goals with stakeholders and design an approach that can be scaled up if your game exceeds expectations and scaled down while game development teams optimize the player experience. Adequately prepare and test ahead of time to meet these requirements, and make sure that business and technical stakeholders are aligned on target objectives for operating the game. With targets defined, the game teams are able to achieve an appropriate balance between cost and performance during planning, designing, provisioning, deploying, and operating the game backend infrastructure.

Use operational runbooks to proactively plan and pre-scale for game launches and special events: Game operations teams should coordinate with business stakeholders to model projections for anticipated peak player concurrency for events and perform proactive planning to pre-scale infrastructure capacity ahead of time. Due to the fluctuating nature of player traffic during events, prior planning and pre-scaling activities should augment your existing automated scaling systems to improve your chance of success during an event and ensure that you have enough resources to provide a positive player experience. Develop operational runbooks to ensure consistency in the process.

Establish an operating model for receiving, investigating, and responding to player support requests: Post-launch, it is important to monitor reports of complaints and issues with the game. Implement appropriate systems to interact with players in a secure and effective manner to adequately resolve player issues and such as community forums, social media, e-mail, ticketing systems, call centers or automated chat bot solutions.

Definition

Extending the operational excellence pillar in the Well-Architected Framework for cloud-based games. Through this lens, we will review how each of the following factors contribute towards operational excellence for cloud-based games.

There are three best practice areas for operational excellence of games workloads in the cloud:

- **Prepare**
- **Operate**
- **Evolve**

Best practices

The following topics are best practices for your cloud architecture.

Topics

- [Prepare](#)
- [Operate](#)
- [Evolve](#)

Prepare

GAMEOPS01 - How do you define your game's live operations (LiveOps) strategy?

GAMEOPS_BP01: Use game objectives and business performance metrics to develop your live operations strategy.

You should consult business stakeholders, such as game producers and publishing partners, to determine objectives and performance metrics for a game. This can help you to develop plans for how you will manage the game, including defining your maintenance windows, software and infrastructure update schedules, and system reliability and recoverability goals.

For example, you might define targets for player concurrency (CCU), daily and monthly active user targets (DAU/MAU), infrastructure budget, financial targets, or other performance metrics such as the frequency for release of new content and features, or the frequency of in-game events and promotions to increase player engagement. These objectives and metrics feed into decisions about the game design, release management, observability, and support that is needed for efficient operations.

Your game might have an objective to release new content updates at least once each month with no downtime during release. This information can help you to define your release deployment strategy and coordinate the scheduling of required maintenance that may require downtime at other times throughout the month and contribute towards your availability SLA.

These metrics can also help you to determine at which stage of your game's lifecycle you should incorporate a live operations team (Live Ops) to monitor game health, collect direct game feedback, and build streamlined and automated release processes. For example, a new game might wait until a certain scale is achieved, measured by active player count, revenue, or another set of metrics, before setting up a dedicated live operations team. An established game development studio might already have live operations experience, perhaps for their other games, so they'd only need to on-board the new game.

GAMEOPS_BP02: Validate and test your existing game software before reusing it in your game

Organizations tend to reuse existing components and source code (from a previous game) to save on development time and cost. These legacy components and/or code may not be subjected to a thorough review, or have detailed integration testing, and instead rely on their past performance.

While reuse helps improve productivity, it can also introduce the risk of reintroducing past performance and stability issues into a new project. Therefore, when reusing existing components and source code from previous games, robust testing should be implemented.

For example, if the source code and components that were designed, written, and tested for Game A are reused in Game B, they might not be able to handle all the conditions that Game B requires. During a production incident, the developers might not have sufficient knowledge to debug and fix that code/component, or the time to rewrite it to alleviate their operational pain. If the original authors of the code are unavailable, that can increase the time to implement appropriate fixes. It is recommended that if previously used code or components had an issue, replacing or fixing them must be a priority before they are used again, without waiting for them to impact operations again.

GAMEOPS02 - How do you structure your accounts for hosting your game environments?

GAMEOPS_BP03: Adopt a multi-account strategy to isolate different games and applications into their own accounts.

A game architecture deployed in AWS should use multiple accounts that are logically organized to provide proper isolation, which reduces the blast radius of issues and simplifies operations as your game infrastructure scales. AWS accounts that host game infrastructure are typically grouped into the following logical environments:

- **Game development environments (Dev)** are used by developers for developing the software and systems for the game.
- **Test or QA environments** are used for performing integration testing, manual quality assurance (QA), and any other automated testing that must be conducted.
- **Staging or pre-production environments** are used for hosting final built software so that load testing and smoke testing can be conducted prior to launching to production.
- **Live or production environments** are used for hosting the live software and infrastructure and serving production traffic from players.
- **Shared services or tools environments** provide access to common platforms, software, and tools that are used by many different teams. For example, a central self-hosted source control repository and game build farm might be hosted in a shared services account.
- **Security environments** are used for consolidating centralized logs and security technologies that are used by teams that focus on cloud security.

For game infrastructure on AWS, it is recommended to create separate accounts for each game environment (dev, test, staging, prod), as well as accounts for security, logging, and central shared services.

Typically, smaller game development studios that manage a limited number of infrastructure resources, usually a few hundred servers or less, may desire to create one AWS account for each of these environments, for example, one production account, one development account, and one staging account. However, as your game infrastructure or team size grows over time, this simplified model may not scale well. When setting up these environments, it is important to consider that many AWS services share resource and API-level [Service Quotas](#) for an entire account within a particular Region. This must be considered when determining how to logically organize accounts. AWS accounts only incur cost for consuming services deployed into them. Therefore, this provides a way to effectively reduce resource contention and service quotas, particularly as your game grows and more developers need access to build and manage resources.

Based on our experience working with larger game development studios that typically operate thousands of servers with hundreds of developers accessing resources, we recommend to design a more fine-grained account structure where individual applications supporting your game have their own development, test, staging, and production accounts for each. Our experience supporting large and successful games shows that it is difficult and time consuming to re-design your AWS multi-account strategy after you have launched your game due to the complexity in planning and migrating live systems. Consider your future scaling needs when determining the right multi-account structure.

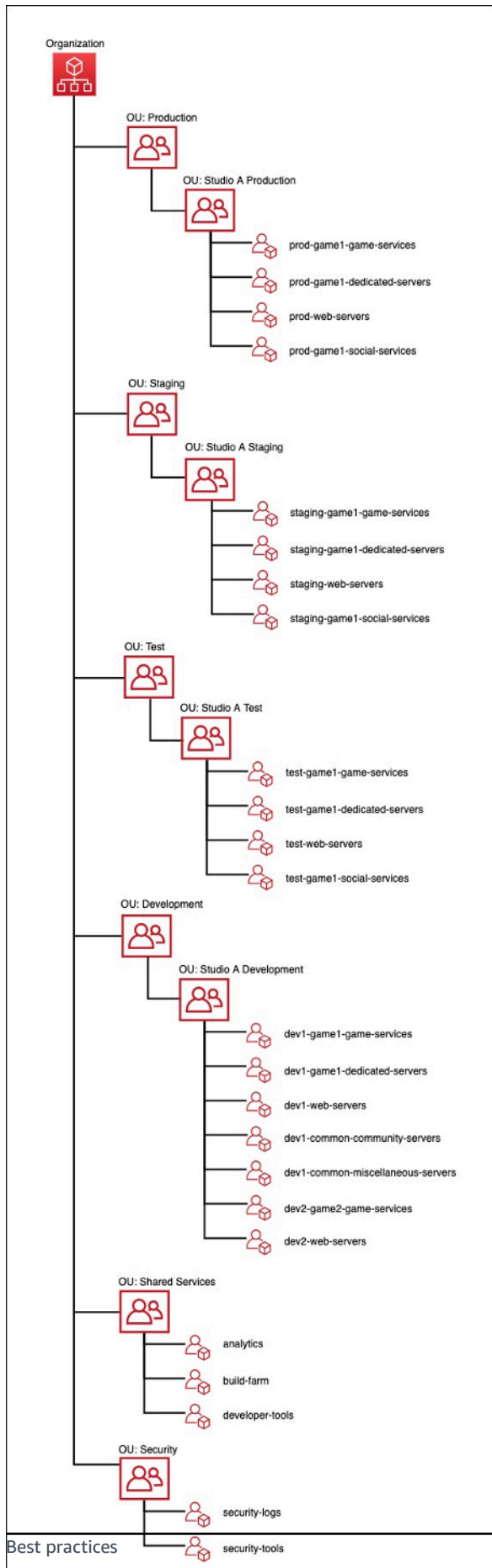
[AWS Organizations](#) can be used to set up a hierarchy and grouping of AWS accounts, and define [organizational units](#) (OUs) to apply common OU-level policies to them through [service control policies](#) (SCPs). AWS Organizations helps you centrally manage and govern your environment as you grow and scale your resources. You can programmatically create new accounts and allocate resources, group accounts to organize your workflows, apply policies to accounts or groups for governance, and simplify billing by using a single payment method for all of your accounts. Additionally, Organizations is integrated with other services so that you can define central configurations, security mechanisms, audit requirements, and resource sharing across accounts in your organization.

[AWS Control Tower](#) provides the easiest way to set up and govern a secure, multi-account environment, called a landing zone. Control Tower creates your landing zone using AWS Organizations, bringing ongoing account management and governance as well as implementation best practices based on AWS's experience working with thousands of customers as they move to

the cloud. [AWS Config](#), [AWS Trusted Advisor](#), and [Security Hub](#) are services that provide with an aggregated or centralized view of your account's hygiene.

Such isolation allows you to set up custom or individual permissions and guardrails to each game environment. Production accounts should have all the necessary guardrails, access restrictions, monitoring and alerting, and security tools, while non-production accounts may not require the same level of guardrails and permissions. Non-production environments can be automated to shutdown resources after hours and save costs. Separation of accounts at this level of granularity makes it easier to monitor infrastructure costs for each of the environments supporting a game.

The following is an example of a multi-account structure for a game company using AWS Organizations, and organizational units (OUs) to logically group AWS accounts into separate environments and studios. In this example, OUs are used to group together accounts based on their environment, and then based on the studio that operates the environment. This demonstrates how you can create a nested hierarchy to allow separate applications and games to be deployed into their own accounts, which can be useful if you develop and operate multiple games. Refer to the documentation and whitepapers provided in the resources section of this pillar to learn about additional strategies that you can consider for organizing your multi-account strategy.



Example of account structure for game environments

GAMEOPS_BP04: Organize infrastructure resources using resource tagging.

To effectively manage and track your [infrastructure resources](#) in AWS, it is recommended to use proper [resource tagging](#) and [grouping](#) that can help identify each resource's owner, project, app, cost-center, and other data. Tagged resources can be grouped together, using [resource groups](#), which assists with operational support.

As a best practice, you should define [tagging policies](#). Typical strategies include resource tags for identifying the resource owner, such as team name or individual name, the name of the game/app/project name, the studio name, environment (such as dev, test, prod, staging, or common), and the role of the resource (such as, database server, web server, dedicated game server, app server, or cache server). You can add any other tags to help with business and IT needs. [AWS Config](#) can also help to enforce a [tagging policy](#) at resource creation and update time. Tags and resource groups are available from the AWS Management Console, the AWS CLI, and API operations.

GAMEOPS03 - How do you manage game deployments?

GAMEOPS_BP05: Adopt a deployment strategy that minimizes impact to players.

You should incorporate a deployment strategy for your game software and infrastructure that minimizes the amount of downtime that keeps players out of your game. While certain types of updates might require installing new updates to the game client, the game backend should be designed to minimize or eliminate the need for downtime during deployments.

One of the most important steps to consider when developing a game deployment strategy is to determine how your game infrastructure will be managed. It is a best practice to manage your game infrastructure using an infrastructure as code (IaC) tool such as [AWS CloudFormation](#) or [Terraform by Hashicorp](#) to reduce human errors during environment preparation. Infrastructure templates can be deployed and tested in automated pipelines, which allows you to create consistency in the configuration of different game environments.

There are several deployment strategies that can be used for a game:

Rolling substitution: The primary objective of a rolling substitution for deployment is to perform the release without shutting down the game and without impacting any players. It is important

that the upgrade or changes that are to be performed are backward compatible and will work adjacent to the previous versions of the system.

As the name suggests, in this deployment the server instances are incrementally replaced (substituted or rolled-out) by instances running the updated version. This rolling substitution can be performed in a few different ways. For example, to implement rolling updates to a fleet of dedicated game servers, a typical approach involves creating a new Auto Scaling group of EC2 instances that contain the new game server build version deployed onto them, and then gradually routing players into game sessions hosted on this new fleet of servers. If there is an associated game client update that is required as a prerequisite in order to use the new game server build, then you must include a validation check to ensure that only players that have this new game client update installed are routed into these game sessions.

Server fleets (for example, EC2 Auto Scaling groups) containing the old game server build version are only removed from service after they are drained of all active player sessions in a graceful manner, typically by setting up per-server metrics that allow game operations teams to automate this process. Alternatively, to reduce the amount of infrastructure and time to conduct a rolling deployment, an approach can be performed where existing production instances are removed from service, updated with the new game server build, and then placed back into the production fleet. This approach reduces the amount of infrastructure that is required, but it also increases risk since the number of available live game servers for players is reduced as servers are being replaced.

This model can also be used for performing rolling deployments to backend services such as databases, caches, and application servers that don't host gameplay. As long as these services are deployed in a highly-available manner with multiple clustered instances, then the complexity of deployments to these services should be less than deployments to dedicated game servers.

Blue/green deployment: The primary objective of a blue/green deployment in a game is to minimize downtime while also allowing safe rollback to the previous deployment if any issues are identified. It is suitable for deployments where two versions of the game backend are compatible and can serve players simultaneously. In the blue/green deployment strategy, two identical environments (blue and green) are set up and the existing game version is labeled as blue, while the new game version that is the deployment target is labeled as green. When the green environment is ready for migration, you can configure your routing layer to flip the traffic over to the green environment while keeping the old environment (blue) available for a period of time in case failback is needed. In this scenario, the routing updates might require updating the matchmaking service to configure it to begin sending game sessions to the new fleet, or in the

case of game backend services this could be updating DNS records in Route 53 for your service or [shifting application load balancer weights](#) to send traffic to your new target group.

One of the drawbacks of the blue/green deployment strategy is the inherent cost of the standby environment due to the additional infrastructure that is required while performing the deployment. An option to mitigate this additional infrastructure cost is to consider adopting a variant of blue/green deployment where new game software is deployed onto the same servers that are already deployed into production. In this scenario, a new green server process can be started with the new software alongside the existing blue server process, with the cut over happening between server processes rather than between completely separate physical infrastructure. This approach can also speed up game deployments across a large amount of infrastructure by removing the need to wait for new servers to be launched in the cloud. Refer to the [Blue/Green Deployments on AWS](#) whitepaper to learn more about best practices for this deployment approach.

Canary deployment: Canary deployment is of particular interest to game developers, as the strategy can be applied to release an early alpha/beta build of a game, or a game feature like a new game mode/map/challenge to a restricted or small set of players in-production. Such a deployment is called a canary. The release may have additional tracking and reporting, so when real players play that game/feature, their game play telemetry is collected and analyzed for anomalies/issues. For new features, the players are not always notified about this, and the game telemetry is the primary source used to determine if players are experiencing issues, and if the release should be rolled-back. At the same time, if no significant issues are identified, the feature can then be further rolled-out to more players for additional data. If the players are notified, then they can be asked to provide regular feedback about their experience. Such test activity would ideally be coordinated by a live operations team.

As a strategy, canary deployment can also be used for standard releases, to gradually make a new feature available to the players. A potential advantage over the standard blue/green environment is that a full-scale second environment does not need to be set up. The capacity of the new scaled-down environment determines how many players are to be onboarded to the new feature. Before adding more players, the capacity has to be scaled appropriately. Even if this customized blue/green technique is expected to cost comparatively lesser than standard blue/green, it is still estimated to incur cost that may be higher than the rolling substitution technique of canary deployments.

