

AWS CERTIFICATION GUIDE AND INTERVIEW Q&A

**100+
PAGES!**

Crack your AWS Data Engineering
Certification and Interviews



AWS Data Engineering Interview Questions

Deepa Vasanthkumar

AWS General Topics	3
Amazon S3	7
AWS Virtual Private Cloud (VPC)	10
Amazon EC2 (Elastic Compute Cloud)	13
Amazon EMR (Elastic MapReduce)	16
AWS Glue	20
Glue & EMR Comparison	23
Glue Data Catalog	24
Glue Streaming ETL	27
AWS Lambda	31
Amazon Kinesis	35
Amazon Athena	39
Amazon Redshift	43
Amazon RDS for PostgreSQL	47
Amazon DynamoDB	49
AWS Step functions	53
AWS Identity and Access Management (IAM)	58
Amazon Simple Notification Service (SNS)	61
Amazon Simple Queue Service (SQS)	64

AWS General Topics

AWS (Amazon Web Services) is a comprehensive, evolving cloud computing platform provided by Amazon. It provides a mix of infrastructure as a service (IaaS), platform as a service (PaaS), and packaged software as a service (SaaS) offerings. AWS services can offer organizations tools such as compute power, database storage, and content delivery services.

Here are some common AWS interview questions along with detailed answers to help you prepare for your interview.

Q. Explain the difference between Amazon S3 and EBS.

Amazon S3 (Simple Storage Service) is an object storage service that offers industry-leading scalability, data availability, security, and performance. This means customers of all sizes and industries can use it to store and protect any amount of data for a range of use cases, such as websites, mobile applications, backup and restore, archive, enterprise applications, IoT devices, and big data analytics.

Amazon Elastic Block Store (EBS) provides block level storage volumes for use with EC2 instances. EBS volumes are highly available and reliable storage volumes that can be attached to any running instance that is in the same Availability Zone. EBS is particularly suited for applications that require a database, file system, or access to raw block level storage.

Q. What is IAM in AWS?

IAM (Identity and Access Management) is a web service that helps you securely control access to AWS resources. You use IAM to control who is authenticated (signed in) and authorized (has permissions) to use resources. IAM provides the ability to control users, security credentials such as access keys, and permissions that control which AWS resources users and applications can access.

Q. Describe a VPC and its components.

A VPC (Virtual Private Cloud) allows you to provision a logically isolated section of the AWS Cloud where you can launch AWS resources in a virtual network that you define. The main components of VPC are:

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar -| LinkedIn](#)

- Subnets: A range of IP addresses in your VPC.
- Route Tables: They determine where network traffic is directed.
- Internet Gateways: Connects a network to the internet.
- NAT Gateways: Allow private subnets to connect to the internet or other AWS services but prevent the internet from initiating a connection with those instances.
- Security Groups: Act as a virtual firewall for your instance to control inbound and outbound traffic.
- Network ACLs: A layer of security for your VPC that acts as a firewall for controlling traffic in and out of one or more subnets.

Q. How do you secure data at rest on AWS?

Securing data at rest on AWS can be achieved by several methods:

- Encryption: Using AWS services such as Amazon S3, EBS, and RDS which support encryption at rest. AWS offers integrated tools like AWS KMS (Key Management Service) and AWS CloudHSM to manage and rotate encryption keys.
- Access Control: Implementing fine-grained access control using IAM roles and policies to ensure only authorized users and systems can access your data.

- Network Security: Utilizing VPCs and associated security features such as security groups and network ACLs to isolate resources and control network access.

Q. What are some strategies to reduce costs in AWS?

- Reserved Instances: Purchase Reserved Instances for services like Amazon EC2 and Amazon RDS to save up to 75% compared to on-demand pricing.
- Auto Scaling: Use Auto Scaling to automatically adjust the amount of computational resources based on the server load, thus reducing costs by minimizing idle resources.
- Right Sizing: Regularly review and adjust your configurations to ensure you are using the optimal resources for your workloads.
- Delete Unattached EBS Volumes: Regularly delete unattached EBS volumes, as you are billed for the storage that you provision.
- Use Cost Explorer: AWS Cost Explorer helps you visualize and manage your AWS spending over time.

Amazon S3

Amazon S3 (Simple Storage Service) is a key service in AWS, offering object storage through a web service interface. Here are some common interview questions about Amazon S3, accompanied by detailed answers to help you prepare for a technical interview.

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar -| LinkedIn](#)

Amazon S3 is an object storage service offered by Amazon Web Services that provides scalable, high-speed, web-based storage for data backup, archival, and analytics. Unlike file storage, which organizes data into a directory hierarchy, S3 stores data as objects within buckets. Each object is identified by a unique key and can contain up to 5 TB of data.

Q. Can you explain the difference between S3 and EBS? - Amazon S3 is object storage good for storing vast amounts of data in a non-structured format. Data is accessible from anywhere and is suited to static files, backups, and big data analytics.

- Amazon EBS (Elastic Block Store) provides block-level storage volumes for persistent data storage for use with EC2 instances. EBS is suitable for applications that require a database, file system, or access to raw block level storage.

Q. What are S3 Buckets?

An S3 bucket is a container for storing objects in Amazon S3. Each object is stored in a bucket and retrieved via a unique, developer-assigned key. Buckets serve as the basic container in which data is stored in S3, and every object must be contained in a bucket.

Q. Describe the durability and availability of Amazon S3.

Amazon S3 provides high durability and availability. S3 is designed to deliver 99.999999999% (11 9's) durability over a given year, ensuring that data is virtually indistinguishable from being lost. Furthermore, S3 Standard aims for 99.99% availability, while the S3 Standard-IA (Infrequent Access) and One Zone-IA provide slightly lower availability at a reduced cost.

