

[illegible]

Provisioned IOPS SSD volumes that is designed to provide 100X durability of 99.999% as well as a 10X higher IOPS to storage ratio of 500 IOPS for every provisioned GB—at the same price as the previous generation (i01) |||

Provisioned IOPS SSD io2 Block Express: Use Cases: io2 Block Express offers the highest performance block storage in the cloud with 4x higher throughput, IOPS, and capacity than io2 volumes, along with sub-millisecond latency. Block Express is the next generation of Amazon EBS storage server architecture purpose-built to meet the performance and latency requirements of the most demanding applications. ||| Throughput Optimized HDD for frequently

[illegible][illegible][illegible]

AMAZON API GATEWAY

[illegible]

Development pipeline: Build-Test-Release-Monitor-Plan-Build | Important features = Lead time, Deployment frequency, Mean time between failure(MTBF), Mean time to recover(MTRR)
Deployment: Source code=Developer workstation/Dev AWS ac | Code reviews = Works on smallest thing possible, explicitly defines dependencies(including those outside of planning), published changes for review, peer review needed
Production: Every Checkpoint (git) of development branch is deployed to production
Example pipeline (Development phase): Code review=Dependencies+Package build and runs unit testing, Bundles code and run-time dependencies into a combined artifact, providence of dependencies is tracked all small and fast, THOROUGH | **Final Phase:** Beta = run destructive or any other testing, smoke testing, will want > Gamma = isolated, "prod-like", integration testing full "end-to-end", load testing, security testing | **Production:** Start small and fast, deploy to users, give every wave time to "bake", always ready for any reason for fail, roll back on test failures
New releases: New releases are planned by release manager, release manager provides regular service, reduce impact of change, give insight as to what's going on, protect customers and the business
How CI/CD affects organizations: Low weekly/monthly: failure rate/yearly: ~60% | Medium(weekly/monthly): ~18% | High(daily/weekly): ~9% | Etkinon demand: ~8% | Only 18% of organizations use sites

[illegible][illegible]

Rolling update. Version B is slowly rolled out in terms of microservices and replaces version A. In a rolling deployment is a deployment strategy that slowly replaces previous versions of an application with new versions of an application by completely replacing the infrastructure on which the application is running. For example, in a rolling deployment in Amazon ECS, containers running previous versions of the application will be replaced one-by-one with containers running the new version of the application. This allows the application to be updated without downtime and without the need to create new infrastructure. This strategy is useful for applications that are stateless and can be scaled horizontally. The old and new application versions. This allows rolling deployments to complete more quickly, but also increases risks and complicates the process of rollback if a deployment fails. Rolling deployment strategies can be used with most deployment solutions. Refer to [CloudFormation Update Policies](#) for more information on rolling deployments with CloudFormation; Rolling Update with Amazon ECS for more details on rolling deployments with Amazon ECS.

Environment (blue) is released alongside the current version and then the traffic is switched to version B. A blue/green deployment is a deployment strategy in which you create two separate, but identical, environments. One environment (blue) is running the current application version and one environment (green) is running the new application version. Using a blue/green deployment strategy increases application availability and reduces deployment risk by ensuring the rollback path is always available. For more information, see [Blue/Green Deployments on AWS](#).

Canary deployment is a deployment strategy in which you release a new version of an application to a small percentage of users and then gradually increase the percentage of users to which you release the new version. AWS deployment services support blue/green deployment strategies including Elastic Beanstalk, AWS OpsWorks, CloudFormation, CodeDeploy, and Amazon ECS. Refer to [Blue/Green Deployments on AWS](#) for more details and strategies for implementing blue/green deployment processes for your application.

Rollback is the process of reverting to a previous version of an application (e.g., a new version of an application is deployed). The purpose of a canary deployment is to reduce the risk of deploying a new version that impacts the workload. The method will incrementally deploy the new version, making it visible to new users in a slow fashion. As you gain confidence in the deployment, you will deploy it to replace the current version in its entirety. Steps: Use a router or load balancer that allows you to send a small percentage of users to the new version. Use a dimension on your KPIs to indicate which version is reporting the metrics. Use the metric to measure the success of the deployment.