It is recommended to run only a single canary on a production environment, and to focus it for its data and feedback. If multiple canaries are deployed, it complicates troubleshooting and isolating of issues in production, and impairs the quality of the datasets and feedback being collected.

A variation in the canary is when one or more experiments (generally UI tests) are run via targeted deployments, where one set of the game backend servers serve one version of a feature, and another same sized set serve another version of the same feature. No additional/special infrastructure is spun up for this, and only the chosen backend servers receive these updates. The outcome of the experiments is to observe how players react to each of the versions of the same feature, and if there is a consensus of overall like or dislike, and to observe if there are any issues identified with its usability or functionality, and other such intended results. Such strategic experiments are also called A/B tests, and the overall process is called **A/B testing**. On completion of these experiments, any necessary test data is collected before reverting to the current version of the game backend system on the servers used for the tests.

Legacy traditional deployment: In the traditional style of deployment, during a scheduled maintenance window the game is shut down and all connected players are dropped or drained before all server instances within the game backend are updated with the latest code builds. This deployment impacts all players each time it is performed, and the players must be notified ahead of the schedule. As a result, this model results in the most player impact and should be avoided whenever possible. After the game update is deployed, the game can be smoke tested prior to opening up the game to the players, who would be waiting for the game to reopen. This can cause a spike of traffic when all players try to login and play within a short period of time. Therefore, if the game is not designed to handle such spikes of traffic, you can choose to gradually allow players back into the game in batches. Alternatively, you can opt to over-provision the infrastructure to sustain the opening spike of traffic, and eventually after the game traffic settles, resources can be scaled down. It is recommended that this type of deployment be conducted during off-peak hours when the number of players is at its lowest. Frequently scheduled maintenance, as well as extended maintenance, inherently carries a risk of player attrition and potential loss of revenue. Players also expect changes after a new release and can lose trust in the game once returning after a period of downtime.

GAMEOPS_BP06: Pre-scale infrastructure required to support peak requirements

You should scale infrastructure ahead of large-scale game events to make sure that you can handle the sudden increase in player demand.

In addition to new game launches, live games typically run in-game events, promotions, new content, and season releases as examples of ways to sustain and improve player engagement. Such

activities experience a high volume of player traffic for the duration of the event or promotion. The business expects to hit or surpass their intended targets for the event, and the game infrastructure must sustain and support them through it. It is important to prepare your infrastructure ahead of time to be able to support the anticipated player load that you will experience during large scale events. To prepare, game operations teams should coordinate with stakeholders in sales and marketing to estimate the projected demand that will be generated in an upcoming event by looking at past player concurrency, engagement metrics, and sales data. If the event is for a new game launch, game operations teams should work with these stakeholders to identify realistic projections for what scale they anticipate. While it may be difficult to predict how successful a game will become, it is important that everyone understands what the expectations are for success so that the infrastructure can be scaled and tested to support those goals.

Many games choose to launch in stages, starting with a soft-launch by opening the game to a small number of players, and then organically scaling the players at every stage, prior to a full public launch. During the soft-launch period, it is recommended to monitor, identify, track, and resolve any issues while refining your projections for the public launch.

To properly estimate infrastructure requirements, you should collect data through load and performance tests run against your game backends running on production, or a production-like staging environment, prior to the game launch. Multiple rounds of these tests should be run to simulate different conditions of the game, and validate that the backend is able to withstand the load under all conditions. To achieve this, developers can write game play bots that traverse various workflows in the game, and emulate different conditions. It is imperative that these tests inspect all the different system layers of the game backend so that each layer and component is tested and the details are recorded. The data collected from these tests is used for the provision planning for the game launch.

Single points of failure (SPOF) should be removed, when possible, by making the application more highly available and fault tolerant. Use load tests to find any SPOFs by emulating failures at different upstream and downstream layers, and verifying game and other component behavior.

Along with the necessary estimated infrastructure to be provisioned for the game launch, or in-game event or promotion preparations, the system should also be set up to automatically scale on-demand. Scaling event thresholds must be defined, configured, and monitored to allow the game backend to scale to sustain a high volume of player traffic. For variable traffic, pre-provisioning is best because there may not be enough time to scale-out. Manual scaling might be required during initial game launches that drive higher than anticipated demand faster than automated systems can scale resources.

On AWS, organizations should request higher [Service Quotas](#) for the services that they use in the game backend. Service Quotas are set up for all accounts to safeguard customers from inadvertently standing up or scaling more infrastructure than intended. There is no cost for higher quotas until additional resources are consumed. When a game running in an account hits the upper limit of the configured service quota in that Region, the service throttles all the requests beyond the provisioned quota, and any burst provisions. Throttles can cause unintended or unexpected errors, and impair the player experience. It is essential to monitor, track, and regularly review service quota thresholds for all the services used by the game in-production to avoid throttling. When the usage crosses a tolerable service quota threshold, an increase in the quota can be requested by raising an [Support Case](#) from the Console Support Center, after logging in to the affected account, or via the [Support API](#).

If you are launching a game hosted on Amazon GameLift, review the [pre-launch checklists](#) to help you prepare.

GAMEOPS_BP07: Conduct performance engineering and load testing before launch by simulating player behavior

To prepare for a launch, you should develop gameplay simulations that can be tested at scale against your infrastructure to ensure that you can scale to meet your peak usage requirements.

Performance engineering is the process of monitoring multiple key operational metrics of an app to discover optimization opportunities that can further improve the app's performance. This is an iterative process that starts with testing, followed by optimizing code, its dependencies, associated processes, its host operating system, and the underlying infrastructure.

To conduct a deeper analysis of the app's performance, it is recommended to integrate an Application Performance Monitoring (APM) or debugging tool in the app code that can isolate issues and reduce troubleshooting time by tracking its behavior for anomalies across all flows of the app. APM tools are also able to identify slow performing methods and external operations.

[AWS X-Ray](#) helps developers with their performance engineering activities, like identifying performance bottlenecks and analyzing and debugging production errors. You can use X-Ray to understand how your application and its underlying services are performing, and identify and troubleshoot the root cause of performance issues and errors. Through numerous rounds of load tests, in which the app and its infrastructure is gradually loaded with synthetic player traffic, various system bottle-necks, app errors, exceptions, OS problems, and other issues are identified that may have not been found during other Quality Assurance (QA) tests.

To simulate artificial player traffic, you need bots that emulate the game client flows and transact with the game backend to simulate real-world player behavior. Generally, this data is captured through game play logs and data generated by human QA tests, as well as through real-world limited scale alpha or beta tests where real players are invited to play the game.

It is also recommended to perform load testing and inject different kinds of failures in the game backend during these various load tests to check how the system behaves for each failure. It is important to record the system's behavior in an operational runbook to assist in troubleshooting possible failures in the future. It is important to have human testers test the game against the same environment that is being load tested while it is being load tested. Humans can catch something during the load test that bots or other metrics do not.

For critical events like game launches, and even major in-game events or promotions, [AWS Infrastructure Event Management \(IEM\)](#) is available. An IEM offers architectural and operational guidance, along with operational support during the preparation and deployment of your planned event like a game launch, or a major in-game event, promotion, or migration. Through the IEM process, assists in assessing operational readiness of your system, identifying and mitigating risks, and deploying your event confidently with all additional experts.

[AWS Fault Injection Simulator](#) is a fully managed service for running fault injection experiments on that makes it easier to improve an application's performance, observability, and resiliency. Fault injection experiments are used in chaos engineering, which is the practice of stressing an application in testing or production environments by creating disruptive events, such as sudden increase in CPU or memory consumption, observing how the system responds, and implementing improvements. Fault injection experiment helps teams create the real-world conditions needed to uncover the hidden bugs, monitoring blind spots, and performance bottlenecks that are difficult to find in distributed systems.

Operate

GAMEOPS04 - How do you monitor the health of the game?

GAMEOPS_BP08: Instrument the game to detect and monitor player-impacting issues.

In addition to responding to social media and player reports of issues, you should instrument your game with monitoring solutions to you detect when there are player-impacting issues that need to be investigated.

No amount of testing is sufficient enough to identify all issues in a game. Games are usually launched with known issues that are planned to be gradually fixed with the next release of the game. Known and reproducible issues are easier to address and fix. It is difficult to preempt all issues or bugs that a player may experience. To assist with identifying such issues, it is recommended that game clients implement logging and reporting in the app at various strategic locations to help the backend team identify any client-side issues. The ability to find such issues early allows the game developers to troubleshoot and fix the issue before it becomes widespread. The data and logs reported by the tracking code should never include any player identifiable information (PII), and they should only contain game specific metadata that help with debugging.

Implement an observability solution for detecting and responding to issues such as game crashes or bugs. You can use [Amazon CloudWatch Synthetics](#) to create canaries that can monitor the health of your player-facing backend game services. You can instrument your backend services with [X-Ray](#) to trace requests across distributed services, and send your custom logs and metrics to [Amazon CloudWatch](#). Third-party solutions, such as [Backtrace.io](#) and [Sentry](#), are popular solutions for error reporting in games. Application performance monitoring (APM) solutions from partners such as [New Relic](#), [Splunk](#), [Datadog](#), and [Honeycomb.io](#) are also popular.

The game's Live Operations' team and community managers should also monitor various social networks and channels to check for player feedback, complaints, and bug reports in addition to the official support channels. Every game specific complaint received should be reviewed and attempted to be reproduced, or sent to the QA team for their review. If reproducible, the issue should be escalated to the game developers for their troubleshooting and a fix before it impacts the larger player base.

Evolve

GAMEOPS05 - How do you optimize your game over time?

GAMEOPS_BP09: Monitor key game metrics that can help identify player trends and patterns; use the information to rebalance the game, optimize game design, and improve the infrastructure.

In addition to game client system usage, app usage, exceptions, and crash data, it is highly recommended that you capture game telemetry data that is sent to a game backend system. This data should represent player activity so that you can understand how players interact with various features in the game.

Depending on its implementation, game clients can collect telemetry data at predefined game features or locations in a game world. The data is sent to the backend ingestion service for processing. If the backend service is unreachable for any reason, the clients can store the data locally on the local device until the backend service is available again. The game designers use this telemetry data to review how players are playing the game, and if there are any anomalies in the game. For example, player movements and interactions with items in a map can be extracted from telemetry data and plotted as a heat map of activities in-game by all players over a set window of time. Such data helps the game designers identify the need to balance various elements in the game, such as the power of a weapon, the power of an in-game character, or the complexity of a map. The raw telemetry data is generally stored and then processed to extract analytics that can be visualized by analysts.

The [Game Analytics Pipeline](#) solution implementation helps game developers launch a scalable serverless data pipeline to ingest, store, and analyze telemetry data generated from games and services. The solution supports streaming ingestion of data, allowing users to gain insights from their games and other applications within minutes.

For custom game telemetry data ingestion, storage, processing and analytics, AWS also offers a number of [specialized services for big data processing and Analytics](#).

Resources

Refer to the following resources to learn more about our best practices related to operational excellence.

Documentation and blogs

- [Architecture best practices for Game Tech](#)
- [Establishing your best practice Environment](#)
- [multi-account strategy for your Control Tower landing zone](#)
- [Game Analytics Pipeline](#)
- [How to setup a CI/CD pipeline on](#)
- [Other relevant CI/CD blogs](#)
- [Game DevOps made easy with Game-Server CD pipeline blog](#)
- [Amazon GameLift Prepare for launch](#)
- [Best practices with Organizational Units](#)
- [AWS X-Ray](#)

- [AWS Infrastructure Event Management \(IEM\)](#)

Partner solutions

- [New Relic](#)
- [Splunk APM](#)
- [Backtrace.io](#)
- [Sentry](#)
- [Datadog APM](#)
- [Honeycomb.io](#)

Whitepapers

- [Organizing your Environment using multiple accounts](#)
- [Scalable gaming patterns](#)

Video content

- [YouTube series: Building games on AWS](#)
- [Re:Invent 2017: How Amazon scales its infrastructure to handle billions of transactions](#)
- [Re:Invent 2019: Scaling up to your first 10 million users](#)

Training materials

- [Curriculum - Game Tech starter pack](#)

Security

The security pillar includes the ability to protect information, systems, and assets while delivering business value through risk assessments and mitigation. Due to their global visibility and large number of players, games are desirable targets for exploiters, hackers and others looking for ways to exploit and abuse systems, which can result in disappointing player experience and increased costs for the game developer. As described in the [Shared Responsibility Model](#), it's important to understand which aspects of security are the responsibility of AWS, and which aspects are the

responsibility of the customer so that you are prepared to maintain a strong security posture. This pillar provides best practice cloud security guidance for you to consider when developing and operating games in the cloud.

Design principles

In addition to the design principles from the security pillar of the Well-Architected Framework whitepaper, the following design principles can help you strengthen security of your game workload in the cloud:

Monitor and moderate player usage behavior: Capture and analyze usage data to help you understand how players interact with your game and social features. By analyzing this data, you can detect and respond to abusive and inappropriate behavior that can degrade the player experience.

Definition

There are five best practice areas for security in the cloud:

- Identity and access management
- Detective controls
- Infrastructure protection
- Data protection
- Incident response

Before you architect any system, you must establish security practices and access control. Additionally, you should be able to identify security incidents, protect your systems and services, and maintain the confidentiality and integrity of data through data protection. You should have a well-defined and practiced process for responding to security incidents. These tools and techniques are important because they support objectives such as preventing financial loss or complying with regulatory obligations.

Best practices

There are five best practice areas for security in the cloud.

Topics

- [Identity and access management](#)

- [Detective controls](#)
- [Infrastructure protection](#)
- [Data protection](#)
- [Incident response](#)

Identity and access management

GAMESEC01: How do you manage player identity and access management?

When developing a game, you must determine how you will provide players with access to your game and related services. This decision is influenced by player acquisition and monetization strategy, player experience, and other factors such as the existing capabilities that might be provided by your game publishing partners. For example, a game might require purchases and require a player to create a user profile to associate real-money payment methods with their account.

Alternatively, a game may desire to reduce the barrier to entry for first-time player experiences by removing the need to create a user account before playing the game, thereby improving the chance that a player will try the game for the first time. Typically, games will implement one or more combinations of player identity and access management approaches for their game.

Unauthenticated or anonymous access: This access level is useful in situations where a game does not require a player to create a new user account or link with their identity on social networks and gaming platforms. This is the simplest and quickest way for player to start to play a game and is particularly useful in mobile games where a game developer may want to reduce the barrier to entry for the initial first time experience. In this access scenario, if you want to identify usage from the game installation, you can do so by programming the game client to generate and store a unique identifier onto the player's device that is used to identify the player across game sessions on their device and allow analytics reporting on usage over time. At a later date, if a player chooses to create an account, then you can associate their new user account to their previously generated unique identifier to link the player's historical usage such as stats and game achievements to their new player identity. If a player does not eventually create and link an account, the device that the player uses to interact with the game can be uniquely identified, but any recoverable information about the player is not collected and stored. Thus, if the player breaks the device, the previous stored data associated with the device is also lost and might not be recoverable.

Authentication with username and password: A game may allow players to create their own user accounts with a username and password that is stored within the game's backend. Typical reasons for this are when a game developer is collaborating with a game publisher who already has an existing player account system that the developer can integrate with, or when a developer who may publish their own games wants to simplify the player experience by allowing players to create a single user account for access across all of the game they publish.

Authentication and account linking with third-party social networks and gaming platforms: It is common for online games and games with social features to provide third-party identity provider federation as a way to simplify the player experience. Instead of asking players to create a username and password combination to authenticate, you can use identity federation to allow players to authenticate using their third-party accounts with social networks and gaming platforms. This login process simplifies sign-in and registration experience and it provides a convenient alternative to mandatory account creation and a frictionless method for players to access games. For game developers, federated login can offer a streamlined player verification workflow and may also provide a more reliable way to manage player data for personalization because you do not need to ask players to provide you with certain data that they likely have already provided to the third-party identity provider. Additionally, these platforms provide integration with additional social features such as allowing you to link players with their friends.