Q. What are some of the storage classes available in S3?

S3 offers several storage classes:

- S3 Standard: For frequently accessed data, offers high durability, availability, and performance.
- S3 Intelligent-Tiering: Moves data automatically between two access tiers when access patterns change.
- S3 Standard-IA: For data that is less frequently accessed but requires rapid access when needed.
- S3 One Zone-IA: Similar to Standard-IA but data is stored in a single availability zone.
- S3 Glacier and S3 Glacier Deep Archive: For archiving data with retrieval times ranging from minutes to hours.

Q. What is S3 Versioning?

S3 Versioning is a feature that allows you to keep multiple versions of an object in the same bucket. This is used to preserve, retrieve, and restore every version of every object stored in your S3 bucket. It helps protect from unintended overwrites and deletions.

Q. How does S3 Encryption work?

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar -| LinkedIn](#)

S3 provides two means of encryption:

- Server-Side Encryption (SSE): Where Amazon handles the encryption process, the decryption, and key management. You can choose among three keys management options: SSE-S3 (uses keys managed by S3), SSE-KMS (uses AWS Key Management Service), and SSE-C (where you manage the encryption keys).
- Client-Side Encryption: Your data is encrypted on the client side before uploading it to S3.

Q. Explain the S3 Lifecycle Policies.

Lifecycle policies in S3 are used to manage your objects automatically as they transition to different stages of their lifecycle. Policies can be used to transition objects to less expensive storage classes or schedule objects for automatic deletion after certain periods, helping reduce costs by managing data according to its lifecycle.

Q. What is the difference between S3 and S3 Glacier?

S3 is used for general, high-speed storage accessible at any time, suitable for a wide range of applications including websites, mobile apps, and backups. S3 Glacier is a lower-cost storage service optimized for data archiving and long-term backup, where retrieval times of several minutes to hours are acceptable.

AWS Virtual Private Cloud (VPC)

AWS Virtual Private Cloud (VPC) is a crucial component of AWS that allows you to provision a logically isolated section of the AWS Cloud where you can launch AWS resources in a virtual network that you define.

Q. What is a VPC? An Amazon Virtual Private Cloud (VPC) is a service that allows you to launch AWS resources into a virtual network that you've defined. This virtual network closely resembles a traditional network that you'd operate in your own data center, with the benefits of using scalable infrastructure on AWS. It enables you to control your virtual networking environment, including selection of your IP address range, creation of subnets, and configuration of route tables and network gateways.

Q. What are subnets?

In AWS, a subnet is a range of IP addresses in your VPC. You can launch AWS resources into a specified subnet. Use a public subnet for resources that need to be connected to the internet, and a private subnet for resources that won't be connected to the internet. Each subnet must reside entirely within one Availability Zone and cannot span zones.

Q. Explain the difference between a security group and a network access control list (NACL).

- Security Group: Acts as a virtual firewall for your EC2 instances to control inbound and outbound traffic. Security groups operate at the instance level, they support allow rules only, and are stateful—responses to allowed inbound traffic are allowed to flow outbound regardless of outbound rules.

- Network Access Control List (NACL): Acts as a firewall for controlling traffic in and out of one or more subnets. NACLs operate at the subnet level, they support allow and deny rules, and are stateless—responses to allowed inbound traffic are subject to the rules for outbound traffic (and vice versa).

Q. What is an Internet Gateway, and why is it used? An Internet Gateway (IGW) is a VPC component that allows communication between instances in your VPC and the internet. It provides a target in your VPC route tables for internet-routable traffic, and performs network address translation for instances that have been assigned public IPv4 addresses.

Q. How do you securely connect your on-premises network to a VPC?

We can securely connect your on-premises network to your Amazon VPC using one of the following options:

- AWS Direct Connect: Establishes a dedicated private connection from a remote network to your VPC. Direct Connect is often used to reduce network costs, increase bandwidth throughput, and provide a more consistent network experience than internet-based connections.

- VPN Connection: Utilizes the public internet to make a secure and encrypted connection between your premises and your VPC. This is accomplished by creating a VPN gateway in your AWS environment, which provides a highly available and secure connection.

Q. What is a NAT Gateway, and why would you use it? A Network Address Translation (NAT) Gateway enables instances in a private subnet to connect to the internet or other AWS services, but prevent the internet from initiating a connection with those instances. This is particularly useful for updating software in private subnets, or enabling secure outbound internet access for processing tasks.

Q. Describe how routing works in a VPC.

Each VPC has a routing table that contains a set of rules, called routes, that are used to determine where network traffic from your VPC is directed. Each subnet in your VPC must be associated with a route table; the table controls the routing for the subnet. The default route table includes a local route for communication within the VPC, and additional routes can be added that direct traffic to specific destinations, such as an internet gateway, NAT gateway, VPN connection, or AWS Direct Connect.

Q. What are VPC Endpoints? VPC endpoints enable you to privately connect your VPC to supported AWS services and VPC endpoint services powered by PrivateLink without requiring an internet gateway, NAT device, VPN connection, or AWS Direct Connect connection. Instances in your VPC do not require public IP addresses to communicate with resources in the service. Traffic between your VPC and the other service does not leave the Amazon network.

Amazon EC2 (Elastic Compute Cloud)

Amazon EC2 (Elastic Compute Cloud) is a central part of Amazon's cloud computing platform, AWS. It allows users to rent virtual computers on which they run their own computer applications.

Q. What is Amazon EC2? Amazon Elastic Compute Cloud (EC2) is a web service that provides resizable compute capacity in the cloud. It is designed to make web-scale cloud computing easier for developers. EC2's simple web service interface allows you to obtain and configure capacity with minimal friction. It provides you with complete control of your computing resources and lets you run on Amazon's proven computing environment.