Cloud Native Continuous delivery | Supports multi-tenant architecture | Updates | Model and visualize your software release process | Builds, tests and deploys your code every time there is a code change | Integrates with third-party tools and AWS
ACM/DNS AWS X-Ray, Amazon CloudWatch, Amazon DevOps-Git | Updates | ML-powered cloud operations service to improve application availability | Select coverage: Select the entire account or specific CFN stacks | Select analysis: metrics-to-log generates insights | Integrate with your workflow: Integrate with OpsCenter/SNS/ServiceNow
CDK CDK, AWS CloudFormation, Cloud Development Kit (AWS CDK, CDKs, CDK-terraform) | AWS Serverless Application Model(AWS SAM), AWS Amplify
Infrastructure as Code(IaC) What? Writing code to create, configure and deploy infrastructure components | Infrastructure includes: networking, compute, databases, security management tools etc. || Why? Makes infrastructure creation repeatable and predictable | Documents your infrastructure | Automates the deployment of infrastructure | Automates the testing of infrastructure | Automates the rollback of infrastructure | Automates the automation IaC with AWS CloudFormation: Code your templates+Upload, test, review changes+Execution of changes creates a stack+Manage stacks and stack sets

with [AWS Cloud Development Kit \(CDK\)](#): `cdk init -l npm run build -> cdk synth -> cdk diff -> cdk deploy`

GENERATIVE AI

Generative AI is powered by foundation models. Proliferated on vast amounts of unstructured data. Contains a large number of parameters that make them capable of learning complex concepts. Can be applied in a wide range of contexts. Customized fine-tuning your data for domain-specific tasks.

Enhance Customer Experiences: Chatbots, Virtual Assistants, Conversation Analytics, Personalization | Boost employee productivity & creativity: Conversational Search, Summarization, Content Creation, Code Generation, Data To Insights | Optimize business processes: Document Processing, Data Augmentation, Cybersecurity, Process Optimization | Healthcare & Life Sciences: Ambient clinical scribe, Medical imaging, Drug discovery, Enhance clinical trials | Financial Services: Fraud detection, Credit risk assessment, Robo-advisors, Personalized financial planning | Retail: Product recommendations, Personalized marketing, Supply chain optimization, Intelligent advisory, Fraud detection, Compliance assistance | Retail: Price optimization, Virtual try-ons review, Marketing Optimization, Product descriptions, Press recommendations | Media & Entertainment: HQ

content at scale. Enrich broadcast content. Automated content distribution. Optimize subscriber experience. Automated highlights gen. [||](#)

Section 4: AI-Driven Infrastructure for ML Training and Inference

GPU's | Translating Inference | SageMaker | UltraClusters | EFA | EC2 Capacity Blocks | Nitro | Neuron

CGI NVIDIA Tesla M2090 "Fermi" GPUs | G2 NVIDIA DGX QK140 "Kepler" K80 | P2 NVIDIA K80 GPUs | G3 NVIDIA Tesla M60 GPUs | P3 NVIDIA V100 Tensor Core GPUs | G4 NVIDIA V100 Tensor Core GPUs | G4 NVIDIA T4 Tensor Core GPUs | P4 NVIDIA A100 Tensor Core GPUs | G5 NVIDIA A40 Tensor Core GPUs | G5 NVIDIA T40 Tensor Core GPUs | G5 NVIDIA H100 Tensor Core GPUs | A4 Higher Throughput | L4 Lower Latency

Amazon SageMaker and deploy ML models at scale | Automatic model fine-tuning & distributed training | Flexible model deployment options | Tools for ML operations like [||](#) Support for Hugging Face models and Hugging Face AWS Deep Learning Containers

Generative AI Stack: Tools to Build with LLMs and Other FMs

Amazon Bedrock | Quantized & Distributed Inference | [||](#)

- **CloudWatch**: The easiest way to monitor AWS resources with LLMs and other APIs. [Choice of industry-leading FPs from A1Z1 Labs (JURASSIC-2), Amazon (AMAZON TITAN), Anthropic (CLAUDE - Turkish support), Cohere (COMMAND + EMERD), Meta (LLAMA 2), Mistral AI (Mistral7B, Mistral v8x), and Stability AI (STABILITY DIFFUSION XL)] Customizes fops using your organization's data (Enterprise-grade security and privacy) [Customizable LLMs] [Customizable LLMs] [Customizable LLMs] [Customizable LLMs] [Customizable LLMs] [Customizable LLMs] [Customizable LLMs] [Customizable LLMs] [Customizable LLMs] [Customizable LLMs]
- **Amazon Bedrock**: Amazon's managed service for generative AI. [AI data is encrypted in transit and at rest.] Data used to customize models resides entirely within VPC. | Support for standards, including GDPR & HIPAA | Amazon Bedrock simplifies: Choice, Customization, Integration.
- **CloudTrail**: CloudTrail logs API calls made to AWS services. [CloudWatch Integration (Track usage metrics and build customized dashboards), CloudTrail Integration (Monitor API activity and troubleshoot issues), SOC (System and Organization Controls) [Compliance and audit trail], Cost Explorer (Analyze costs and optimize spending).]
- **AWS IAM**: Identity and Access Management (IAM) manages user access to AWS resources. [AWS IAM provides fine-grained control over who can access what resources in your account. It helps you manage users, groups, roles, permissions, and policies. You can use IAM to restrict access to AWS services and resources based on conditions such as IP address ranges, geographic locations, or specific times of day.]
- **Amazon GuardDuty**: Detects malicious activity and threats to your AWS resources. [GuardDuty uses machine learning to analyze log data and identify potential threats. It can detect suspicious activity such as unauthorized access attempts, malware infections, and data exfiltration.]