The following best practices can help you to incorporate secure access control for your game:

GAMESEC_BP01: Requests to game backend services are authenticated.

Any requests sent to game backend services should be authenticated to prevent unwanted requests to your game backend.

You should provide an authentication service for players to login, which should return secure short-lived tokens, such as JSON Web Token (JWT), to the game client when a player successfully authenticates. These tokens can include claim assertions containing player attributes and other relevant metadata that can be used in subsequent requests sent from the game client to your game backend to authenticate requests and authorize them in the context of the authenticated player. You have the option to either design and build your own player authentication system, which would require ongoing improvement and maintenance, or you can take advantage of the scalable and secure user sign-up, sign-in, and access control features provided by [Amazon Cognito](#). Amazon Cognito User Pools allows you to create a user directory and provides APIs that you can integrate into your game for sign-up, sign-in, password reset workflows, and it can be integrated with third-party identity providers. [Application Load Balancers](#) and [Amazon API Gateway](#) both

provide integrations with Cognito to easily integrate user authentication for requests sent to your custom game backends hosted with these services.

If your game supports anonymous access and you cannot authenticate a player, you can use a client authentication approach to provide a secure experience when integrating with your game backend. If your game client uses AWS services, then requests to these services must be signed using credentials. To provide credentials to your game client for unauthenticated users, you can use the AWS SDK to retrieve short-lived credentials from [Amazon Cognito Identity Pools](#) that can be used to sign your requests to AWS services. These credentials can be refreshed from your game client. In addition to directly integrating with the AWS SDK from the game client, you can also build your own game backend, using a service such as Amazon API Gateway, which supports custom authorization. By designing your own game backend service, you can gain authoritative control over all requests with custom server-side logic. For more information on building a backend service for games hosted using Amazon GameLift, refer to [Design your backend service](#).

GAMESEC_BP02: Player requests to join a multiplayer game should be validated by your game backend service.

Typically, in multiplayer games, a player will join a game session by selecting an option directly from a list of available sessions, or they will submit a request to find a match, which places the responsibility on the game developer to locate an eligible game session and provide the connection information, usually an IP address and port, back to the player's game client. The implementation may vary depending on the genre of game you are developing, but regardless, it is a security best practice to perform server-side validation of a player's game join request.

For example, in a session-based multiplayer game, a request from a player to join a game session should be validated by your game server software with your game backend matchmaking service before authorizing their connection to the server. When a player requests to join a game session, the game server should check the request for a unique identifier, such as a player session ID and server-generated ticket that was previously provided to the game client by your game backend matchmaking service. Upon initiating the connection to the game server, your server-side software can use this information to verify with the matchmaking service that the player's connection request is valid, and ensure that the player is not joining a spot previously reserved in the game session for another player. For games hosted using Amazon GameLift, refer to the [Amazon GameLift documentation](#) for an example of how this type of server-side validation can be implemented.

GAMESEC_BP03: Player user accounts that require passwords should enforce a strong security policy.

If a game provides players with the ability to create a user account with a password, you should require that passwords adhere to strong policies. For example, Amazon Cognito User Pools provides support for [defining password requirements](#) for user accounts.

GAMESEC_BP04: Provide an option for players to setup multi-factor authentication (MFA) on their accounts.

Player accounts can be a valuable asset to bad actors, particularly in games that support in-game currency and purchases. Due to the pervasiveness of player account hacking and social engineering against game developer player support teams, it is important to provide players with the option to configure multi-factor authentication (MFA) for their user account to enhance the security of their accounts.

When a player accesses their account using MFA, they are challenged with a temporary code that is sent to their email, phone number, or a purpose-built multi-factor authentication mobile app, which they must enter to login within a limited time frame to successfully authenticate. MFA can also be used to protect an account that is attempting to authenticate from a new geo-location, accounts that have been flagged by player support for potential malicious activity, and even for accounts that have not logged into the game for an extended period of time. For example, Amazon Cognito User Pools provides support for [configuring multi-factor authentication](#) on user directories.

GAMESEC02: How do you prevent unauthorized access to game content?

Modern games include a significant amount of content, including downloadable content (DLC), which is an important aspect of player engagement and game monetization. Players expect an ongoing stream of new characters, levels, and challenges, requiring game developers to keep up with this constant demand for fresh content in order to retain players. The variety and the size of the content can vary greatly depending on the type of the game and also the device the game is played on whether it is PC, console, or mobile. It is important to ensure that game content is protected from unauthorized access.

GAMESEC_BP05: Restrict access of downloadable content to authorized clients and users.

Access to game content should be restricted to authorized applications and clients.

Consider using Amazon S3 as a cost-effective and scalable origin for storing downloadable game content, and Amazon CloudFront to provide globally performant content delivery to players. Both services provide built-in mechanisms for restricting access to data that is stored.

When you need to grant access to S3 content, there are several best practices that should be considered. By default, only the AWS account that created an S3 bucket can access the objects stored within it. To grant access to your internal applications and to manage content stored in Amazon S3 buckets, use [Identity and Access Management \(IAM\)](#) to create policies that provide appropriate access. [IAM Roles](#) can be associated with federated users, systems, or applications hosted in services, such as Amazon EC2, Lambda, and container-based applications hosted in Amazon EKS and Amazon ECS. For example, you might use the AWS SDK or AWS CLI to publish and manage game content assets in S3 buckets. To support this use case, you can create an IAM Role with appropriate access to read and write game content to your S3 buckets and associate it the EC2 Instances that host your software and scripts.

Resource-based policies can also defined for your bucket, and for specific objects. [S3 Bucket Policies](#) are associated with the S3 bucket and can be used to restrict access to the bucket and objects within it, as well as grant access to your S3 resources from other accounts. For example, in scenarios where multiple teams or separate game development studios are working on the same game content and require the same access to centrally hosted content in S3, you can use an S3 bucket policy to define permissions for cross-account access to the S3 resources. Consider using [S3 Access Points](#) which can simplify managing data access to shared data by creating access points with names and permissions specific to each application or sets of applications. The Amazon S3 documentation contains additional [best practices for access control in Amazon S3](#).

It is recommended to generate temporary URLs that grant short term access to your content. Amazon S3 provides support for generating [presigned URLs](#), which allow object owners to share objects with others by generating a presigned URL using their own security credentials within their backend that grant time-limited permission to download the objects. By doing so, the end user or application that is being granted access is not required to have an account or IAM permissions, and instead uses the presigned URL to access the content. This is a best practice that is commonly used in a variety of games use cases, such as granting authorized players access to downloadable content that they have been entitled to, and providing temporary access to limited time game content. Presigned URLs can also be used to provide temporary permissions to upload content to your S3 bucket, for example to provide a player with access to upload client logs to assist your player support team during troubleshooting a player support case. provides best practices for [limiting presigned URL capabilities](#).

[S3 Block Public Access](#) is a set of security controls that ensures S3 buckets and objects do not have public access by overriding the settings defined by specific users to enable centralized control across your S3 access points, buckets, and AWS accounts.

While your applications, game developers, artists, and other personnel may need direct access to the content in S3 buckets for development and management purposes, it is recommended to use a content delivery network to provide access to content available publicly to players or other users over the internet. This improves download performance and reduces costs by caching frequently accessed content. Amazon CloudFront can globally distribute your content by caching and delivering it closer to your players while reducing the load on your game's download origin, such as Amazon S3.

Rather than serving your public content directly from S3 buckets, it is recommended to keep this content private and serve it publicly using CloudFront, which can be configured to require that players access your private content, for example a new game download for paid players only, using either [signed URLs](#) or [signed cookies](#). You then develop your application either to create and distribute signed URLs to authenticated users, or to send set-cookie headers that set signed cookies for authenticated users. When you create signed URLs or signed cookies to control access to your files, you can specify an ending date and time, after which the URL is no longer valid. Optionally, you can also specify the IP address or range of addresses of the computers that can be used to access your content, which is useful if you want to restrict access to specific game development studio partners or contractor networks. Use signed URLs in the cases where you want to restrict access to individual files or your users are using a client that doesn't support cookies. Use signed cookies in the cases where you want to provide access to multiple restricted files, or you don't want to change your current URLs. Signed URLs take precedence over signed cookies.

GAMESEC_BP06: Limit origin access to authorized content delivery networks (CDNs).

You should prevent users from circumventing your content delivery networks to directly access content from your origin, such as your Amazon S3 buckets. It is important to restrict access to your origin to only your authorized CDNs, which helps to reduce data transfer costs from unnecessarily serving content out of the origin. It also improves security posture by ensuring that all public access to your origin content flowed through the same entry point, where you can deploy edge security controls such as AWS WAF layer-7 filtering, injection and inspection of security-related HTTP request parameters, and Distributed Denial of Service (DDoS) protections. To implement these controls for an S3 origin, you can use [Amazon CloudFront origin access identity \(OAI\)](#) which ensures that all requests to your S3 objects are originating from your CloudFront Distribution. It is recommended to associate AWS WAF with your CloudFront distribution in order to provide

layer-7 filtering. However, if you are serving content from additional CDNs, you can configure the CDN to insert one or more custom HTTP headers into origin requests which can be inspected by AWS WAF to ensure that the incoming traffic originated from your authorized CDN provider. This approach is also useful to prevent users from circumventing your CDN providers when your origin is hosted behind an [Application Load Balancer \(ALB\)](#). ALBs can be associated with AWS WAF for layer-7 protections, and AWS WAF can be configured to insert a custom HTTP header that can be inspected by your ALB to ensure that incoming traffic to the load balancer was first processed and inspected by AWS WAF.

GAMESEC_BP07: Implement geo-restrictions to prevent unauthorized access.

When a player requests your content, CloudFront serves the requested content from the nearest edge location, regardless of where the player is located. However, there may be scenarios where you need to restrict how your content is accessible by users in specific parts of the world. For example, you may have a rolling game deployment strategy that releases content in phases on a country-by-country basis, or you may have to abide by country-specific access controls. You can use geo restriction, also known as geo blocking, to prevent players in specific geographic locations from accessing content that you're distributing through a CloudFront distribution. You can use the CloudFront geo restriction feature to restrict access to all of the files that are associated with a distribution and to restrict access at the country level. Alternatively, you can use a third-party geolocation service to restrict access to a subset of the files that are associated with a distribution or to restrict access at a finer granularity than the country level.

Using CloudFront geo restriction, you can allow your players to access your content only if they're in one of the countries on an allow list of approved countries and prevent your players from accessing your content if they're in one of the countries on a deny list of banned countries. If a request is received from a blocked geographic location, CloudFront will return 403 Forbidden HTTP status code to the player.

GAMESEC_BP08: Restrict access to content with digital rights management (DRM) solutions.

In addition to the access-control based approach, you can also take an encryption-based approach by encrypting your private content and distributing the decryption keys to authorized players using a digital rights management (DRM) solution. DRM solutions are recommended in situations where you want to allow players to download game content early, but you do not want them to be able to be able to access or play the content until a predetermined time. For example, this is common in situations, typically in PC games, where players are allowed to pre-order a game and configure their game client to automatically begin downloading the encrypted files early so that it is downloaded and ready to be played when the game is officially released. After the game is

released, the player's game client can request decryption keys from the DRM backend solution so that it can decrypt the previously downloaded files and begin playing the game. DRM systems are also used as a way to prevent unauthorized re-distribution and manipulation of games after they have been downloaded and installed by an authorized player. DRM systems require integration with the origin for exchanging encryption keys and authorizing players to retrieve the decryption key. Commercial DRM systems providers offer a range of solutions with particular features and support for different devices.

Detective controls

GAMESEC03: How do you monitor and analyze player usage behavior within your game?

To maintain a positive player experience, you should have a process for capturing, storing, and analyzing relevant data that can help you understand how players engage with your game's features and with other players.

GAMESEC_BP09: Collect, store, and analyze player usage logs to detect inappropriate behavior.

Instrument your game to collect logs that help you understand how players use the features of your game and how they interact with other players so that you can prevent unauthorized activity which can degrade the player experience. This can be done by sending structured log events to the [Game Analytics Pipeline](#), or by using a logging solution such as [Amazon CloudWatch Logs](#), [Amazon OpenSearch Service](#), or a solution from an AWS Partner such as [Datadog](#), [Sumo Logic](#), [New Relic](#), [Honeycomb](#), or [Splunk](#). These player usage logs should be structured so that they can be used to detect when specific actions by players need to be investigated.

After you have captured the data, you should consider implementing tools to help you detect inappropriate usage behavior. For example, if your game has social features such as in-game player messaging and voice chat, or online forums, it is recommended to save logs from these player engagements in a format that can be analyzed for moderation purposes. Configure your game's voice chat feature to export recordings to Amazon S3 and use [Amazon Transcribe](#) to convert the audio speech to text format which can be stored for processing. Alternatively, you can perform real-time streaming transcription by integrating your game backend voice chat service directly with the Transcribe API to [transcribe streaming audio](#) in real-time. Moderation teams can manually review the content, and once the content is in a standard format, you can also use AWS AI/ML services to perform moderation automatically. [Amazon Comprehend](#) can be used to perform natural language processing (NLP) to uncover information from the unstructured text, which can

help you classify and organize the conversations into relevant topics and identify inappropriate behavior such as profanity.

If your game allows players to generate or upload content, consider using [Amazon Rekognition](#) to identify the content of the images for moderation. For video use cases such as player live streaming, you can send video streams to [Amazon Kinesis Video Streams](#) which you can integrate with [Amazon Rekognition Video](#) or your own custom application to analyze and moderate in real-time. Your game may provide players with the ability to contact player support agents through a call center such as [Amazon Connect](#), or chat bots using Amazon Lex. Amazon Connect provides support for [monitoring live and recorded conversations](#). To analyze interactions between players and player support chat bots built with Amazon Lex, you can store the [conversation logs](#) from these interactions in Amazon CloudWatch Logs which can be exported to S3 and analyzed as described previously.

You can also integrate your game with [Amazon Fraud Detector](#), a fully managed service that uses machine learning to identify potentially fraudulent activity so customers can catch online fraud quickly. You can use Fraud Detector to detect potentially fraudulent activity and flag that activity for review so that you can prevent fraudulent in-game purchases in real-time, detect compromised accounts by looking for behavioral changes and anomalies, and distinguish between legitimate and high-risk new account registrations.

[Amazon Lookout for Metrics](#) uses machine learning to automatically detect and diagnose anomalies in your business and operational data, and monitors the metrics that are most important to your businesses with greater speed and accuracy. The service also makes it easier to diagnose the root cause of anomalies such as sudden dips in revenue, logins, transactions, and retention. It does not require game developers to have any ML experience to setup and can connect to popular data sources including Amazon S3, Amazon CloudWatch, Amazon RDS, Amazon Redshift, as well as many SaaS applications. For example, you can [integrate Amazon Lookout for Metrics with the Game Analytics Pipeline](#) and other data sources to begin analyzing behavior to detect anomalies.

Alternatively, you may choose to build, train, and host a custom machine learning model using [Amazon SageMaker](#) to address use cases such as content moderation, toxicity detection, cheat detection, fraud detection, and more.

In addition to generating custom game usage logs, it is also recommended to capture and store system-level logs from relevant services, such as [S3 server access logs](#), [CloudFront access logs](#), and [ALB access logs](#). These logs can be stored in an Amazon S3 bucket in your account and are useful for associating your player usage information from within the game with system-level information including connection details such as IP addresses, request headers, and any relevant

request manipulation and filtering that you may have configured within your game backend. These logs can be sent to the same logging solutions mentioned earlier, and can also be [analyzed using SQL queries with Amazon Athena](#) without requiring the logs to be moved out of Amazon S3.