Q. Can you explain the different types of instances available in EC2? EC2 provides a variety of instance types optimized to fit different use cases. Instance types comprise varying combinations of CPU, memory, storage, and networking capacity and give you the flexibility to choose the appropriate mix of resources for your applications. The main categories include:

- General Purpose: Balanced CPU, memory, and networking, suitable for a variety of applications.
- Compute Optimized: Ideal for compute-bound applications that benefit from high-performance processors.
- Memory Optimized: Designed to deliver fast performance for workloads that process large data sets in memory.

- Storage Optimized: Optimized for workloads that require high, sequential read and write access to very large data sets on local storage.
- Accelerated Computing: Use hardware accelerators, or co-processors, to perform functions such as floating-point number calculations, graphics processing, or data pattern matching more efficiently than is possible in software running on CPUs.

Q. What are Amazon Machine Images (AMIs)? An Amazon Machine Image (AMI) provides the information required to launch an instance, which is a virtual server in the cloud. You can think of an AMI as a template of the operating system (OS), the application server, and applications that configure the EC2 instance. AMIs are region-specific and can include one or more EBS snapshots, instance store-backed volumes, permissions to specify which AWS accounts can use the AMI, and a block device mapping.

Q. How do you secure an EC2 instance?

Securing an EC2 instance involves several steps:

- Security Groups: Control inbound and outbound traffic to instances. You should configure security groups to allow the minimum necessary traffic for your instances.
- Key Pairs: Used for secure SSH access to your instances. Always keep your private keys secure.

- IAM Roles: Use IAM roles to securely give applications that run on your EC2 instances permission to AWS API requests.
- Network ACLs: Acts as a firewall for associated subnets, controlling both inbound and outbound traffic at the subnet level.
- Patch Management: Regularly apply security patches to your instance's operating system.
- Encryption: Use AWS services like EBS encryption to protect data at rest.

Q. Explain Elastic Load Balancing. Elastic Load Balancing automatically distributes incoming application traffic across multiple targets, such as EC2 instances, containers, and IP addresses. It can handle the varying load of your application traffic in a single Availability Zone or across multiple Availability Zones. Elastic Load Balancing offers three types of load balancers that all feature high availability, automatic scaling, and robust security necessary to fault-tolerantly distribute application loads.

Q. What are Spot Instances?

Spot Instances allow you to request spare Amazon EC2 computing capacity for up to 90% off the On-Demand price. Spot Instances are recommended for applications that have flexible start and end times, applications that are only feasible at very low compute prices, and users with urgent computing needs for large amounts of additional capacity.

Q. Describe how you can vertically scale an EC2 instance. Vertical scaling refers to increasing the size of an EC2 instance. For example, if you start with a smaller instance

type and realize that your application needs more CPU, memory, or IO, you can stop your instance and change its instance type to a more powerful one, then start it again. This is a simple way to scale for applications that are not designed to scale out across multiple instances.

Amazon EMR (Elastic MapReduce)

Amazon EMR (Elastic MapReduce) is a cloud-native big data platform, allowing businesses to process vast amounts of data quickly and cost-effectively across resizable clusters of Amazon EC2 instances. Here are some common interview questions related to Amazon EMR, along with detailed answers, that might be asked in a technical interview:

Q. What is Amazon EMR? Amazon EMR is a managed cluster platform that simplifies running big data frameworks, such as Apache Hadoop and Apache Spark, on AWS to process and analyze large datasets. The service manages the provisioning, configuration, and tuning of the cloud infrastructure so that users can focus on processing their data. EMR is designed to be cost-efficient and to reduce the complexity of the hardware provisioning, cluster setup, configuration, and tuning of big data frameworks.

Q. How does EMR handle data processing? EMR helps in processing large amounts of data quickly by distributing the data across a resizable cluster of Amazon EC2 instances. It uses popular data processing frameworks like Hadoop, Spark, HBase, Presto, and Hive. Users can write their data processing jobs in various languages like Python, Scala, or Java. EMR takes care of distributing the code and data to the instances, running the jobs, and storing the results in an Amazon S3 bucket or passing them on to other AWS services for further processing.

Q. What are the main components of an EMR cluster?

An EMR cluster typically consists of the following components:

- Master node: Manages the distribution of data and tasks to the other nodes in the cluster. There is only one master node.
- Core nodes: These nodes run tasks and store data in the Hadoop Distributed File System (HDFS) across your cluster.
- Task nodes: Optional nodes that only process data and do not store data in HDFS. They are used to enhance processing power.

Q. What file systems are supported by Amazon EMR?

EMR supports several file systems:

- HDFS (Hadoop Distributed File System): The default, distributed file system used by Hadoop components.

- EMR File System (EMRFS): An extension of Amazon S3, allowing Hadoop to directly interact with data stored in S3. EMRFS is used for scenarios where you want the data to be stored in S3 for durability, cost-effectiveness, or easy scalability.
- Local File System: The disk file system of the EC2 instances in the cluster.

Q. How does Amazon EMR integrate with other AWS services?

Amazon EMR integrates with several AWS services to enhance its capabilities:

- Amazon S3: For durable, cost-effective storage of big data processing results or as a data lake.
- Amazon RDS and Amazon DynamoDB: For direct database queries within EMR jobs.
- AWS Data Pipeline: For workflow management of data between various AWS compute and storage services.
- Amazon CloudWatch: For monitoring cluster performance and running automated actions based on specific conditions.
- AWS Identity and Access Management (IAM): For securing access to EMR resources.

Q. What are the cost optimization strategies for EMR?

To optimize costs in EMR, consider the following strategies:

- Use Spot Instances: Utilize Spot Instances for task nodes where applicable, as they can reduce the cost significantly.
- Right-size the cluster: Choose the right number and type of nodes based on the workload to avoid overprovisioning.
- Shut down idle clusters: Terminate clusters when not in use, or use auto-scaling to scale down resources automatically.
- Use Reserved Instances: For long-running clusters, Reserved Instances provide a significant discount over On-Demand pricing.