Quadrants for Amazon Sagemaker: Safeguard your Amazon AI applications with your responsible AI policies | Easily configure harmful content filtering based on your responsible AI policies | Apply Guardrails to any PM or agent toolchain | Personalize identifiable information in PM responses (coming soon)

Amazon AI at Scale: APPLICATIONS THAT LEVERAGE LLMs AND OTHER PMs

Amazon Q | Amazon Q in Amazon QuickSight | Amazon Q in Amazon Connect | Amazon Q ChatWhisperer

QuickSight | **Amazon Q** | AI-powered code suggestions in the IDE and the command line | Customization capabilities: Better and more relevant code suggestions | Your data is never used to train the underlying model, only Amazon Q's model is trained on your data | Onboard developers faster | Can co-develop personalized CodeWhisperer by fine-tuning it if you have 100MB code

Amazon Q | A generative AI-powered assistant for work that is tailored to your business | Provides interactive answers, solves problems, generates content, and takes action | Understands your company information, code, and systems | Personalizes interactions based on your role and permissions | Built to be secure and private | Trained on 10 years of AWS knowledge | Assists you everywhere you work with AWS - in the console, IDE, and systems

Amazon Q :: Your AWS Expert, Business Expert | Amazon Q :: Your Intelligence Expert | Amazon Q :: Your Contact Center Expert

AWG Cloud DATABASE
AWG Data is a fully managed serverless environment where you can extract, transform, and load (ETL) data at scale. With AWS Glue, you can categorize data, clean it, enrich it, and move it reliably across various data stores and streams in a cost-effective manner. (in the scope of Analytics)
 The problem: Data grows exponentially: From new resources! Increasingly driven! Used by people with diverse facilities! Accessed by many applications! | Data preparation requires huge amount of time (80%) | It is time consuming to extract, clean, and normalize data | Data is scattered in different formats and load data at scale & schedule right tools to be integrated | Costly user licenses & skilled tools causing network & moving large data out of VPC | Hard to operationalize and build repeatable workflows & needs a lot of code-based heavy-lifting for it to work at scale.
AWG Data Solution – Database is a no-code data preparation tool that you can use to visually explore, clean, and transform data. You can choose from more than 250 prebuilt transformations to automate data preparation tasks without writing any code (no SQL). Visualize data quality (understands data quality (patterns & anomalies) | Cleans and normalizes data | Visually maps data lineage (understands how data has been through) | Automates data scale (prepared for new data).

Functionality: Connect to data sources (S3, Redshift, RDDBase, Oracle DataPilot) 256+ built-in transformations | For users of all technical levels (visuals) | Advanced data profiling/get the statistics | Visual data lineage | Integrate with data pipelines | serverless usage | **Visualizations:** Interactive charts | Interactive Visual ETL IDE for non-Spark experts & Dev endpoints: interactive development for Glue jobs for Spark experts | Streaming: Process IoT streams in real time | Glue Studio - Develop streaming jobs visually & Glue Streaming - processes real time data stream from IoT devices | Data Replication: Replicate data across purpose built data catalog (AWS Elastic View - replicate data across multiple data stores without code) | Data Prep: Transform data without code (Glue Data Brew - build data without coding) | **Integrations:** Amazon Athena, Amazon Redshift, Local file | S3 & data catalog stores | **Integrations:** AWS Glue Data Catalog - AWS Glue Data Warehouse - AWS Glue Data Transfer - S3 output bucket - Amazon QuickSight | [Set up data quality rules with AWS Lambda] Recurring raw data feed - S3 - AWS Lambda - AWS Glue Data Warehouse - Amazon EventBridge - AWS Lambda - Amazon SNS - Email notification | Data preprocessing for ML | [Orchestrating data preparation in workflows] | **Highlighted features:** Enrich data with unions and joins | update schema | generate data profiles | feature engineering-plugin in Notebooks | remove outliers | analyze aggregated data