[Access Analyzer for S3](#) is a feature that monitors your bucket access policies, ensuring that the policies provide only the intended access to your S3 resources. Access Analyzer for S3 evaluates your bucket access policies and allows you to discover and swiftly remediate buckets with potentially unintended access.

To continuously monitor for malicious activities and unauthorized behaviors within your AWS environment, consider using [Amazon GuardDuty](#). GuardDuty identifies threats by monitoring account behavior, network activity, and data access patterns within your environment. It analyzes tens of billions of events across multiple data sources, such as CloudTrail event logs, Amazon VPC Flow Logs, and DNS logs for potential threats. By integrating with Amazon CloudWatch Events and Lambda, GuardDuty alerts can be automatically forwarded to relevant security teams for further analysis.

[AWS Security Hub](#) provides a comprehensive view of your security state in AWS and helps you to check your environment against security industry standards and best practices. Security Hub collects security data from across AWS accounts, services, and supported third-party partner products and helps you to analyze your security trends and identify the highest priority security issues. The [Amazon GuardDuty integration with Security Hub](#) enables you to send findings from GuardDuty to Security Hub. Security Hub can then include those findings in its analysis of your security posture.

It's common for bad actors to employ bots to take over accounts and cheat in games. [WAF Bot Control](#) gives you visibility and control over common and pervasive bot traffic that can consume excess resources, skew metrics, cause downtime, or perform other undesired activities.

Ransomware is malicious code designed to gain unauthorized access to systems and datasets and encrypt that data to block access by legitimate players. After ransomware has locked players out of their systems and encrypted their sensitive data, cyber criminals demand a ransom before providing a decryption key to unlock the data. Organizations can be completely shut down by an attack, incurring significant costs and loss of business productivity. Refer to [Securing your Cloud Environment from Ransomware](#) for best practices you can apply to strengthen your ability to fight ransomware before, during, and after an incident takes place.

Refer to the Well-Architected Framework whitepaper for additional best practices in the detective controls area for security.

Infrastructure protection

Refer to the Well-Architected Framework whitepaper for best practices in [infrastructure protection](#) for security that apply to games workloads.

Data protection

Refer to the Well-Architected Framework whitepaper for best practices in [data protection](#) for security that apply to games workloads.

Incident response

GAMESEC04 - How are you defining and enforcing policies to respond to player misconduct and abusive behavior?

GAMESEC_BP10 - Implement an incident response plan to handle bad actors and abusive behavior.

Refer to the Well-Architected Framework whitepaper for best practices in [incident response](#) for security that apply to games workloads.

GAMESEC_BP11 - Ban accounts associated with bad actors.

If left unmitigated, abusive behavior in a game can continue to cause impact to the gaming experience for others and should be mitigated as soon as possible. You should implement a process to impose bans, or other forms of restrictions, on bad actors who are confirmed to be in violation of your terms of service. Typically, the rules and evaluation process for determining the circumstances for imposing these types of restrictions will be determined by personnel such as a player community team, or trust and safety team, within your organization. After bad actors have been flagged, you should have a pre-determined workflow that can be run to take action on the identified players. AWS Step Functions and Lambda functions can be used to run an automated workflow that accepts a batch of player accounts as input and updates entries in a DynamoDB table called Bans, which can include details on the player account, the ban reason, and duration. Depending on the way your game and account management system is designed, and the type of abuse, it is valuable to have a banning system of record that is separate than your account management system. You may not want to turn off the player's account from your account management system, opting instead to simply turn off their ability to play your game. This can be

useful in situations where the player's account credentials are used to access multiple games with different terms of service or policies.

Resources

Refer to the following resources to learn more about our best practices related to security.

Documentation and blogs

- [Common Amazon Cognito Scenarios](#)
- [Using signed URLs](#)
- [Use channel flows to remove profanity and sensitive content from messages in Amazon Chime SDK messaging](#)
- [Security in Amazon GameLift](#)

Whitepapers

- [Secure content delivery using Amazon CloudFront](#)
- [Security Response Guide](#)
- [Best Practices for DDoS Resiliency](#)
- [Securing your Cloud Environment for Ransomware](#)

Partner solutions

- [Datadog](#)
- [Sumo Logic](#)
- [Splunk](#)
- [Honeycomb](#)
- [New Relic](#)
- [AWS Marketplace - DRM Solutions](#)

Training materials

- [Getting Started with Amazon Cognito](#)
- [Security self-paced training](#)

Reliability

The reliability pillar includes the ability of a system to recover from infrastructure or service disruptions, dynamically acquire computing resources to meet demand, and mitigate disruptions such as misconfigurations or transient network issues.

Design principles

In addition to the design principles in the AWS Well-Architected Framework whitepaper, the following are design principles that can help you increase reliability in the cloud for games workloads:

Agree on the peak player concurrency and system scalability targets required to meet business projections: Prior to launching a game and during live game operations, develop estimates for the number of concurrent players expected at peak, and to establish target goals for system scalability to meet these projections. This helps create a baseline for your game's reliability. Define scaling policies to accommodate changes in demand automatically without impact availability, such as by ensuring that your scaling systems gracefully manage active player sessions.

Measure your reliability and the impact on player experience: Define key performance indicators (KPIs) that represent the health of your game. Monitor the impact of changes in infrastructure and game features on your reliability.

Definition

There are three best practice areas for reliability in the cloud:

- Foundations
- Change management
- Failure management

To achieve reliability, a system must have a well-planned foundation and monitoring in place with mechanisms for handling changes in demand or requirements. The system should be designed to detect failure and automatically heal itself.

Best practices

The following provide best practices for your cloud architecture.

Topics

- [Foundations](#)
- [Change management](#)
- [Failure management](#)

Foundations

Refer to the Well-Architected Framework whitepaper for best practices in [foundations](#) for reliability that apply to games workloads.

Change management

GAMEREL01 — How does your game infrastructure scale to changes in player demand?

As your player demand fluctuates over time, your game infrastructure should be able to adaptively scale to handle these changing requirements. While it is difficult to predict the popularity of a game ahead of time, you should design an architecture approach that allows you to easily add and remove infrastructure capacity to accommodate fluctuations in player population.

GAMEREL_BP01 - Implement a scaling strategy that incorporates the state of active player game sessions.

Implement a solution for automatically scaling your game infrastructure in a manner that incorporates the stateful nature of your actively connected player sessions and gracefully handles scaling activities without disrupting gameplay.

One of advantages of developing a game in the cloud is the elasticity that can be achieved by automatically scaling server infrastructure as needed to meet demand. While stateless or asynchronous games and backend services can be dynamically scaled using [Amazon EC2 Auto Scaling policies](#) or similar techniques typically adopted for scalable web applications, game developers typically require a more customized approach for scaling stateful or synchronous games to prevent disruptions to active player sessions.

For stateful games, it is a best practice for your game backend to generate custom metrics that can be used to monitor the state of your player sessions and available game server capacity, which can be reported to Amazon CloudWatch as custom metrics. Using this data, you can implement

game server scaling software, for example, as a serverless application using AWS Lambda function or AWS Fargate, to manage the fleet of dedicated game server instances by using the AWS SDK to make API calls to update the minimum, maximum, and desired capacity settings for the EC2 [Auto Scaling groups](#) hosting your game server build.

Alternatively, you can use Amazon GameLift to host your game servers and use the [out-of-the-box game server auto scaling capabilities](#) to manage this scaling process for you. Amazon GameLift's automatic scaling capabilities are aware of active player sessions and can be configured to prevent the termination or scale-in of game server instances and can be configured to prevent the termination or scale-in of game server instances that are actively hosting players. For more information, refer to [Monitor Amazon GameLift with Amazon CloudWatch](#).

GAMEREL_BP02 - Support the use of multiple EC2 instance types for your game.

When hosting your game using EC2 instances, or if you use containers hosted on EC2 instances in your AWS account, then you should use multiple instance types in your hosting strategy. By using multiple instance types, you increase the number of compute options that can be used when your game is scaling to add more servers to support player growth, which improves reliability in case your preferred instance type is unavailable. This is also a best practice when using Spot Instances to host your game, since the availability of Spot Instances fluctuates based on customer demand. You should test your game on multiple instance types to meet your cost and performance requirements and determine a prioritized ranking of instance types. Amazon EC2 Auto Scaling supports using multiple instance types and sizes as well as [assigning weights to each instance type](#) in your configuration so that you can implement prioritized ranking of compute options.

If you host your game using Amazon GameLift managed hosting, it uses Amazon EC2 instances to deploy your game servers and host game sessions for your players. When setting up a new fleet, you decide what type of instances your game needs and how to run game server processes on them (using a runtime configuration). When choosing resources for a fleet, you must consider several factors, including game operating system, instance type (the computing hardware), and whether to use On-Demand Instances, Spot Instances, or both. Hosting costs with Amazon GameLift primarily depend on the type of instances you use. For more information, refer to [Choosing computing resources](#).

Failure management

GAMEREL02 — How do you minimize the impact of infrastructure failures on active players?

You should monitor game server failure metrics and the impact of these failures on player behavior over time so that you can adjust your game server hosting strategy to meet your game's reliability requirements. Game server infrastructure that is determined to be degraded should be removed from service immediately if it is impacting players, or proactively replaced when there are no active player sessions hosted on the server.

For scenarios where games are hosted as REST APIs, system reliability can be managed similar to traditional web application architectures on where traffic can be load balanced across multiple servers in a distributed manner to mitigate the risk of server failures.

For real-time synchronous gameplay, a game session is usually hosted on a game server process running on a virtual machine, or game server instance, since gameplay state needs to be maintained in a performant manner and replicated to all connected game clients. This implementation means that a player's experience is tightly coupled to the performance and reliability of the game server process that hosts their game session. This type of architecture makes managing the reliability of game servers more complex than traditional approaches.

To mitigate the impact of a game server failure, you can configure your game to continuously perform asynchronous updates of a player's game state to a highly-available cache or database such as [Amazon ElastiCache for Redis](#), or [Amazon MemoryDB for Redis](#), respectively. If a server failure occurs, the player's last saved game state can be fetched from the external data store and their session can be restored on a new game server instance. However, this approach adds additional cost and complexity to manage this external state, and may not be suitable for fast-paced or competitive games where the state changes are so frequent and happening at such a significant scale that introducing even a performant in-memory cache data store would result in replication lag that is too significant to be useful to restore a session from. For games of this nature, the optimal approach is to accept the loss of the server and send the player back to a game lobby to find another session or you can automatically redirect them into another game session.

You should make sure to capture as much useful log data about what caused the server disruption so that you can investigate the issue later. Amazon GameLift provides guidance for [debugging fleet issues](#), and provides the ability to [remotely access Amazon GameLift fleet instances](#).

GAMEREL_BP03 - Implement loose coupling of game features to handle failures with minimal impact to player experience.

Decoupling components refers to the concept of designing server components so that they can operate as independently as possible. Some aspects of gaming are difficult to decouple since data needs to be as up-to-date as possible to provide a good in-game experience for players. However,

many components and gaming tasks can be decoupled. For example, leaderboards and stats services are not critical to the gameplay experience, and the reads and writes to these services can be performed asynchronously from the game.

You should consider how to develop features in your game that can be disabled automatically or by an administrator if issues are detected, as well as configure upstream services that depend on the feature to be able to gracefully handle the failure. For example, if specific player data is not properly loading within your game client, you should consider whether this data is critical to the gameplay experience. If not, configure the game client to gracefully handle this failure without disrupting the experience for the player, optioning to retry fetching this data at a later time when the player revisits the screen. Employ logic such as timeouts, retries, and backoff to handle errors and failures. Timeouts keep systems from hanging for unreasonably long periods. Retries can provide high availability of transient and random errors.

Define non-critical components which can be loosely coupled to critical components. Loose coupling allows systems to be more resilient since failure in one component does not cascade to others. When game features do not require stateful connections to your game servers or backend, you should implement stateless protocols to scale dynamically and easily recover from transient failures. Develop your non-critical components where it can be loosely coupled with stateless protocols using an HTTP/JSON API. It is also recommended to implement network calls from the game client to be asynchronous and non-blocking to minimize the impact to players from slow-performing game features or other dependent services.

To further improve resiliency through loose coupling, use a messaging service such as a queuing, streaming, or a topic-based system between components that can be handled asynchronously. This model is suitable for any interaction that does not require an immediate response or where an acknowledgment that a request has been registered is sufficient. This solution involves one component that generates events and another that consumes them. The two components will not integrate through direct point-to-point interaction but through an intermediate such as a durable storage or queuing layer. This also helps to improve system's reliability by preserving messages when processing fails. Selection of an appropriate messaging service is key, since various messaging services have different characteristics such as ordering and delivery mechanisms. Design operations to be idempotent so the chosen message system delivers messages at least once. As an example, consider a typical game use case where your game needs to track player playtime, stats, or other relevant data which can lead to a high write-throughput use case at times of peak player concurrency.

To implement a reliable architecture, consider whether the use case requires read-after-write consistency as perceived by the player. Typically, scenarios such as these are suitable for asynchronous processing and can be achieved by implementing a write-queueing pattern where the requests are ingested into a scalable and durable message queue such as Amazon SQS, and can be inserted into your backend database in batches using a consumer service, such as a Lambda function. This approach is more reliable than synchronous communication between multiple distributed components including the player's game client, your backend web and application servers, and your internal database system. It also reduces costs because the backend database does not need to be scaled to meet peak write throughput since the consumer processing from the write queue can be used to slow down this ingestion rate as needed.

For more information, refer to the following documentation:

- [Build highly scalable and reliable workloads using microservice architecture](#)
- [Integrating microservice by using Serverless services](#)
- [Asynchronous messaging for microservices](#)
- [Introduction to Scalable Game Development Patterns on](#)

GAMEREL_BP04 - Monitor infrastructure failures over time to measure impact on player behavior.

Ensure that your game server process and game server instance metrics are being monitored to determine the root cause of issues. In addition to monitoring CPU and memory, you can also setup monitoring for network metrics related to [network limitations of EC2 instances](#) to alert you of issues such as exceeding bandwidth, packets-per-second, or other network-level issues that may indicate your server resources are under provisioned. For game servers hosted using Amazon GameLift, consider monitoring [metrics](#) such as `GameServerInterruptions` and `InstanceInterruptions` which can help you understand how limitations in Spot instance availability are impacting your game servers deployed using Spot, and `ServerProcessAbnormalTerminations` can be used to detect abnormal terminations in your game server processes.

It is recommended to maintain historical metrics data of your game server reliability. Use this historical data for reporting purposes and join it with other datasets in order to uncover potential trends to assess the impact on player behavior over time that may be due to game server issues. Amazon CloudWatch does not retain metrics indefinitely, and the [storage resolution of metrics](#) is increased over time, so you should consider exporting these metrics to cost-effective long-term

storage such as Amazon S3. You can configure [CloudWatch Metric Streams](#) to automatically deliver your metrics from CloudWatch to your own S3 bucket where they can be stored long-term in a storage tier such as S3 Intelligent-Tiering and eventually archived using Amazon S3 Glacier. By placing your metrics in Amazon S3, they are readily available to be joined with other datasets in your data lake for interactive querying with [Amazon Athena](#).

For more information, refer to the following documentation:

- [Use Amazon EC2 instance-level network performance metrics](#)
- [CloudWatch Metric Streams](#)

GAMEREL_BP05: Adjust the number of game sessions hosted on each game server instance to reduce blast radius.

Establish a risk tolerance level for the number of players that you would be comfortable with being impacted if a game server experienced an infrastructure or software issue. Use this information to help you determine the maximum number of game sessions that you are comfortable hosting per game server instance.