Q. Explain how security is managed in EMR. Security in Amazon EMR is managed through:

- IAM Roles: To control access to AWS resources.
- Security Groups: Acts as virtual firewalls to control inbound and outbound traffic to the EC2 instances.
- Encryption: At rest using Amazon S3 with EMRFS, or in transit using TLS across nodes.

- Kerberos Authentication: Provides a mechanism for authentication of users accessing the cluster.

AWS Glue

AWS Glue is a fully managed extract, transform, and load (ETL) service that makes it easy for customers to prepare and load their data for analytics. AWS Glue is a serverless data integration service that makes it easy to discover, prepare, and combine data for analytics, machine learning, and application development. AWS Glue provides both visual and code-based interfaces to make data integration simple. It automatically discovers and catalogs metadata about your data stores into a central catalog that makes the data readily searchable and queryable by services like Amazon Athena and Amazon Redshift Spectrum.

Q. Can you explain the components of AWS Glue?

AWS Glue consists of several key components:

- Data Catalog: A central repository to store structural and operational metadata for all your data assets.
- Crawler: A tool that connects to your data source, progresses through a prioritized list of classifiers to determine the schema for your data, and then creates metadata tables in your data catalog.

- ETL Jobs: Scripts generated automatically by AWS Glue or written by you to transform, flatten, and enrich data in various formats across data stores.
- Triggers: Conditions that you specify for starting ETL jobs and crawlers. Triggers can be time-based or event-based.
- Development Endpoint: An environment to develop and test your ETL scripts interactively.

Q. How does AWS Glue handle schema evolution? AWS Glue can handle schema evolution. When schema changes are detected, AWS Glue updates the schema in its Data Catalog and automatically revises ETL jobs accordingly. This means that when your data structures change, Glue can adapt to accommodate these changes without manual intervention, ensuring that downstream processes are not disrupted.

Q. What are AWS Glue Crawlers and what do they do? AWS Glue Crawlers are used to populate the AWS Glue Data Catalog with metadata tables. These are generated based on the data source's schema it assesses. Crawlers can connect to a source or target data store, classify the data format, infer schemas, and create metadata tables in the Data Catalog. They can be scheduled to run periodically, ensuring that the Data Catalog is kept up-to-date with changes in the data store.

Q. What is the AWS Glue Data Catalog? The AWS Glue Data Catalog is a central repository to store structural and operational metadata for all your data assets. It is fully managed and serves as a persistent metadata store for all your ETL and data discovery

activities across all your data silos. It integrates with services like Amazon Athena, Amazon Redshift Spectrum, and Amazon EMR, providing a unified view of all your data assets.

Q. Explain the types of triggers in AWS Glue.

Triggers in AWS Glue can be classified into three types:

- Schedule-based Triggers: These triggers start jobs at specific times using a cron-like syntax.
- Event-based Triggers: These are activated by specific events, such as the successful completion of another job.
- On-demand Triggers: These are manual triggers that start jobs whenever they are activated by the user.

Q. What is a Job Bookmark in AWS Glue? Job bookmarks in AWS Glue help manage state information and prevent the reprocessing of old data during subsequent runs of an ETL job. When enabled, Glue keeps track of data that has already been processed, which means that only new and changed data is processed in subsequent runs. This feature is particularly useful when dealing with incremental loads.

Q. How does AWS Glue integrate with other AWS services?

AWS Glue integrates with various AWS services to enhance its data integration capabilities:

- Amazon S3: Used as a data source and target for ETL jobs.
- AWS Lambda: Can trigger ETL jobs in response to events.
- Amazon Redshift: Can directly run transformation jobs and output the results into Redshift.
- Amazon Athena: Uses the Data Catalog as a central schema repository.
- Amazon RDS and Amazon DynamoDB: Can be sources or targets for ETL jobs.

Glue & EMR Comparison

AWS Glue is a managed ETL service focused on simplifying data preparation and integration tasks, while Amazon EMR is a fully managed big data platform designed for running distributed data processing frameworks at scale. AWS Glue is a fully managed extract, transform, and load (ETL) service designed to automate data preparation tasks. Amazon EMR is a fully managed big data platform that provides frameworks like Apache Hadoop, Apache Spark, and Presto for processing large datasets.

AWS Glue is primarily used for data integration, cataloging, and ETL tasks. It helps automate the process of discovering, preparing, and combining data for analytics. Amazon EMR is used for processing and analyzing large datasets using distributed

computing frameworks. It's ideal for running big data processing tasks like batch processing, data transformation, and machine learning.

AWS Glue provides Serverless architecture where you define ETL jobs using a visual interface or code. Glue manages the underlying infrastructure, including provisioning resources and scaling based on demand. Amazon EMR has a Managed cluster platform where you provision and configure clusters of EC2 instances to run big data frameworks. You have more control over the cluster configuration and can customize the environment to meet specific requirements.

AWS Glue Provides a simpler and more streamlined approach to ETL tasks, suitable for users who prefer a managed service and don't require fine-grained control over the underlying infrastructure. Amazon EMR offers more flexibility and customization options, allowing users to choose from a wide range of big data frameworks, adjust cluster configurations, and install custom software packages.

Glue Data Catalog

The AWS Glue Data Catalog is a centralized metadata repository that stores structural and operational metadata for all your data assets. It provides a unified view of your data, making it easy to discover, catalog, and query data across different data stores and AWS services.

Q. What are the key features of AWS Glue Data Catalog?

Key features of AWS Glue Data Catalog include:

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar -| LinkedIn](#)

- Metadata Storage: Stores metadata such as table definitions, partitions, schemas, and statistics.
- Data Discovery: Allows users to search and explore data assets using a simple interface.
- Schema Evolution: Supports schema evolution to handle changes in data structures over time.
- Integration: Integrates seamlessly with other AWS services like Amazon Athena, Amazon Redshift Spectrum, and Amazon EMR.
- Custom Metadata: Allows users to define custom metadata attributes and tags for their data assets.