AMAZON CLOUDWATCH

[Monitor] Derives observability on a single platform applications and infrastructure | Easiest way to collect metrics from AWS and on-premises | Improve operational performance and resource optimization | Get operational insights and insight | Correlate actionable insights from logs & CloudWatch Metric-Monitor-Architecture Actions across all your resources, applications and services (Run on AWS and On-premises servers)

[Monitor]: Visualize applications and infrastructure with CloudWatch dashboards; correlate logs and metrics side by side to troubleshoot and alert with CloudWatch Alarms | Act: Automate response to operational changes through CloudWatch Events and Auto Scaling | **Analyze:** Up to 1-second metrics, extended data retention(15minutes) and real-time analysis with CloudWatch Metric Math ||| **User Prof.**: Application Monitoring, System-wide Visibility, Resource Utilization

[Collect] Collect logs from: EC2 instances, on-premises servers, VPC Flow Logs, CloudTrail, Lambda, other AWS Services ||| Log data can be stored and accessed indefinitely in highly durable, cost-optimized storage so you don't have to worry about filling up hard drives ||| Built-in metrics: Collecting metrics is time consuming. CloudWatch allows you to collect default metrics from more than 70 AWS services, such as: EC2, DynamoDB, S3, ECR, Lambda, API Gateway etc.

Your own applications to monitor operational performance, troubleshoot issues, and spot trends. User activity is an example of a custom metric you can collect and monitor over a period of time. Publish metrics using the AWS CLI or the API | Standard realizations, with a one-minute granularity | High resolution, with a granularity of one second | Aggregate data before you publish to CloudWatch | Statsd and collect support via CloudWatch Agent | Collect metrics from your servers, aggregate and send them to CloudWatch | Metrics and logs from your containerized applications and microservices. Collects metrics from each container; CPU, Memory, Disk Network | Automatically generated dashboards | Set alarms on metrics

Monitor. Amazon CloudWatch dashboards enable you to create re-usable graphs and visualize your cloud resource usage and applications in a unified view. [A single view for selected metrics and alarms] Multiple AWS accounts and regions are supported. Playbooks help you define and execute remedial actions based on detected anomalies. You can also create custom dashboards to visualize metrics and alarms across multiple AWS accounts. Alerts trigger an action. [When a single alarm or the result of a logical expression] Add alerts to dashboards to visualize them | **Perform actions** based on the value of metrics (Send a notification to SNS topic, Auto Scaling action, EC2 Action (Stop, Terminate, Restart or Recover)) | Logs and metric correlation. CloudWatch also makes it easy to correlate metrics and logs. [Manage logs and metrics in a single platform.] Use metric filters to convert log events to metrics.

[illegible][illegible]

Evolute patterns in structured log events • Display on CloudWatch dashboards • Add to CloudWatch alarms

AMAZON FOR GAMES

• AWS for Games solution areas and use cases: Build, Run, Grow. • **BUILD:** Cloud game development (Workstations, Build pipelines, Version control, 3D world building) **RUN:** Game servers/Hosting session-based games, Global game infrastructure () • Game security (Defend against DDoS attacks, Protect against data breaches) **SHOW:** LiveOps (Game Services for LiveOps) () • Game analytics (Centralized game analytics) () | AI & ML, Community health/safety, and more

• AWS For Games Provides Solutions For Games Customers Wherever They Are In The AWS-Builder-Continuum: Native AWS services, game-specific AWS services, AWS solutions & Quick Starts, AWS ProDev & AWS Partner Solutions & AWS Marketplace

• **AMAZON GAMES SERVERS (GameServer)** (Amazon build cloud products) () | **AMAZON LITMUS (Ready-to-deploy solutions assemble AWS Services code and configurations)** () | **PARTNER ROI TOOLS (Software, Staff, of AWS)** ()

managed services from AWS Partners) | [GUIDANCE](#) (Preventive architectural diagrams, sample code, and technical support) | [Amazon and AWS Game Services](#) (Dedicated game server hosting solution for multiplayer games) | [Open 20 Engine](#) (ODD/AAR-capable, cross-platform, open-source game engine available under an Apache 2.0 license) | [Amazon Luna](#) (Amazon Game Studio) | [Amazon GameLift](#) (Low-cost game servers, scale your servers automatically based on player demand and leverage low cost Spot Instances for short-lived game sessions) | [Amazon GameLift](#) (Low latency, Deploy a single fleet globally and leverage built-in latency-based matchmaking. High availability (fully managed game servers hosted across AWS Regions on multiple Availability Zones, as well as Local Zones, providing availability and scalability). Flexibility (Select the features you need. Use the built-in matchmaking or build your own. Select between fully managed and self-managed 11 hosting).