When it comes to determining the amount of game sessions to host on a server instance, game developers typically start by focusing on cost, but it is important to consider the impact of this decision on reliability when designing your game architecture. While it is important to increase density of game sessions hosted per game server to improve the utilization of infrastructure, you want to ensure that you are comfortable with the blast radius of a potential failure of a single game server instance. If a single instance fails, all of the game server processes hosting active game sessions on that instance would be lost, so the number of players that experience a disruption should be tolerable based on your requirements.

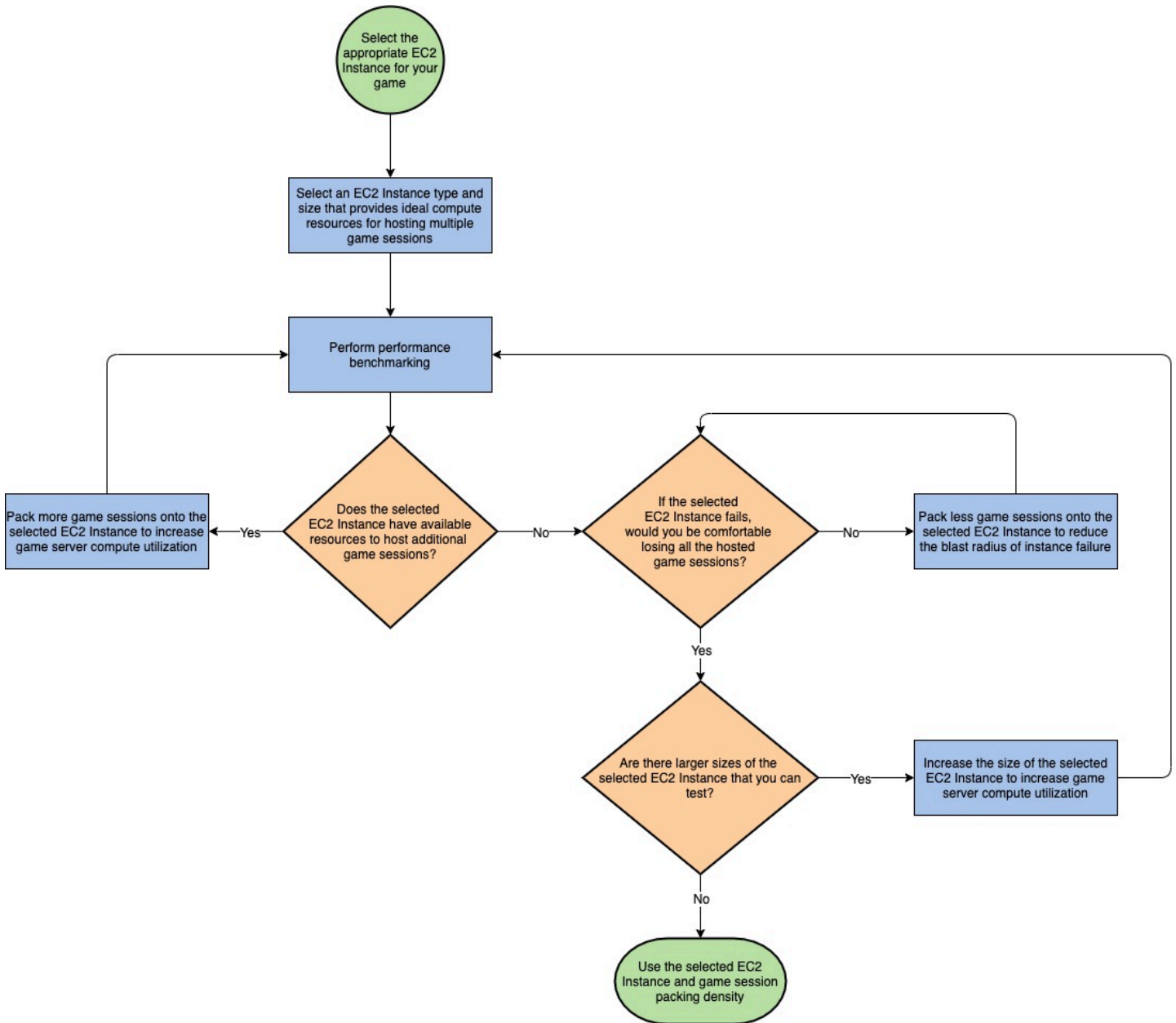
If your game requires large match sizes, such as with Battle Royale or other MMO style games, you might not have much flexibility to reduce the number of players hosted in a single game session since the style of game requires it and this requirement is more of a game design and player experience decision than an infrastructure decision. It is important to remember that typically the costs of EC2 instances scale linearly as you increase in size within a particular instance type, such as moving from a 2xlarge to 4xlarge.

Therefore, to improve reliability and reduce the impact of a single game server instance failure, you should consider increasing the number of game server instances that are hosting your game sessions. For example, as an alternative to hosting 50 game sessions on a 4xlarge EC2 instance,

you can reduce the number of players impacted by a potential instance failure if you split those 50 game sessions evenly between two 2xlarge EC2 instances with 25 game sessions hosted on each instance.

In most cases, the compute cost of these two deployment architectures is the same, but even if the costs weren't the same, it is still important to consider how this type of a change in game server hosting strategy can change the impact to players from a potential game server failure. It is also important to note that this approach assumes that the game server processes that host your game sessions each use a fixed amount of pre-allocated resources on the game server that can be evenly divided in this manner.

The following diagram provides an example process that you can use to help you to determine the number of a game sessions to host per game server instance.



Example of how to determine the number of game sessions to host per game server instance

GAMEREL_BP06: Distribute game infrastructure across multiple Availability Zones and regions to improve resiliency

To minimize the impact of localized infrastructure failures on your players, you should distribute your infrastructure deployment uniformly across enough independent locations to be able to withstand unexpected failures while still having enough capacity to meet the needs of your player demand.

When deploying your game infrastructure, it is recommended to uniformly distribute your capacity across multiple Availability Zones in a Region so that you can withstand disruptions to one or more Availability Zones without disrupting the player experience. Game backend services such as web applications should be load balanced across multiple Availability Zones or should be built using managed service such as AWS Lambda and Amazon API Gateway which provide Regional high availability by design. Similarly, components that maintain state such as caches, databases, message queues, and storage solutions should all be designed to provide durable persistence of data across multiple Availability Zones, which is provided by design in services such as Amazon S3, DynamoDB, and Amazon SQS, and can be configured in other services.

When designing your game server hosting architecture for resiliency, you should deploy your fleets of game servers uniformly across all Availability Zones within an AWS Region to maximize your access to all available compute capacity in the Region as well as reduce the blast radius of availability zone failures. For example, you can configure [Amazon EC2 Auto Scaling](#) to use all Availability Zones. If an EC2 instance becomes unhealthy, EC2 Auto Scaling can replace the instance, as well as launch instances into other Availability Zones if one or more of the Availability Zones becomes unavailable.

It is a best practice to deploy your game infrastructure into multiple Regions in order to maximize high availability. While you are encouraged to do this for your game backend services to achieve high availability, this recommendation is especially important for your game servers. In a multiplayer game for example, your infrastructure capacity for game servers is likely to outpace the capacity needs for your other services, since game servers are used to host long-lived game sessions with players. Many games choose to shard players into logical game Regions such as North America East, North America West, Europe, Asia Pacific and so on. To simplify the player experience and make it easier for you to utilize global infrastructure to host games, you should consider de-coupling the name of your player-facing game Regions, such as "North America West", from the underlying cloud provider region or data center location that is physically hosting the game servers, which might include the Oregon (us-west-2) and N. California (us-west-1) Regions, along with other infrastructure such as [Local Zones](#) or your own data centers that are all hosting game server instances supporting that player game Region.

When designing your matchmaking service, you should deploy a multi-Region architecture with separate software deployments across regions. Decouple your matchmaking service deployment from the fleets that host your game server instances so that you can route players to a game server in any Region regardless of which regional deployment of your matchmaking service handled the matchmaking request. Design logic in your matchmaking implementation to favor the game server regions that meet your latency and other rules, with the ability to fallback to routing

players to other regions if your fleets are low on capacity or there are other regional infrastructure disruptions.

For more information, refer to the following documentation:

- [Best practices for Amazon GameLift game session queues](#)
- [Amazon GameLift Multi-Region fleets](#)

Resources

Refer to the following resources to learn more about our best practices related to reliability.

Documentation and blogs

- [Practicing Continuous Integration and Continuous Delivery on AWS](#)
- [Autoscaling Asynchronous Job Queues](#)
- [Design Your Workload Service Architecture](#)
- [Timeouts, retries, and backoff with jitter](#)
- [Well-Architected Framework - Reliability Pillar](#)
- [Architecting for Reliable Scalability](#)
- [The Amazon Builder's Library](#)
- [Massive Scale Real-Time Messaging for Multiplayer Games](#)

Whitepapers

- [Introduction to Scalable Game Development Patterns on AWS](#)
- [Running Containerized Microservices on AWS](#)
- [Web Application Hosting in the Cloud](#)
- [Building a Scalable and Secure Multi-VPC Network Infrastructure](#)

Video content

- [re:Invent 2020: Ubisoft - Building a multi-platform multiplayer game on AWS](#)
- [re:Invent 2018: Supercell - Scaling Mobile Games](#)
- [re:Invent 2019: How CAPCOM builds fun games with containers, data, and ML](#)

- [re:Invent 2018: Globalizing Player Accounts at Riot Games While Maintaining Availability](#)
- [re:Invent 2020: GameLoft - A zero downtime data lake migration deep dive](#)

Training materials

- [Using Amazon GameLiftFleetIQ for Game Servers](#)
- [Game Server Hosting with Amazon EC2](#)

Performance efficiency

The performance efficiency pillar focuses on the efficient use of computing resources to meet requirements and maintaining that efficiency as demand changes and technologies evolve.

Design principles

In addition to the design principles in the AWS Well-Architected Framework whitepaper, the following design principles can help you achieve performance efficiency for your games:

Measure game performance from end-to-end including the game client, internet, and your game infrastructure: It is important to measure performance as it is perceived from the perspective of your players. This means you should measure the performance of the game client, your game infrastructure, the internet connectivity that connects your players to the infrastructure. This will help you to understand where you can make performance improvements at any level of your architecture.

Definition

There are four best practice areas for performance efficiency in the cloud:

- Selection
- Review
- Monitoring
- Tradeoffs

Take a data-driven approach to selecting a high-performance architecture. Gather data on all aspects of the architecture, from the high-level design to the selection and configuration of resource types.

By reviewing your choices on a cyclical basis, you will ensure that you are taking advantage of the continually evolving platform. Monitoring will ensure that you are aware of any deviance from expected performance and can take action on it. Finally, your architecture can make tradeoffs to improve performance, such as using compression or caching, or relaxing consistency requirements.

Best practices

The following provide best practices for your cloud architecture.

Topics

- [Selection](#)
- [Review](#)
- [Monitoring](#)
- [Tradeoffs](#)

Selection

GAMEPERF01 — How do you determine which geographic regions to host your game infrastructure?

GAMEPERF_BP01 — Review feedback from players and business stakeholders.

For an initial game launch, you should determine where to deploy infrastructure based on discussions with your business stakeholders, such as publishing teams who can help you determine where the game is expected to be made available to players, and where they are focusing their pre-launch marketing and advertising efforts.

Your business stakeholders should also have mechanisms to stimulate demand to help gain a better understanding of player reception and viability. For example, these teams will have mechanisms such as game pre-orders, marketing events and campaigns, public email lists for players to register interest before launch, and other approaches to establish relevant signals to help them determine where the game will likely have the most players at launch. The game may also use a pre-determined regional roll out strategy with which you can play test and soft-launch to further help you determine Regional player demand.

GAMEPERF_BP02 - Design an approach that supports placing latency-sensitive game infrastructure close to players to improve performance.

When first launching a game, you may not yet have enough information about your player base to adequately know where best to deploy infrastructure closest to the players that are most interested in playing your game. This is a common challenge, and you should prepare for this scenario by designing an architecture that allows you to rapidly adjust your hosting strategy to deploy servers where they are needed closer to players. It is typical for game developers to regularly assess their game infrastructure deployment as a recurring activity post-launch in order to incrementally invest in improvements over time with an iterative approach.

A best practice is to use infrastructure-as-code templates, such as AWS CloudFormation or Terraform, for the configuration of your infrastructure such as VPCs, subnet configurations, and any dependencies required to launch critical game services so that you can refer to these templates, quickly customize them if needed, and deploy them into locations where additional infrastructure is needed to support your players.

You should also make sure you understand how your current deployment strategy could be evolved to allow future expansion. For example, make sure to consider the size of the subnets you are creating for the hosting of game servers and be sure they are large enough to accommodate growth. You should also consider how game servers deployed across multiple locations will connect to your game backend, which may be hosted in a central location or in multiple locations, and may require additional configuration to support private connectivity. These considerations should be continuously evaluated over time so that you can make changes to your game hosting strategy as your game's requirements evolve over time or your player requirements change.

When determining how many game hosting locations to use for your game, you should consider the following factors:

- **Quality of player experience improvement:** How much of a player experience improvement can you introduce by adding additional game hosting locations? What is the incremental performance gain that you can achieve by doing so? How will you measure this performance improvement?
- **Which player populations to prioritize:** How many players can you improve the experience for if you add additional game hosting locations? Which player populations, or geographic locations, will you prioritize?
- **Downstream impacts of change:** If you change your game hosting strategy, how will this influence your matchmaking wait times for players? Does the match size, or number of players

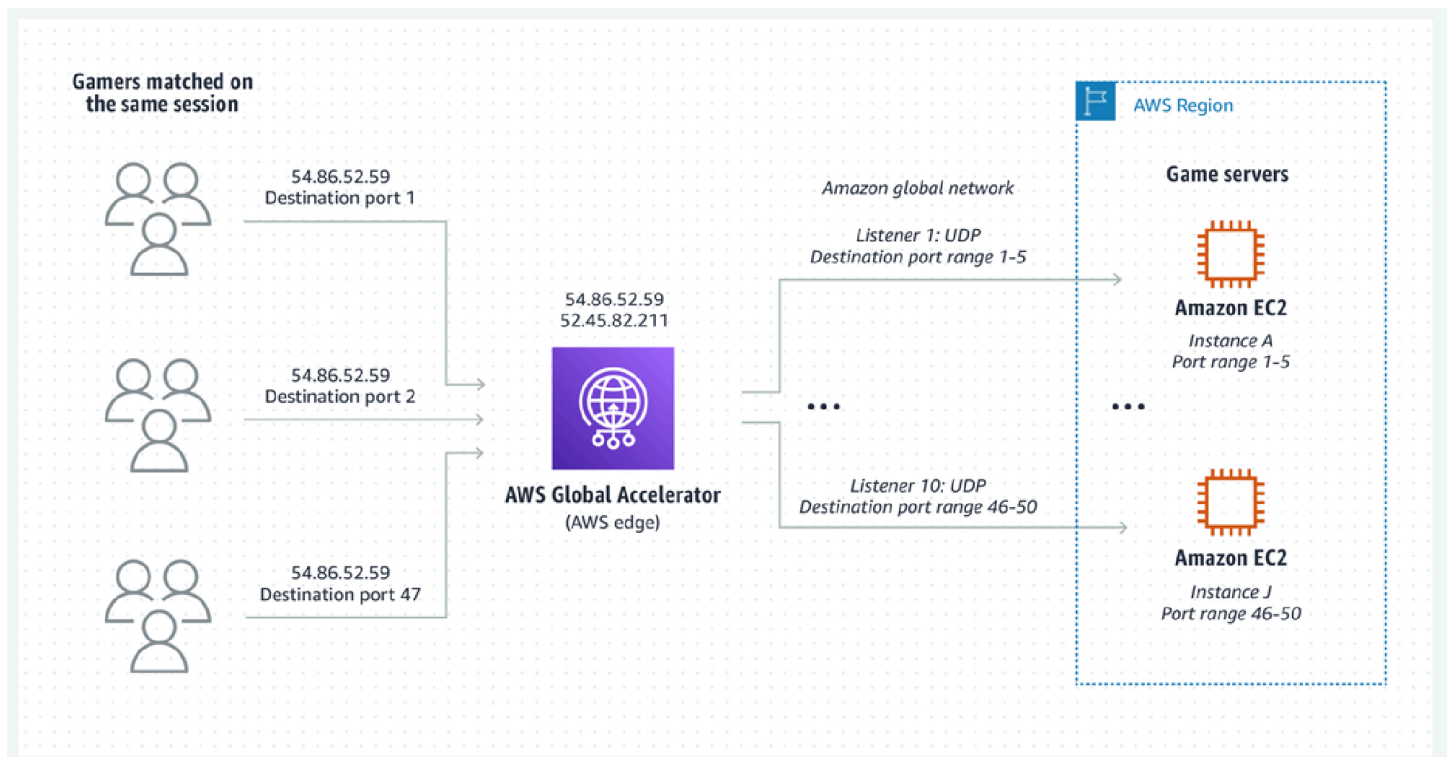
required in order to form a game session, impact your ability to build sizable player populations in different parts of the world if you introduce change?