Q. How does AWS Glue Data Catalog compare to traditional metadata management solutions?

Compared to traditional metadata management solutions, AWS Glue Data Catalog offers several advantages:

- Fully Managed: AWS Glue Data Catalog is a fully managed service, eliminating the need for infrastructure management.
- Scalability: It scales automatically to handle large volumes of metadata.

- Integration with AWS Services: It seamlessly integrates with other AWS services, simplifying data integration and analysis workflows.
- Serverless Architecture: It follows a serverless architecture, enabling users to focus on data management tasks rather than infrastructure maintenance.

Q. What are the benefits of using AWS Glue Data Catalog for data governance?

AWS Glue Data Catalog provides several benefits for data governance:

- Centralized Metadata Repository: It serves as a centralized repository for storing metadata, making it easier to manage and govern data assets.
- Data Lineage: It tracks the lineage of data assets, providing insights into data movement and transformation processes.
- Access Control: It supports fine-grained access control, allowing administrators to control who can access and modify metadata.
- Data Quality Monitoring: It enables data quality monitoring by capturing statistics and metrics about data assets.

Q. How can you integrate AWS Glue Data Catalog with other AWS services?

AWS Glue Data Catalog integrates seamlessly with other AWS services through its APIs and native integrations. For example:

- It can be used as a metadata repository for Amazon Athena, enabling users to run SQL queries directly against data stored in Amazon S3.
- It can be integrated with Amazon Redshift Spectrum to query data in Amazon S3 using Redshift.
- It can be used with AWS Glue ETL jobs to transform and load data into various data stores.

Q. What is the significance of custom metadata in AWS Glue Data Catalog? Custom metadata in AWS Glue Data Catalog allows users to define additional attributes and tags for their data assets. This can include business metadata, such as data ownership, data lineage, and data classification, as well as technical metadata, such as data format, compression type, and encryption settings. Custom metadata enhances data governance, data lineage tracking, and data discovery capabilities.

Glue Streaming ETL

AWS Glue Streaming ETL is a feature of AWS Glue that enables real-time data processing and analytics by continuously ingesting and transforming streaming data. It

allows you to build scalable, serverless streaming data pipelines for processing data in motion.

Q. How does AWS Glue Streaming ETL differ from batch ETL?

AWS Glue Streaming ETL processes data in real-time as it arrives, while batch ETL processes data in discrete batches. With streaming ETL, data is processed as soon as it becomes available, enabling low-latency processing and real-time analytics. Batch ETL, on the other hand, processes data in fixed-size batches, which may introduce additional latency.

Q. What are some use cases for AWS Glue Streaming ETL?

Some use cases for AWS Glue Streaming ETL include:

- Real-time analytics and monitoring
- Fraud detection and anomaly detection
- Real-time recommendation engines
- Clickstream analysis and user behavior tracking
- IoT data processing and analysis

Q. How does AWS Glue Streaming ETL handle fault tolerance and scalability?

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar - LinkedIn](#)

AWS Glue Streaming ETL automatically handles fault tolerance and scalability by distributing data processing across multiple worker nodes. It monitors the health and performance of worker nodes and automatically scales the resources up or down based on the workload. In case of failures, it retries failed tasks and ensures that data processing continues without interruption.

Q. Can you explain the architecture of AWS Glue Streaming ETL?

The architecture of AWS Glue Streaming ETL consists of the following components:

- Streaming Sources: Data streams from sources such as Amazon Kinesis Data Streams, Apache Kafka, or Amazon Managed Streaming for Apache Kafka (Amazon MSK).
- AWS Glue Streaming ETL Job: A Glue job that ingests streaming data, applies transformations, and outputs the results to a destination.
- Worker Nodes: EC2 instances provisioned by AWS Glue to process data in parallel.
- Data Destination: The target data store where transformed data is stored, such as Amazon S3, Amazon Redshift, or Amazon DynamoDB.

Q. How does AWS Glue Streaming ETL integrate with other AWS services?

AWS Glue Streaming ETL integrates seamlessly with other AWS services to build end-to-end streaming data pipelines. For example:

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar -| LinkedIn](#)

- It can ingest data from Amazon Kinesis Data Streams or Apache Kafka.
- It can output transformed data to Amazon S3, Amazon Redshift, or Amazon DynamoDB.
- It can trigger AWS Lambda functions or AWS Step Functions based on specific events or conditions.

Q. What are some best practices for using AWS Glue Streaming ETL?

Some best practices for using AWS Glue Streaming ETL include:

- Designing fault-tolerant and scalable data pipelines.
- Optimizing data partitioning and compression for efficient processing.
- Monitoring job performance and resource utilization.
- Using schema validation and data validation to ensure data quality.
- Leveraging encryption and access controls to secure sensitive data.

Q. How does AWS Glue Streaming ETL handle late-arriving data?

AWS Glue Streaming ETL provides built-in support for handling late-arriving data. It maintains a watermark to track the progress of data ingestion and processing. If data arrives late, it is processed with a timestamp that falls within the watermark window. This ensures that late-arriving data is correctly processed and integrated into the output.

AWS Lambda

AWS Lambda is a serverless computing service provided by Amazon Web Services (AWS) that allows you to run code without provisioning or managing servers. You can upload your code to Lambda and AWS takes care of automatically scaling and managing the underlying infrastructure required to run your code in response to events.

Q. How does AWS Lambda work? AWS Lambda works by executing code in response to events triggered by other AWS services or custom events. You upload your code to Lambda and define the event sources (e.g., Amazon S3, Amazon DynamoDB, Amazon Kinesis) that trigger the execution of your code. When an event occurs, Lambda automatically provisions the necessary compute resources, runs your code, and then scales down to zero when the code has finished executing.