[illegible]

Amazon GameLift **Core** provides fully managed game server hosting. Upload your build, deploy globally. Managed latency-based session placement. Supports on-demand and spot fleets for cost optimization with spot capacity allocation. Fully managed game server hosting. Upload your build, deploy globally. Managed latency-based session placement. Supports on-demand and spot fleets for cost optimization with spot capacity allocation.

Amazon GameLift **Flex** - More flexible. Hosted on EC2 on your AWS account. Simple API layer for game session management. Use your existing tools and software to develop and running game servers. Access GameLift Spot visibility anywhere independent of other GameLift services. AWS SDK on the server side; use the programming language of your choice.

Amazon GameLift **Custom** - Use GameLift Core, GameLift Flex, or GameLift Spot. Use GameLift Services (client communicates with GameLift FlexMatch (use the free build in matchmaking based on your needs) - GameLift Queue (Queues route traffic to the lowest latency region) - Fleet (A single fleet can span across up to 23 different AWS Regions).

Amazon GameLift **Components** - Builds | Flavors | Alphas | Clients | FlexMatch

Amazon GameLift **Tools** - Amazon AWS CLI - Amazon CloudWatch (for monitoring) - Amazon CloudFormation (for provisioning) - Amazon IAM (for access control) - Amazon Instance Configuration | Supported Operating Systems - Amazon Linux - Windows - Server SDKs - C++ - Unreal Plugin (SDK) - Unity Plugin (Fully-featured) | Quick Dashboard - Access key information about builds stored in Amazon GameLift - Status - Version - OS - Size - Number of fleets using this build

[illegible][illegible]

Queues - Fleet Policies: Prioritize your destinations (Fleet Alises) to indicate the preferred order. | Amazon GameLift uses the priority to locate the best available session, starting with the first priority and moving down. | **⚠️** When using Spot Instances set this Fleet to a higher priority, then set an On-Demand Fleet below for cases where Spot capacity is not available.

Queues - Fleet Policies: Set a **Maximum Session Duration** to specify the maximum amount of time a session can last. | **⚠️** For On-Demand Fleets, this value must have a value of 1 or greater. | Specify the amount of time to enforce the policy | **⚠️** Multiple policies can be specified | **⚠️** FleetMatch also does this but Quesue latency configuration allows you to control overflow to other Regions in case of no available sessions.

Queues - Best Practices: Use Multi-region Fleets registered to a Queue to deploy globally with a single build upload and Fleet configuration. | Use a Spot Fleet on high priority, and prepare a fallback on-demand Fleet. | For disaster recovery, create a second Queue in a different Region and use a Spot Fleet on high priority. | All Fleets must have a T2S certificate, regardless whether enabled or disabled.

Client: Amazon GameLift supports any game client or game service that is able to make use of one of the AWS supported SDKs, languages include: C++ - C# - Go - Python - JavaScript/Node.js - Java. | Always build a backend service that calls GameLift APIs instead of the game client!

© 2023 Amazon.com, Inc. or its affiliates. All rights reserved. Amazon, the Amazon logo, GameLift, and the GameLift logo are either registered trademarks or trademarks of Amazon.com, Inc. or its affiliates in the United States and/or other countries.

GameList | Capture intelligence from Clients by using AWS endpoints. This can be an example reference to [DynamoDB endpoints](#) available across the regions