Each of these considerations should be evaluated as you determine where you add or remove game hosting locations. For example, you may choose to prioritize improving the experience for players in geographic locations with the least performant gameplay experience, or for players who express the most vocal feedback publicly or to your community management teams. You might also choose to factor the player monetization into your priorities, for example by focusing attention on improving the experience for players in geographic locations that generate a significant source of revenue for your game, or have the potential to generate incremental revenue if you introduce performance improvements.

In addition to hosting infrastructure in AWS Regions, you can use [Local Zones](#), which are an extension of an AWS Region, to host your game servers and other latency sensitive applications such as voice chat servers closer to your players. You might also choose to run game development infrastructure in Local Zones to improve the experience for your game development teams. For example, you can use Local Zones to address use cases such as hosting replicas of your self-managed source control servers closer to your game developers, and to offer game development virtual workstations and content storage to users using Amazon EC2 instances, EBS volumes, and Amazon FSx file systems deployed into one or more Local Zones near your development studios without requiring you to host the infrastructure on-premises.

You can also extend the capabilities of into your existing on-premises data centers and co-location facilities by using [Outposts](#), which is a fully managed service that provides access to the same services and APIs using purpose-built racks and rack-mountable server options. This can help you to create a consistent deployment model across Regions, Local Zones, and Outposts deployed in your facilities. If you are building games using containers and want the flexibility to adopt a hybrid deployment architecture using open-source software that can be deployed on your own infrastructure, you can use [ECS Anywhere](#), or [EKS Anywhere](#) if you want to operate a Kubernetes-based infrastructure.

GAMEPERF_BP03 - Use network acceleration technology to improve performance across the internet.



Enhanced network performance for gaming using Global Accelerator

In addition to physically placing latency-sensitive game infrastructure closer to players, you can also improve the player experience by optimizing the network performance for your game. Use technologies that can improve your game infrastructure's connectivity to the networks, or internet service providers (ISPs), that your players are connecting to your game from. Network acceleration improves performance by optimizing the networking path that is used to route player traffic from their game client across the internet to your game infrastructure, including your game servers and game backend services. For example, [AWS Global Accelerator](#) is a networking service that improves your application's network performance using the global network, which can be used to accelerate your gameplay traffic, voice chat, and real-time messaging traffic, as well as other latency-sensitive applications. Global Accelerator [custom routing accelerators](#) can be integrated with your matchmaking service to provide deterministic routing of multiple players to the same game session IP address and port.

Your game development teams may be distributed around the world and require performant access to shared content or assets. To improve the performance for shared content stored in Amazon S3 buckets, you can setup bi-directional replication of your data across regions using [S3 Cross-Region Replication](#) so that users can access data from buckets closer to them. To simplify this access pattern, use [S3 Multi-Region Access Points](#) which accelerates requests to S3 over the global

network using Global Accelerator. For more information, refer to [Improving the Player Experience by Leveraging Global Accelerator and Amazon GameLift FleetIQ](#).

GAMEPERF02 — How do you prevent game sessions from over-utilizing resources and impacting other players running on the same game server instance?

GAMEPERF_BP04 - Monitor game server processes to detect issues.

You might run multiple game server processes per instance to efficiently utilize the resources on your game server instances. If so, you should make sure to design your architecture so that an individual game server process hosting a game session cannot cause adverse impact to other game sessions hosted on the same game server instance.

Implement monitoring of the limited resources available to your game server instances so that you can generate alerts when individual game server processes are breaching pre-determined resource budget thresholds. When thresholds are breached, you may want to configure your game server software to dump relevant system and game server logs out to durable storage, such as a central logging solution, so that your game server engineers can investigate this behavior. Additionally, your game server instance should be configured to report metrics from each of the game server processes running on the instance so that you can monitor these individual game server processes in addition to the overall metrics for the game server instance. For example, Amazon GameLift provides metrics for [monitoring game sessions](#), which can be augmented with custom game-specific metrics and logs collected using the [Amazon CloudWatch Agent](#) which you can configure on your game server instance. Your metrics can be viewed in CloudWatch or exported to other tools such as [Amazon Managed Grafana](#) which is integrated with Single Sign-On to make it easier to access metrics by users who may not have access to the Management Console. Refer to the following best practices for [managing logs and metrics using Amazon GameLift](#), which also provides support for viewing individual [game session logs](#).

GAMEPERF_BP05 - Performance test your game server with simulated and real gameplay scenarios.

You should conduct performance testing and evaluate various gameplay scenarios to determine whether the game server process handles the utilization of fixed resources appropriately, such as EC2 instance memory, CPU, and network bandwidth.

You should create simulated gameplay tests with bots that can mirror common gameplay paths and behaviors of your players so that you can determine how your game server processes handle this under different usage scenarios. For example, you can implement a solution, such as [Distributed Load Testing on AWS](#) that you can customize to run game client simulations or game client builds to generate gameplay scenarios. Run internal play tests and use QA teams to stress test the various features of your game so that you can develop confidence that your game is designed to perform optimally. [AWS Device Farm](#) can be used to perform mobile and web testing for your iOS, Android, and browser games.

GAMEPERF03 — How do you select the appropriate compute solution for your game?

GAMEPERF_BP06 - Benchmark your game performance across multiple compute types.

For game server workloads, there is no one-size-fits-all approach to identifying the optimal compute solution for hosting your game server. A common choice for game server is to use the compute-optimized EC2 instances, because this instance family is optimized for workloads like game servers that are computationally intensive. Alternatively, if your game requires a significant amount of memory in order to implement specific features, the memory-optimized instances may be most suitable.

For use cases where your workload utilizes significant network resources, consider implementing instances that are network-optimized which is typically indicated by the use of an "n" in the instance name. Games are sensitive to latency and dropped packets, so it is recommended to use EC2 enhanced networking to improve the network performance of your game servers. Enhanced networking uses single root I/O virtualization (SR-IOV) to provide high-performance networking capabilities on [supported instance types](#). SR-IOV is a method of device virtualization that provides higher I/O performance and lower CPU utilization when compared to traditional virtualized network interfaces. Enhanced networking provides higher bandwidth, higher packet per second (PPS) performance, and consistently lower inter-instance latencies. Enhanced networking with Elastic Network Adapter is available for most recent EC2 instance types.

If your game performs similarly across multiple EC2 instance types, then you should consider using multiple instance types to host your game servers so that you can monitor performance over time and perform further optimization after you have hosted enough production game sessions to be able to identify performance trends over time. Remember that your resource requirements may change over time as you add new features into your game that require different allocation of

resources. You can [configure EC2 Auto Scaling groups](#) to use multiple instance types, or you can use separate Auto Scaling groups to host game server instances that run separate instance types which may make it easier to manage correlation and aggregation of metrics.

You should also evaluate how your game performs on different types of processors such as Intel-based instances, AMD-based instances, and ARM-based Graviton instances.

You should also benchmark how your game performance is impacted when it is hosted using containers and Lambda functions. For use cases where long-lived game server processes are not required, such as asynchronous games and for game backend services, you should consider using a serverless architecture with Lambda which can simplify management and operations for game operations teams, as well as allow you to more quickly deploy your game globally to many AWS Regions. For serverless best practices, refer to the [Serverless Applications Lens - Well-Architected Framework](#). For more information, refer to [Choose the right compute strategy for your global game servers](#).

GAMEPERF_BP07 - Use graphics instances for game development virtual workstations.

Game designers, engineers, artists, QA, and other personnel may require the use of virtual workstations. To support these use cases, use graphics-optimized instances, which are denoted with an instance name starting with "g", as are built using GPUs that are purpose-built to support graphics use cases such as game development and game streaming.

Evaluate and benchmark performance across different graphics-optimized instance types using the same tools that your end users typically require. For example, offers a variety of graphics-optimized instances that come with different GPUs from manufacturers such as AMD and NVIDIA. When benchmarking these instances, make sure that your software is compatible with the GPU and associated drivers that are supported. For graphical artist use cases, consider [Amazon Nimble Studio](#) which provides access to cloud-based virtual workstations, file storage and tools needed for operating a cloud-based studio.

If you are developing your own custom virtual workstation using EC2, you must consider how end users will access these virtual workstations. While many options exist for connectivity, consider using [NICE DCV](#), a high-performance remote display protocol for connecting to remote desktops from a locally installed client software or web browser, which is no cost and is available in the [AWS Marketplace](#) as an Amazon Machine Image (AMI) with NICE DCV server pre-installed.

For more information, refer to the following:

- [Virtual workstations on powered by NVIDIA Quadro technology](#)

- [Game Production in the Cloud - Workstations: Stream Game Development from anywhere with NICE DCV](#)
- [Stream a remote environment with NICE DCV over QUIC UDP for a 4K monitor at 60 FPS](#)
- [Putting bitrates into perspective](#)

GAMEPERF_BP08 - Push non-latency-sensitive compute tasks to asynchronous workflows.

When you are optimizing the performance for your game, it is important to keep in mind that not all interactions between the client and the game backend must be performed in a synchronous manner. You should consider each feature from the perspective of the player experience and determine whether or not certain interactions require synchronous communications, which are blocking and resource intensive, or whether those features can be implemented in an asynchronous manner. Ensure that when you implement network calls, you are doing so using an asynchronous non-blocking approach. Additionally, your game backend should also be configured to perform work in an efficient manner by offloading tasks to queues and prioritizing fast responses to clients where possible.

For example, updating a leaderboard at the end of a player session can be implemented asynchronously so that the client does not need to wait for the leaderboard update to complete. Instead, implement this asynchronously on the game client, and also consider designing your backend service to push these types of operations into queues, such as Amazon SQS. With this architecture, your backend should be configured to accept the request, enqueue it in SQS which durably stores messages for asynchronous processing, and promptly reply back to the client. When the leaderboard update is completed, the backend can send an update to the game client so that the player's view of the leaderboard is updated. Alternatively, the player can simply visit your game's leaderboard screen to retrieve the latest data, which can issue a web request to your backend to retrieve the latest data from cache.

For more information, refer to the following documentation:

- [Understanding Asynchronous Messaging for Microservices](#)
- [Lambda - Using service integrations and asynchronous processing](#)

Review

Refer to the Well-Architected Framework whitepaper for best practices in [performance efficiency](#) that apply to games workloads.

Monitoring

GAMEPERF04 — How do you design your matchmaking service to optimize performance?

GAMEPERF_BP09 - Define network latency thresholds for your game.

When developing a multiplayer game, ensure that your game infrastructure does not introduce unnecessary latency for players. If your game is sensitive to network latency, then you should set latency thresholds in your matchmaking logic to prioritize placing players on game server instances that are hosted in Regions where their connection to the game session meets your objective for ideal player experience.

In many latency-sensitive games it is common to instrument the game clients to ping each of the game's infrastructure Regions to gather performance data such as network latency, jitter, and packet loss, and report this data to your metrics collection backend so that it can be analyzed. When matching players into game sessions, you can configure your game to incorporate the game client's perceived network latency to your game server infrastructure as one of the inputs used in your matchmaking service when selecting a game for a player.

GAMEPERF_BP10 - Run a separate matchmaking service for each gameplay mode and game hosting Region.

If your game offers multiple gameplay modes for players to choose from, you should separate the matchmaking systems for each of them so that you can independently tune the performance for each gameplay mode based on its unique requirements, and reduce resource contention. Each gameplay mode will likely have unique requirements for acceptable latency, match size, as well as other customize game-specific matchmaking logic. They will also likely attract different types of players. Run each game mode's matchmaking service as a separate software deployment so that you can more easily performance test and operate the game modes independently. For example, you might run these as separate Lambda functions for each game mode, or you might operate them as separate container-based service deployments.

Deploy your matchmaking services to multiple Regions, preferably the same Regions that you host your game servers, so that players can be routed to a matchmaking service that is closest to them which can improve the efficiency of finding game servers with low latency. Amazon GameLift FlexMatch provides additional guidance for selecting Regions for matchmakers, and includes the ability to integrate your matchmakers with [multi-Region game session queues](#).

GAMEPERF_BP11 - Regularly monitor matchmaking performance.

One of the most effective ways to optimize the performance of a game for players is to reduce the time that they must wait before they can enter into a game session. Long wait times can cause players to lose interest and lead to attrition, so it is important to consider this when designing your matchmaking solution.

When you are designing your game matchmaking configuration for your game, you will need to create rules that determine the conditions that are applied in order to form a match. You should consider the impact that these rules will have on the performance of the system, particularly the wait times for players. Before deploying changes to your matchmaking implementation, such as the addition of new matchmaking conditions or filters, you should properly test this beforehand, or consider releasing this change gradually to a small sample population of players as a canary or A/B test to gather performance metrics before introducing this change to the entire player population.

Configure your matchmaking service to generate detailed logs that can help you to understand what conditions or rules were applied to each matchmaking request so that you review and adjust matchmaking implement as necessary. For example, Amazon [Amazon GameLift FlexMatch](#) provides a fully-managed matchmaking service which can be used as a standalone service with your own game server hosting or used with game servers hosted on Amazon GameLift. FlexMatch can generate event notifications to Amazon EventBridge (formerly CloudWatch Events) and Amazon Simple Notification Service (Amazon SNS) in JSON format so that you can automatically process and store this data for analysis to help you improve matchmaking performance.

Setup metrics to track how long your matchmaking service takes to find a suitable game session for players. Regularly review matchmaking duration metrics and correlate these times with player behavior and community sentiment in order to develop suitable thresholds for matchmaking timeouts that you can include in your matchmaking rule configuration. For example, Amazon GameLift FlexMatch provides support for defining matchmaking request timeouts as well as creating matchmaking rules that can [allow requirements to relax over time](#). This feature allows you to create matchmaking that can adapt to make it easier to create matches and place players into game sessions when matches are difficult to find.

Tradeoffs

Refer to the Well-Architected Framework whitepaper for best practices in [tradeoffs](#) for performance efficiency that apply to games workloads.

Resources

Refer to the following resources to learn more about our best practices related to performance efficiency.