Q. What are the benefits of using AWS Lambda?

Some benefits of using AWS Lambda include:

- No server management: Lambda automatically scales and manages the underlying infrastructure for you.
- Cost-effective: You only pay for the compute time used by your code.
- Scalable: Lambda can scale automatically to handle high volumes of requests.
- Fully managed: AWS takes care of patching, monitoring, and maintaining the infrastructure.
- Integrated with AWS services: Lambda integrates seamlessly with other AWS services, enabling you to build serverless architectures easily.

Q. What languages are supported by AWS Lambda?

AWS Lambda supports several programming languages, including:

- Node.js
- Python
- Java
- Go

- .NET Core (C)

Q. What is a Lambda function?

A Lambda function is the code that you upload to AWS Lambda to be executed in response to events. It contains the code logic that defines what happens when the function is invoked.

Q. How do you trigger a Lambda function?

Lambda functions can be triggered by various event sources, including:

- Changes in data stored in Amazon S3
- Updates to records in Amazon DynamoDB
- Messages published to Amazon SNS topics
- Events generated by AWS services such as AWS CloudWatch Events or AWS IoT

Q. Can you explain the concept of cold starts in AWS Lambda?

Cold start is the term used to describe the latency that occurs the first time a Lambda function is invoked or when it hasn't been invoked for a while. During a cold start, AWS Lambda needs to provision the required compute resources to run the function, which

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar - LinkedIn](#)

can result in increased latency. Subsequent invocations of the function reuse the existing resources, resulting in lower latency (known as warm starts).

Q. How can you optimize the performance of AWS Lambda functions?

Some ways to optimize the performance of Lambda functions include:

- Reducing the size of the deployment package
- Minimizing initialization code and dependencies
- Using provisioned concurrency to keep functions warm
- Implementing efficient error handling and retries
- Tuning memory allocation and function timeouts

Q. What are the limitations of AWS Lambda?

Limitations of AWS Lambda include:

- Maximum execution duration (default is 15 minutes)
- Maximum memory allocation (up to 10 GB)
- Maximum deployment package size (up to 250 MB)
- Concurrency limits (default is 1000 concurrent executions per region)
- Statelessness (Lambda functions are stateless and have no persistent storage)

Q. How can you monitor and debug AWS Lambda functions?

We can monitor and debug Lambda functions using AWS CloudWatch logs, which capture log output from your function. We can also use CloudWatch metrics to monitor function performance and errors.

Amazon Kinesis

Amazon Kinesis is a platform on AWS for real-time streaming data ingestion, processing, and analysis. It enables you to collect, process, and analyze large streams

of data in real-time from various sources such as website clickstreams, IoT devices, logs, and social media feeds.

Q. What are the main components of Amazon Kinesis? Amazon Kinesis consists of the following main components:

- Kinesis Data Streams: Allows you to build custom applications that process or analyze streaming data.
- Kinesis Data Firehose: Automatically loads streaming data into AWS data stores and analytics services for near real-time analysis.
- Kinesis Data Analytics: Allows you to process and analyze streaming data using SQL or Apache Flink.

Q. What is a Kinesis Data Stream?

A Kinesis Data Stream is a scalable and durable real-time data streaming service that enables you to continuously collect and process large streams of data records in real-time. Data records in a stream are distributed across multiple shards, allowing for parallel processing and high throughput.

Q. What is a Kinesis Data Firehose?

Kinesis Data Firehose is a fully managed service that automatically loads streaming data into AWS data stores and analytics services. It can capture, transform, and load streaming data into destinations such as Amazon S3, Amazon Redshift, Amazon Elasticsearch Service, and Splunk.

Q. What is a Kinesis Data Analytics application? Kinesis Data Analytics application is a real-time analytics application that allows you to process and analyze streaming data using SQL or Apache Flink. You can use Kinesis Data Analytics to run SQL queries, aggregate data, detect anomalies, and perform real-time analytics on streaming data.

Q. How does Amazon Kinesis ensure scalability and fault tolerance? Amazon Kinesis ensures scalability and fault tolerance by automatically distributing data records across multiple shards in a data stream. Each shard can handle a specific throughput of data records, and Kinesis dynamically scales the number of shards based on the incoming data rate. Additionally, Kinesis replicates data records across multiple availability zones to ensure durability and fault tolerance.

Q. What are the use cases for Amazon Kinesis?

Some common use cases for Amazon Kinesis include:

- Real-time analytics and dashboarding
- Clickstream analysis and user behavior tracking
- Log and event data ingestion and analysis

- Internet of Things (IoT) data processing and analytics
- Fraud detection and anomaly detection

Q. How does Kinesis Data Firehose differ from Kinesis Data Streams? Kinesis Data Streams is a low-level service that allows you to build custom applications for processing and analyzing streaming data. It provides full control over data processing and allows for custom data transformations. On the other hand, Kinesis Data Firehose is a fully managed service that simplifies the process of loading streaming data into AWS data stores and analytics services. It automates the data ingestion process and supports out-of-the-box integrations with various AWS services.

Q. How can you monitor Amazon Kinesis? We can monitor Amazon Kinesis using Amazon CloudWatch, which provides metrics, logs, and alarms for Kinesis Data Streams, Kinesis Data Firehose, and Kinesis Data Analytics. CloudWatch allows to monitor data ingestion rates, throughput, latency, and error rates for your Kinesis applications.

Q. What are the best practices for designing Amazon Kinesis applications? Some best practices for designing Amazon Kinesis applications include:

- Properly sizing shards to handle the expected data throughput.
- Implementing retries and error handling mechanisms to handle transient failures.
- Using fine-grained partition keys to evenly distribute data across shards.

- Monitoring application performance and scaling resources based on workload patterns.
- Encrypting data in transit and at rest to ensure data security and compliance.