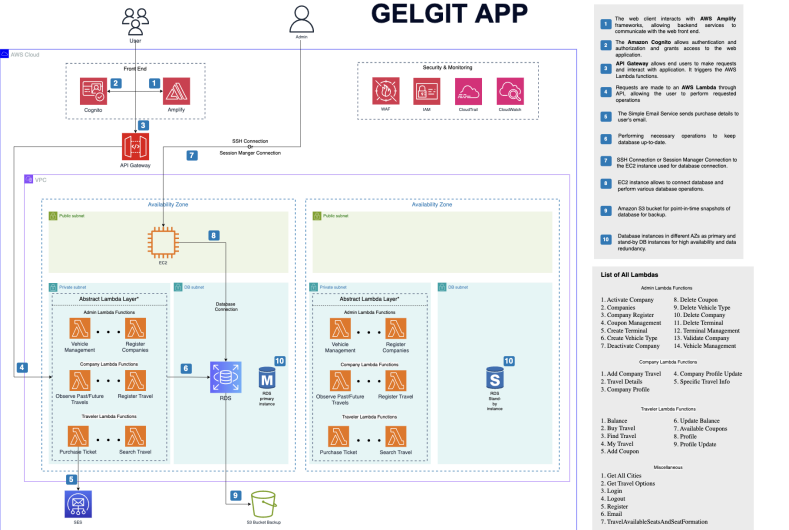
Amazon ElastiCache | More flexibility: Game Server Group hosted on EC2 on your AWS account • Simple API layer for game session management • Manages Fleet scaling for you • Use your existing tools and software to deploy and manage game server processes • Access GameList API: Spot visibility independent of other Amazon features • AWS SDK on the server side: use the programming language of your choice • It keeps load same for the all

How GameList Elastic works | Use AWS SDK, CLI or CloudFormation to create a Game Server Group (an EC2 Auto Scaling Group is created with the configuration of instances, FleetIt automatically scales this and manages spot and on-demand balance based on spot visibility | You control how the game server processes are run on the instances | Game Servers report their state to FleetIt. Your Backend calls [ClientGameServer\(\)](#) to request a few sessions |

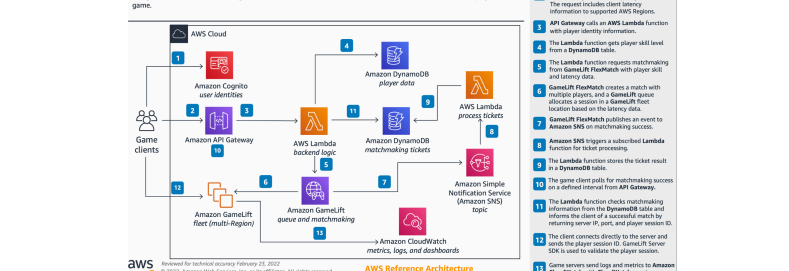
GameList Elastic | GameList Elastic is a managed service that provides a scalable and reliable way to manage your game server instances. It is designed to be used with AWS Lambda, as well as standalone with your own game server hosting solution | [Player teams support](#) | [Latency-based matching](#) | [Role retaining](#) | [Match acceptance](#) | [Best region placement](#) | [Player drop in/out support](#) with [FlexMatch](#) backfill

[illegible]

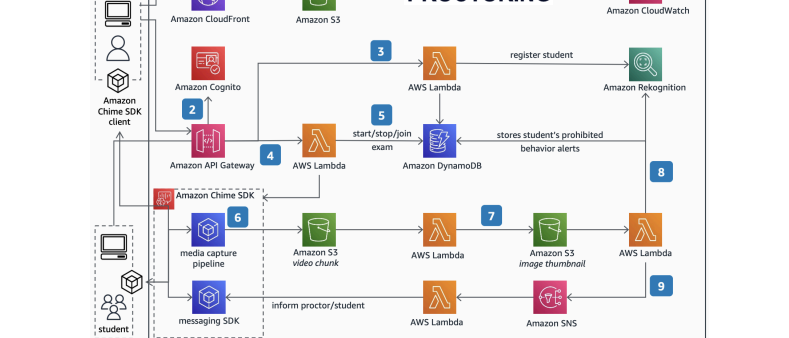
GELGIT APP



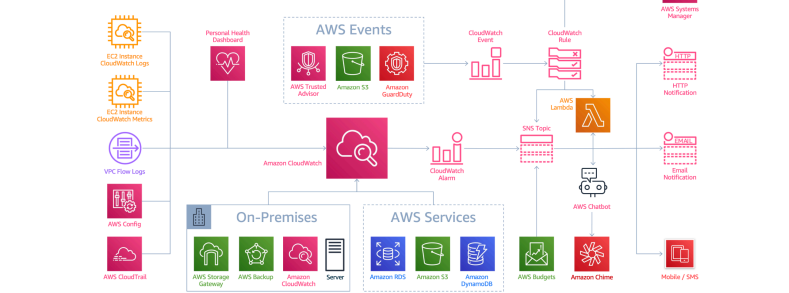
Multiplayer Session-based Game Hosting on AWS



PROCTORING



Integration



TAY