Documents and blogs

- [Games Industry Architecture Center](#)
- [Well-Architected framework Performance Efficiency Pillar](#)
- [Comparing your on-premises storage patterns with Storage services - Storage Blog](#)
- [Amazon EC2 instance store - Amazon Elastic Compute Cloud](#)
- [Collecting Metrics and Logs from Amazon EC2 Instances and On-Premises Servers with the CloudWatch Agent](#)
- [Allow and configure enhanced networking for EC2 instances](#)
- [Global Accelerator and Amazon GameLift FleetIQ](#)
- [Riot Games Technology Blog: Scalability and Load Testing For Valorant](#)
- [Hyper-scale online games with a hybrid Solution](#)

Whitepapers

- [Optimizing multiplayer game server performance on AWS whitepaper](#)
- [Performance at Scale with Amazon ElastiCache](#)
- [Database Caching Strategies Using Redis](#)
- [Amazon Virtual Private Cloud Connectivity Options](#)
- [Best Practice Design Patterns: Optimizing Amazon S3 Performance](#)

Third-party tools

- [Unreal Engine Performance and Profiling](#)
- [Unity Profiler](#)
- [Open 3D Engine \(O3DE\) Profiler](#)
- [Monitoring Amazon GameLift](#)

Video content

- [Amazon EC2 foundations \(CMP211-R2\)](#)
- [Powering next-gen Amazon EC2: Deep dive into the Nitro system](#)
- [Getting Started with Amazon GameLift FleetIQ](#)
- [Riot Games: Outposts Customer Testimonial](#)

Training materials

- [Game Server Hosting with Amazon EC2](#)
- [Using Amazon Amazon GameLift FleetIQ for Game Servers](#)
- [Getting Started with Game Tech](#)
- [Game Server Hosting on AWS](#)
- [Amazon GameLift Primer](#)

Cost optimization

The cost optimization pillar includes the continual process of refinement and improvement of a system over its entire lifecycle. From the initial design of your first proof of concept to the ongoing operation of production workloads, adopting the practices in this paper will allow you to build and operate cost-aware systems that achieve business outcomes and minimize costs, thus allowing your business to maximize its return on investment.

Design principles

In addition to the design principles from the cost optimization pillar of the Well-Architected Framework whitepaper, the following design principles can help you optimize the costs of running your game workload in the cloud.

Measure the infrastructure cost per player, platform, and game feature: The costs to develop and operate the game, and acquire and retain players have significant impact on the financial success of a game. Therefore, it is important to be able to understand and track the infrastructure costs that are required for specific player experiences and features across game platforms so that you can identify the resources of your architecture that may require cost optimization.

Assess the tradeoff of cost optimization versus player experience: Determine whether or not you should focus more on player experience features and improvements, or cost optimization. Typically, after a game has reached a critical mass and player population stabilizes it is time to focus on optimizing the cost for operating the game.

Definition

There are four best practice areas for cost optimization in the cloud:

- Cost-effective resources
- Matching supply and demand
- Expenditure awareness
- Optimizing over time

Prior to launch, game developers typically do not have a clear understanding of how popular, successful, or long-lasting their game will become after launch. Some games are highly anticipated but are unable to retain players for a long period of time, while other games may immediately or gradually increase their player base and develop into sustainable and profitable businesses. Depending on the game's monetization strategy, business priorities, and where the game is in its lifecycle, a game developer will need to make trade-offs to evaluate cost optimization decisions. For example, a game developer may be preparing to release a highly anticipated new game that has media publicity and widespread industry awareness.

During this pre-launch phase, a developer will likely be prioritizing their focus on time-to-market, feature development, and game performance rather than cost optimization. When a new game launches, game developers want to be sure that their infrastructure can scale to meet their peak player demand. This typically leads to the over-provisioning of resources to account for peak player forecast projections, or best-case sales scenarios. Alternatively, if a game is not successful, or if it's in a stage where development on the game is beginning to slow down, a developer may want to prioritize reducing costs as much as possible so that they can continue to operate the game for existing players for years to come.

A game developer may operate multiple games simultaneously, which requires additional considerations to be factored in. For example, a game developer may reuse resources including technical infrastructure, software, and staff for multiple live games, where the operating costs are shared between the games.

Games are unique workloads due to their business models, scale, and unpredictability. The following questions can help you make cost-optimization decisions for your game.

Best practices

The following provide best practices for your cloud architecture.

Topics

- [Cost-effective resources](#)
- [Matching supply and demand](#)
- [Expenditure awareness](#)
- [Optimizing over time](#)

Cost-effective resources

GAMECOST01 - How are you choosing the right compute solution for your game servers?

One of most unique aspects of a game workload compared to other types of workloads is the game server, which is critical to the player experience. Because players connect to game servers from their game client to play a game session, it is also one of the biggest drivers of cost for operating a multiplayer game so it is important to make sure that you optimize how you utilize the compute infrastructure for your game to reduce costs.

GAMECOST_BP01: Benchmark your game server on multiple compute types.

During the initial planning and testing phase of game development, you should perform benchmarking to determine the appropriate type of compute to use for your game. Typically, session-based multiplayer and other types of low-latency games will use Amazon EC2 Instances for hosting game servers. Each EC2 instance type provides a mixture of compute resources that are optimized for different workload profiles. You should perform benchmarking of your game server code to determine what resources such as CPU, Memory or Network bandwidth that your game session utilizes and select the option that provides the right balance of performance at the lowest cost. Most of the popular commercially available game engines such as Unreal Engine, Unity, and Lumberyard provide performance profiling utilities that you can allow in the engine editor so that your game server builds will emit log and/or metric data to help you benchmark performance and

resource utilization. This telemetry can help you evaluate and select the appropriate EC2 Instance types to use.

As part of benchmarking your game server across multiple EC2 instance types, you should determine what type of operating system and processor requirements are needed to run your game. For best cost optimization, it is recommended to run your game compute infrastructure on Linux instances to eliminate the licensing costs that is incurred with Windows. Additionally, [Graviton instances](#) are 64-bit Arm-based EC2 instances that can be used to run your game servers, including [Unreal Engine dedicated servers](#).

GAMECOST_BP02. Optimize the number of game sessions hosted on each game server Instance to reduce costs.

Optimize the number of game sessions hosted per server instance in order to achieve better compute utilization and reduce compute infrastructure costs.

To reduce costs, game developers should maximize the number of game sessions hosted on the same physical or virtual server, also known as the packing density of their game servers. This is achieved by increasing the number of game server processes that can be simultaneously hosted on an EC2 instance. A single game server process should not usually require the use of the entire resources available on the EC2 instance. This is one of the most important ways to reduce compute costs for a game and requires the use of software that can spawn and manage multiple server processes on the EC2 instance on separate ports. For example, Amazon GameLift has a [quota on the maximum number of game server processes per instance](#), which you should strive to utilize so that you can reduce hosting costs. Refer to the documentation for Amazon GameLift for details on the current quota for maximum game server processes per instance.

As an alternative to deploying game server processes on virtual machines such as EC2 instances, it is becoming popular for game developers to run their game servers as container-based applications using container orchestration platforms such as Amazon Elastic Container Service (Amazon ECS) or Amazon Elastic Kubernetes Service (Amazon EKS), or by [hosting the game server using Fargate](#). Container platforms provide job scheduling functionality that can automatically find an available container instance in the cluster to host your game server container based on resource requirements and other placement logic that you specify. However, as discussed in the reliability pillar of this lens, it is important to consider how you will manage the scaling and player placement behavior in way that doesn't disrupt active player sessions.

GAMECOST_BP03. Select the appropriate compute pricing option to reduce costs.

Run performance tests of your game server software across a variety of instance types and compute options to determine which option is most cost-effective for your game.

In addition to efficiently utilizing the right EC2 instance types for your workload, consider which of the available compute pricing options is most suitable for your cost optimization goals. There are several pricing options available, including On-Demand Instances, Spot Instances, Reserved Instances, and Savings Plans.

Spot Instances are ideal for running game servers because they offer the largest compute discounts, do not require usage commitments, and they provide flexibility for unpredictable and spiky workload types. However, Spot Instances can be interrupted, so they are best suited for game server workloads with short game session durations or situations where the tolerance for interruption is higher. For example, the [Running Your Game Servers at Scale for up to 90% lower compute cost blog post](#) provides guidance for running game servers using Kubernetes on Amazon EKS with EC2 Spot Instances. When using Spot, it is also recommended to run game server workloads across multiple EC2 instance types and Availability Zones in an AWS Region to diversify your usage of capacity and reduce interruption risk. It is also recommended to use Spot Instances in combination with On-Demand Instances to minimize the impact of potential disruptions to active game sessions, and to consider using capacity optimized allocations strategy to further reduce the risk of interruption. Refer to the [Best Practices for EC2 Spot](#) for additional best practices. [Amazon EC2 Auto Scaling Capacity Rebalancing](#) can be used to proactively monitor and add additional capacity when Spot Instances are at increased risk of interruption. [Amazon GameLift FleetIQ](#) integrates with Spot Instances to optimize the use of low-cost Spot Instances while reducing the risk of interruptions. If you are hosting your game using Amazon GameLift, you should review the [Amazon GameLift documentation](#) for choosing computing resources.

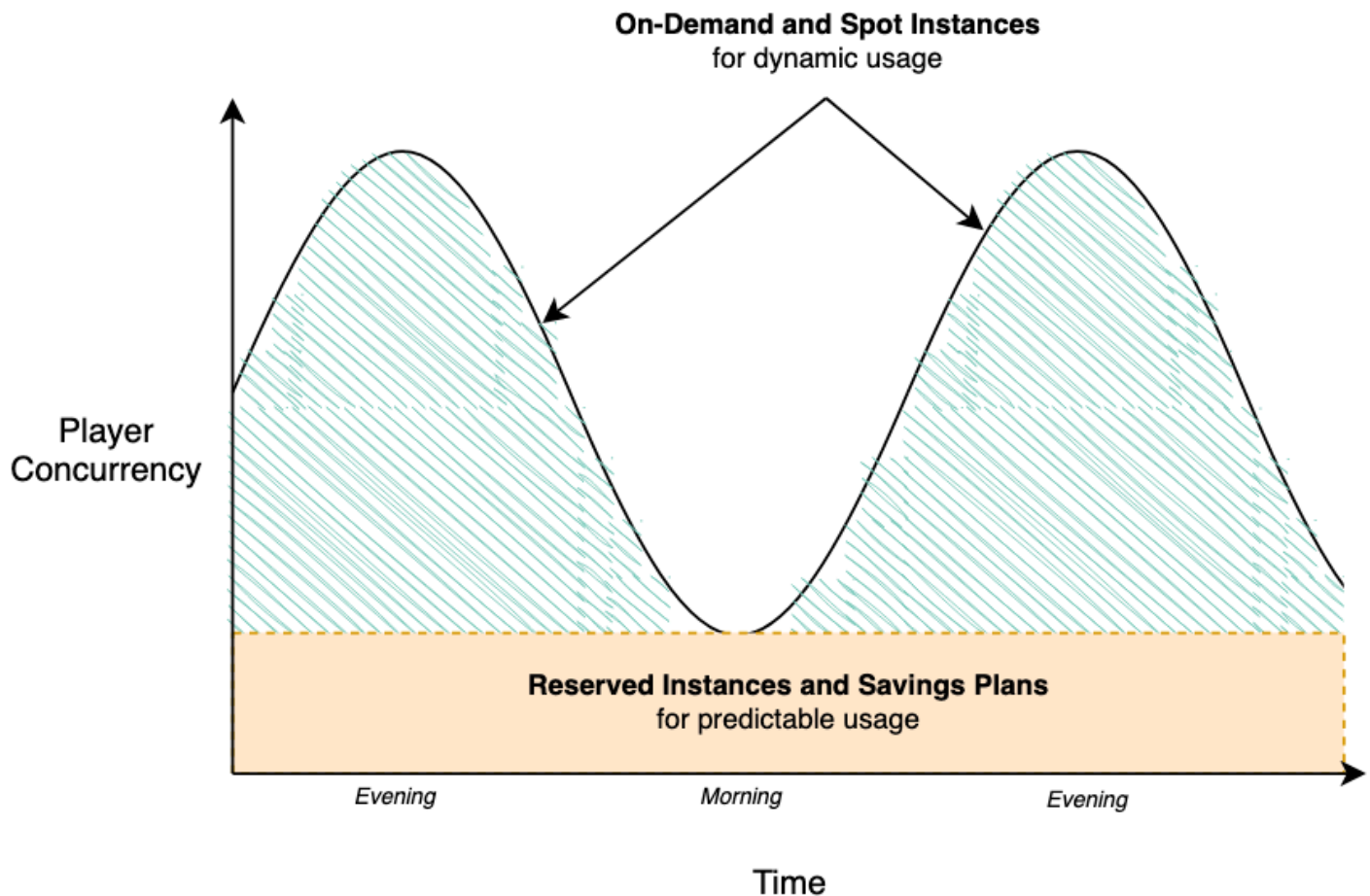
[EC2 Reserved Instances](#) allow you to receive a discount for compute by making usage commitment to a particular Region and instance type, and as an alternative to Reserved Instances (RIs), [Savings Plans](#) provide discounts similar to RIs with flexibility to apply the discounts across Regions, instance family, operating system, tenancy, and can be applied to other compute services such as Fargate and Lambda. Because Savings Plans provide regional flexibility, they are particularly ideal in situations where your game has unpredictable usage across geographies such as with new game launches. This provides a significant discount compared to On-Demand pricing and is ideal for scenarios when you can forecast your expected usage for a 1-year or 3-year period.

The flexibility to apply the discount across different compute services can be a useful benefit to allow you to apply your commitment-based usage discount across the infrastructure for your game servers running on EC2 instances, and your game backend services which may be operating on

other services such as Lambda. Unlike Spot Instances which can be interrupted, Savings Plans and Reserved Instances are simply a billing benefit and provide access the same usage characteristics as On-Demand capacity. Typically, in game server workloads, Reserved Instances are introduced after a game has been running in production for an extended period of time, at least several weeks or months, where daily usage patterns are well understood. Since Reserved Instances and Savings Plans require a usage commitment, it is recommended to maximize the utilization of pre-purchased Reserved Instances and Savings Plans. They can be augmented with other purchase options that provide more flexibility for unpredictable game server usage spikes, such as On-Demand and Spot Instances.

For example, if your daily player usage pattern always requires at least 20 servers to support your player base, but periodically requires up to 40 servers, then you should consider purchasing 20 Reserved Instances or an equivalent Savings Plan commitment, because that usage demand is predictable and consistent, and will result in maximum utilization of the usage commitment that you have purchased. The additional capacity that is required to support your players can be hosted using Spot and On-Demand Instances.

The following diagram provides an example to illustrate the use of multiple compute pricing options for game server workloads.



Hosting game servers with multiple EC2 pricing options

In the diagram, the player concurrency fluctuates over time which makes it difficult to manage utilization and achieve cost optimization. To address this fluctuation, consider adopting a mixture of different compute pricing options, using Reserved Instances and EC2 Savings Plans to meet the needs of your minimum usage requirements while relying on EC2 On-Demand and EC2 Spot Instances for dynamic usage.

GAMECOST02 - How are you optimizing the data transfer costs for your game infrastructure?

Games can transfer a significant amount of data across the internet between your players' game client devices and your game infrastructure to provide the gameplay experience, as well as between the components of your game infrastructure. For example, data transfer occurs when players download game content updates to their game clients, save their game progress state to the cloud, engage in real-time multiplayer game sessions with their friends, and when your

game infrastructure transfers data between Regions and Availability Zones. It is important to understand where the data transfer occurs in your game workload so that you can optimize your architecture choices to reduce this data transfer cost. To optimize the data transfer costs for your game, consider the following best practices.

GAMECOST_BP04: Optimize the cost of data transfer across the internet.

Implement solutions that reduce the cost of transferring data from your game backend to your players.

Use CloudFront to reduce the cost of content delivery and heavily used public-facing web applications. Game content and assets that are stored in the cloud are typically stored in Amazon S3 and delivered to the game client either directly from S3 or from web servers hosted in Amazon EC2 that retrieve the content from Amazon S3 and deliver it to clients. To reduce the data transfer costs of content downloads, consider using Amazon CloudFront in front of your cloud storage to deliver content to users. Using CloudFront can reduce the cost of data transfer because it costs less to deliver your content from CloudFront points-of-presence than directly from Regions, and CloudFront does not charge origin retrieval fees for AWS-based origins, such as Amazon EC2 and Amazon S3. If your content is cacheable, CloudFront can be used to cache content closer to users which can further reduce costs. CloudFront is also beneficial for placement in front of public-facing web applications and services, even if caching is not used, since the cost of data transfer between your servers and clients can be reduced by routing traffic through the CloudWatch network.