Amazon Athena

Amazon Athena is an interactive query service provided by AWS that allows you to analyze data in Amazon S3 using standard SQL queries. It enables you to query data stored in S3 without the need for infrastructure provisioning or data movement.

Q. How does Amazon Athena work?

Amazon Athena works by executing SQL queries against data stored in Amazon S3. It uses a serverless architecture, which means that you don't need to provision or manage any infrastructure. Athena uses Presto, an open-source distributed SQL query engine, under the hood to execute SQL queries in parallel across multiple data files stored in S3.

Q. What types of data formats does Amazon Athena support? Amazon Athena supports various data formats, including:

- Apache Parquet
- Apache ORC

- Apache Avro

- JSON

- CSV

- TSV

Q. What are the main use cases for Amazon Athena? Some common use cases for Amazon Athena include:

- Ad-hoc querying and analysis of data stored in Amazon S3

- Log analysis and troubleshooting

- Data exploration and visualization

- Business intelligence (BI) and reporting

- ETL (Extract, Transform, Load) data processing

Q. How does Amazon Athena handle data partitioning?

Amazon Athena supports data partitioning, which can significantly improve query performance and reduce costs. Data partitioning involves organizing data in Amazon S3 into directories based on one or more partition keys. Athena automatically prunes partitions based on query predicates, allowing it to scan only the relevant partitions when executing queries.

Q. What are the benefits of using Amazon Athena? benefits of using Amazon Athena include:

- Serverless architecture: No need to provision or manage any infrastructure.
- Pay-per-query pricing: You only pay for the queries you run, with no upfront costs or commitments.
- Scalability: Athena can automatically scale to handle large volumes of data and concurrent queries.
- Integration: It integrates seamlessly with other AWS services, such as Amazon S3, AWS Glue, and AWS Identity and Access Management (IAM).

Q. Can you explain the difference between Amazon Athena and Amazon Redshift?

Amazon Athena and Amazon Redshift are both data query services provided by AWS, but they have different use cases and characteristics. Amazon Athena is ideal for ad-hoc querying and analysis of data stored in Amazon S3 using SQL queries, while Amazon Redshift is a fully managed data warehousing service optimized for analytics workloads, with support for large-scale data storage, data loading, and complex queries.

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar -| LinkedIn](#)

Q. How can you optimize query performance in Amazon Athena?

Ways to optimize query performance in Amazon Athena include:

- Partitioning data based on query patterns
- Using columnar data formats like Parquet or ORC
- Reducing the number of columns scanned by projecting only the necessary columns
- Filtering data early in the query using WHERE clauses
- Using the appropriate data types and data compression techniques

Q. How does Amazon Athena handle security and access control? Amazon Athena integrates with AWS Identity and Access Management (IAM) to control access to data and resources. You can use IAM policies to grant or deny access to specific Athena actions and resources based on users, groups, or roles. Athena also supports encryption of data at rest and in transit for enhanced security.

Q. What are some limitations of Amazon Athena?

Some limitations of Amazon Athena include:

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar -| LinkedIn](#)

- No support for updates or deletes: Athena is read-only and does not support modifying data.
- Limited data types: Athena supports a subset of SQL data types and functions.
- Performance variability: Query performance can be impacted by factors such as data volume, complexity of queries, and data partitioning strategies.

Amazon Redshift

Amazon Redshift is a fully managed data warehousing service provided by AWS. It allows you to analyze large datasets using standard SQL queries and provides high-performance, scalable data warehousing capabilities.

Q. How does Amazon Redshift differ from traditional relational databases? Amazon Redshift differs from traditional relational databases in several ways, including:

- Columnar storage: Redshift stores data in a columnar format, which improves query performance and reduces storage requirements.
- Massively parallel processing (MPP): Redshift distributes data and query processing across multiple nodes, enabling high concurrency and scalability.
- Optimized for analytics: Redshift is optimized for running complex analytics queries on large datasets, whereas traditional relational databases are more suited for transactional workloads.

Q. What are the main components of Amazon Redshift?

The main components of Amazon Redshift include:

- Cluster: A Redshift cluster consists of one or more compute nodes, each containing CPU, memory, and storage.
- Leader node: Manages client connections and query execution plans.
- Compute nodes: Store and process data in parallel.
- Node slices: Each compute node is divided into slices, which represent a portion of the node's CPU and storage.

Q. What types of data sources can you load into Amazon Redshift? We can load data into Amazon Redshift from various sources, including:

- Amazon S3: Load data from files stored in Amazon S3 using the COPY command.
- Amazon DynamoDB: Use the COPY command to load data from DynamoDB tables.
- Amazon EMR: Use the COPY command with the EMRFS protocol to load data from Hadoop Distributed File System (HDFS) or other EMR-compatible storage.

Q. How does Amazon Redshift handle data compression and distribution?

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar -| LinkedIn](#)

Amazon Redshift uses data compression and distribution techniques to improve query performance and reduce storage requirements. It supports multiple compression encodings and automatically selects the most appropriate encoding based on the data type and distribution. Redshift also distributes data across slices based on a distribution style (e.g., key distribution, even distribution) to optimize query execution.

Q. What are some best practices for optimizing query performance in Amazon Redshift?

Some best practices for optimizing query performance in Amazon Redshift include:

- Choosing the appropriate distribution style for tables based on query patterns.
- Analyzing query execution plans and using sort and distribution keys to minimize data movement.
- Using column compression and encoding to reduce storage requirements.
- Regularly vacuuming and analyzing tables to reclaim disk space and update statistics.
- Monitoring system performance using Amazon CloudWatch metrics and query monitoring features.

Q. How does Amazon Redshift handle concurrency and scalability?

Amazon Redshift uses a shared-nothing architecture and massively parallel processing (MPP) to handle concurrency and scalability. It can scale horizontally by adding additional compute nodes to a cluster and supports high levels of concurrency by parallelizing query execution across nodes and slices.