[CloudWatch](#) can be used to monitor your Amazon CloudFront usage. For use cases where you use multiple content delivery networks (CDN), [CloudFront Origin Shield](#) can provide an additional layer of caching to consolidate and reduce the number of origin requests from different providers. For more best practices for content delivery, refer to the [Content Delivery for Games whitepaper](#).

[VPC Flow Logs](#) can be used to monitor the network traffic in your environment and help you to identify the sources and destinations of traffic to help you optimize your data transfer costs.

GAMECOST_BP05: Optimize costs to reduce data transfer between services, Availability Zones, and Regions.

In addition to optimizing the data transfer between your game infrastructure and the internet, you should also optimize the data transfer between the internal components of your game infrastructure to reduce the amount of traffic sent between Availability Zones in the same Region, and between Regions, which each incur data transfer costs.

Prioritize keeping internal traffic in the same Availability Zone as the application. To optimize data transfer in your game backend services, you can deploy your database and cache clusters with

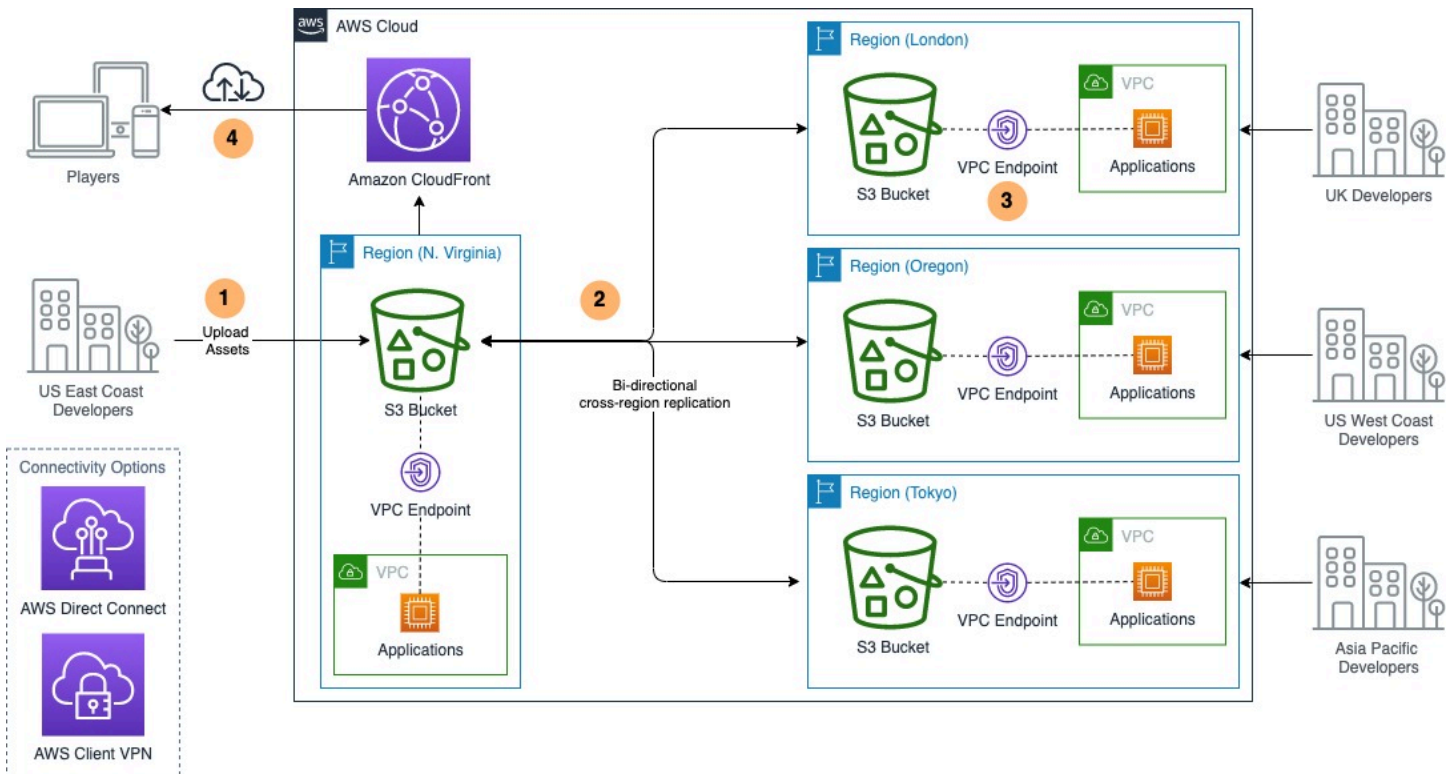
instances into multiple Availability Zones in a Region and configure your applications to prioritize reading data from instances that are in the same Availability Zone as the application server. Although this setup still incurs data transfer costs for the data replication between Availability Zones, this is recommended in use cases where applications heavily utilize databases and caches, such as read-heavy workloads that can achieve cost benefits from having local copies of the data in the same availability zone.

You should replicate copies of data to other Regions if there are applications in those regions that require regular access to the data. It is more cost-effective to replicate the data across Regions so that applications can access a local copy of data as much as frequently as needed, rather than relying on those applications to access data across regions which is not cost-effective at scale, less performant, and requires more complex cross-Region networking configurations in order to provide appropriate security controls.

For example, your game backend services might be deployed in the N. Virginia Region with game servers deployed globally into multiple Regions closest to your players to reduce gameplay latency. If your game servers need to access objects that are stored in an S3 bucket or cache data in Amazon ElastiCache for Redis that is hosted in N. Virginia, it is more cost effective to replicate the cache data to the Regions where the game servers are located to reduce the ongoing data transfer cost for those servers to retrieve the data. AWS offers features that make it easier to set up multi-Region replication of data, such as [Amazon Aurora global databases](#), [Amazon ElastiCache Global Datastore for Redis](#), and [Amazon DynamoDB Global Tables](#). For use cases where objects stored in Amazon S3 needs to be frequently accessed by applications that are hosted in another Region, consider using [Amazon S3 Cross-Region Replication \(CRR\)](#) to reduce cost. CRR can reduce costs by automatically replicating copies of objects to destination buckets hosted in one or more Regions where your applications are deployed. This configuration would still incur the cost of replicating the object to another Region, but it would eliminate the data transfer costs that would otherwise be incurred each time the cross-region application retrieves the object from S3, since it would retrieve it from an S3 destination bucket in the same Region.

It is recommended to use VPC endpoints to integrate with services to reduce data traffic and processing charges through NAT Gateways. Similarly, for public facing applications hosted in Public Subnets, traffic may not need to traverse a NAT Gateway and can be configured to send outbound traffic directly to an internet gateway to avoid the data processing and transfer costs of the NAT Gateway where it isn't needed.

The following diagram illustrates an architecture that can be used to reduce the cost of accessing data from applications hosted in other Regions that require low latency access to shared datasets.



Optimizing costs for accessing latency-sensitive game content from global users

1. Your game development teams may be globally distributed and require access to copies of the same content in Amazon S3. In this scenario, a game developer located in US East Coast can upload content to an Amazon S3 bucket either directly or from an application they are hosting in that Region.
2. **S3 Cross-Region Replication** is configured to replicate copies of objects to buckets in other Regions so that applications hosted in those Regions can retrieve objects from the local Region without needing to send requests across regions to access them. Replication can be configured to be bi-directional so that updates made in any of the other Regions can be updated in the rest of the Regions.
3. **VPC Endpoints** provides private access to Amazon S3 from your VPC so that applications do not need to route traffic through a NAT Gateway, which could be used by other high throughput application traffic and can cause congestion. Game development teams such as other global studios, remote and contract workers can access copies of datasets by connecting to the region that is most performant for them. Use **Direct Connect** to set up a dedicated connection between your studio locations or data centers and Regions. Use **Client VPN** to provide remote workers with secure remote access to your VPCs.

4. Player game clients and other internet-based applications integrate with CloudFront, which provides content caching for objects stored in S3 and reduces the cost of data transfer for static and dynamic content over the internet.

[Multi-Region Access Points in Amazon S3](#) can be used to simplify this access pattern for applications hosted in Regions where you do not host S3 buckets. Applications can interact with a multi-Region access point which can determine the lowest-latency bucket location to serve their request. Multi-Region Access Points have an additional cost.

GAMECOST03 - How are you optimizing the data storage costs for your game infrastructure?

Games can generate large amounts of data that needs to be stored and made available to developers, players, and to the game itself. For example, you may be constantly generating new source code, game content, and assets that need to be stored, your players may be generating new user generated content, and your game clients and servers may be generating game analytics telemetry data that needs to be stored in a data lake and made available to analytics teams. Your game also generates structured data.

GAMECOST_BP06. Choose the appropriate type of storage to reduce costs.

Each type of data that you generate and store has unique characteristics that you should consider when determining the right storage solution to use for your workload.

Use **S3 Object Lifecycle Management** to store object data in the most cost-effective storage class. Amazon S3 provides multiple [storage classes](#) and [object lifecycle management](#) to make it easy to setup simple and fine-grained policies to automatically transition data between storage tiers to reduce costs. Instead of simply storing all data in S3 standard storage class by default, consider setting up a lifecycle configuration to transition data between tiers automatically over time, or use S3 Intelligent-Tiering storage class for unknown or changing access patterns. Alternatively, S3 Intelligent-Tiering can cost-effectively and automatically transition data between tiers and is recommended as a default storage class since it provides cost optimization without the need to manually setup lifecycle policies, and is now the [best choice for small and short-lived objects](#). Common use cases for Amazon S3 include storage of game assets, static content, game logs, data lake storage, and backups. For use cases where file systems are required, such for attaching shared file systems to workstations during development, consider using [Amazon Elastic File System](#)

[\(Amazon EFS\), which provides different storage classes](#) and automatically grows and shrinks as you add and remove files with no need for manage the infrastructure.

Matching supply and demand

Refer to the Well-Architected Framework whitepaper for best practices [managing supply and demand](#) for cost optimization ply to games workloads. Additionally, for more best practices related to dynamic scaling for game infrastructure, refer to the reliability and operational excellence pillars of this lens.

Expenditure awareness

GAMECOST04 - How do you measure the cost of your game environments?

Understand the cost per player, game feature, and environment so that you can manage and forecast your spend as the number of players changes over time and features are added and improved. Consider the following best practices in order to manage your costs of your different game environments.

GAMECOST_BP07: Implement attribution of cost per player, game feature, and environment.

Cost attribution for game servers is usually easier to perform than game backend services because a game server is usually optimized to be able to host a specific number of concurrent players per instance which can be amortized across the cost of running the instance. For game backend services, it is recommended to de-couple the components of your game into distinct features that can be managed as separate logical or physical resources to make it easier to analyze costs. For example, although it may seem easier to implement a single monolithic application to host game backend services, this pattern makes it hard to derive the total cost per player and game feature over time as you add more features because the compute, networking, and storage costs of resources are shared across all of them.

Consider adopting a serverless architecture for your game backend services with services such as Amazon API Gateway and Lambda or Fargate for compute, Amazon SQS and Amazon SNS for messaging, Amazon S3 for object storage, and Amazon DynamoDB for database storage. These services are just a few examples of products that offer pricing that is usage-based and primarily driven by request volume so that costs can be broken down more granularly. Individual resources such as Lambda functions, Fargate services, DynamoDB tables, and S3 buckets can all be associated

with cost allocation tags so that you can attribute the costs of these services with game feature names that make it easier for you to understand the costs for each of your services.

It is also recommended to separately manage each of your game development environments so that you can easily attribute costs for the different environments. Typically, game developers will manage separate environments for development, test, staging and production environments, as described in the operations pillar of this games industry lens. Each environment usually has different scalability, performance, and usage requirements and may be managed by separate teams. To control costs, organize these environments so that you can properly monitor and attribute the costs of each environment.

For more information, refer to the following documentation:

- [Building a serverless multi-player game that scales](#)
- [Standalone Game Session Servers with a Websockets-based backend](#)
- [Standalone Game Session Servers with a Serverless backend](#)

Optimizing over time

Refer to the [Well-Architected Framework Cost Optimization Pillar](#) for guidance on optimizing over time.

Resources

Refer to the following resources to learn more about best practices for cost optimization.

Documentation and blogs

- [How can I reduce data transfer charges for my NAT gateway?](#)
- [Amazon GameLift FleetIQ adapter for Agones](#)
- [How can I find the top contributors to traffic through the NAT gateway in my VPC?](#)
- [Choose the right compute strategy for your global game servers](#)
- [Well-Architected Labs -- Cost effective resources](#)
- [Amazon VPC CNI plugin increases pods per node limits](#)
- [Architecture Best Practices for Cost Optimization](#)
- [Reducing player wait time and right sizing compute allocation using Amazon Sagemaker RL and Amazon EKS](#)

- [Compute Optimizer](#)
- [Electronic Arts optimizes storage costs and operations using Amazon S3 Intelligent-Tiering and S3 Glacier](#)
- [Escape unfriendly licensing practices by migrating Windows workloads to Linux](#)
- [Overview of Data Transfer Costs for Common Architectures](#)

Whitepapers

- [Cost Optimization Pillar - Well-Architected Framework](#)
- [Amazon EC2 Reserved Instances and Other Reservation Models](#)
- [Right Sizing: Provisioning Instances to Match Workloads](#)

Conclusion

Games are designed to deliver entertainment experiences to a global audience of players, and have usage characteristics that are typically unpredictable and variable. The Games Industry Lens describes the common types of scenarios that typically comprise a game architecture, and provides a set of questions and best practices to consider when you are building and operating games in the cloud. By applying this framework to your game architecture, you will be able to build reliable, secure, efficient, and cost-effective games in the cloud.

Contributors

The following individuals contributed to this document:

- Kyle Somers – Principal Solutions Architect, Amazon Web Services
- Hyobin An – Solutions Architect, Amazon Web Services
- Peter Chapman – Mgr, Solutions Architecture, Amazon Web Services
- Nirav Doshi – Principal Technical Account Manager, Amazon Web Services
- Chris Finch – Senior AMER Digital User Engagement Specialist Solutions Architect, Amazon Web Services
- Jackie Jiang – Senior Solutions Architect, Amazon Web Services
- Sungsoo Khim – Solutions Architect, Amazon Web Services
- Byungsu Kim – Solutions Architect, Amazon Web Services
- Minsuk Kim – Solutions Architect, Amazon Web Services
- Pawan Matta – Senior Solutions Architect, Amazon Web Services
- Jinsung Park – Associate Solutions Architect, Amazon Web Services
- Tanya Rhodes – Senior Solutions Architect, Amazon Web Services
- Bruce Ross – Senior Media Solutions Architect, Well-Architected, Amazon Web Services

Document history

To be notified about updates to this whitepaper, subscribe to the RSS feed.

Change	Description	Date
Minor update	Definitions section added.	November 19, 2021
Initial publication	Whitepaper first published.	November 19, 2021

Note

To subscribe to RSS updates, you must have an RSS plug-in for the browser you are using.

Notices

Customers are responsible for making their own independent assessment of the information in this document. This document: (a) is for informational purposes only, (b) represents current product offerings and practices, which are subject to change without notice, and (c) does not create any commitments or assurances from and its affiliates, suppliers or licensors. products or services are provided “as is” without warranties, representations, or conditions of any kind, whether express or implied. The responsibilities and liabilities of to its customers are controlled by agreements, and this document is not part of, nor does it modify, any agreement between and its customers.

© 2021 Amazon Web Services, Inc. or its affiliates. All rights reserved.

AWS Glossary

For the latest AWS terminology, see the [AWS glossary](#) in the *AWS Glossary Reference*.