Q. What is the difference between dense compute and dense storage node types in Amazon Redshift?

Dense compute and dense storage are two types of node configurations available in Amazon Redshift:

- Dense compute: Optimized for compute-intensive workloads and provide high-performance query processing capabilities.
- Dense storage: Optimized for storage-intensive workloads and provide high storage capacity at a lower cost per terabyte.

Q. How does Amazon Redshift handle backups and data durability? Amazon Redshift automatically takes incremental backups of your cluster and stores them in Amazon S3. You can also manually trigger snapshots to create point-in-time backups. Redshift uses replication and checksums to ensure data durability and automatically replaces failed drives or nodes to maintain high availability.

Q. How does Amazon Redshift integrate with other AWS services?

Amazon Redshift integrates with other AWS services to enhance its capabilities, including:

- Amazon S3: Load data into Redshift from S3 and export query results to S3.
- AWS Glue: Catalog metadata and orchestrate ETL workflows for Redshift.
- Amazon EMR: Analyze data stored in Redshift using EMR clusters running Apache Spark or other big data frameworks.
- AWS Lambda: Trigger Lambda functions in response to events in Redshift, such as data loading or query completion.

Amazon RDS for PostgreSQL

Amazon RDS for PostgreSQL is a managed relational database service provided by AWS. It allows you to deploy, operate, and scale PostgreSQL databases in the cloud.

Q. How does Amazon RDS for PostgreSQL differ from Amazon Redshift?

- Unlike Amazon Redshift, which is optimized for analytics workloads, Amazon RDS for PostgreSQL is a general-purpose relational database service. It is suitable for OLTP (Online Transaction Processing) and OLAP (Online Analytical Processing) workloads.

Q. What are the key features of Amazon RDS for PostgreSQL?

- Key features of Amazon RDS for PostgreSQL include automated backups and snapshots, high availability with multi-AZ deployments, read replicas for read scaling, and support for various PostgreSQL extensions.

Q. How does data replication work in Amazon RDS for PostgreSQL?

- Amazon RDS for PostgreSQL supports synchronous and asynchronous replication. Multi-AZ deployments use synchronous replication to replicate data to standby instances in different availability zones for high availability.

Q. What are some best practices for optimizing performance in Amazon RDS for PostgreSQL?

- Best practices for optimizing performance in Amazon RDS for PostgreSQL include choosing the appropriate instance type and storage configuration, tuning PostgreSQL parameters, and monitoring database performance using CloudWatch metrics.

Amazon DynamoDB

Amazon DynamoDB is a fully managed NoSQL database service provided by AWS. It offers seamless scalability, high availability, and low-latency performance for applications requiring single-digit millisecond response times.

Q. What are the key features of Amazon DynamoDB?

Key features of Amazon DynamoDB include:

- Fully managed: AWS handles administrative tasks such as hardware provisioning, setup, configuration, monitoring, and scaling.
- Scalable: DynamoDB automatically scales to accommodate growing workloads by partitioning data across multiple servers.
- Performance: Delivers consistent single-digit millisecond latency for read and write operations.
- Flexible data model: Supports both document and key-value data models with JSON-like syntax.
- Built-in security: Offers encryption at rest and in transit, fine-grained access control, and integration with AWS Identity and Access Management (IAM).

Q. What are the different types of primary keys supported by DynamoDB?

DynamoDB supports two types of primary keys:

- Partition key (or hash key): A simple primary key composed of a single attribute. DynamoDB uses the partition key's value to determine the partition in which the item is stored.
- Composite primary key (or hash-and-range key): A primary key composed of two attributes: a partition key and a sort key. DynamoDB uses both the partition key and sort key values to determine the partition and sort order of the item.

Q. How does DynamoDB ensure scalability and high availability?

DynamoDB ensures scalability and high availability through partitioning and replication:

- Partitioning: DynamoDB automatically partitions data across multiple servers based on the partition key. This allows DynamoDB to handle large volumes of traffic by distributing the workload evenly across servers.
- Replication: DynamoDB replicates data across multiple availability zones within a region to ensure high availability and fault tolerance. Each write is synchronously replicated to at least three availability zones.

Q. What is the difference between provisioned throughput and on-demand capacity modes in DynamoDB?

[Deepa Vasanthkumar – Medium](#)

[Deepa Vasanthkumar -| LinkedIn](#)

- Provisioned throughput: In this mode, you specify the read and write capacity units (RCUs and WCUs) required for your DynamoDB table upfront. You are billed based on the provisioned capacity, regardless of the actual usage.

- On-demand capacity mode: In this mode, DynamoDB automatically scales read and write capacity based on your actual usage. You pay per request for the capacity consumed.

Q. How does DynamoDB handle consistency?

DynamoDB offers two types of consistency models:

- Eventual consistency: In this model, DynamoDB prioritizes availability over consistency. It ensures that all copies of data are eventually consistent within seconds.

- Strong consistency: In this model, DynamoDB provides immediate consistency for read operations. It guarantees that a read will return the most up-to-date data.

Q. What are DynamoDB streams? DynamoDB streams capture a time-ordered sequence of item-level modifications in a DynamoDB table. Each stream record represents a data modification event, such as an insert, update, or delete operation. DynamoDB streams can be used to trigger AWS Lambda functions, replicate data across tables or regions, and implement cross-region replication.

Q. How can you monitor and manage DynamoDB? DynamoDB provides several tools for monitoring and managing tables, including:

- AWS CloudWatch: Monitors and logs metrics such as read/write capacity utilization, error rates, and throttling events.
- AWS DynamoDB Console: Provides a graphical user interface for managing tables, monitoring performance, and configuring settings.
- AWS Command Line Interface (CLI) and SDKs: Allow programmatic access to DynamoDB for automation, scripting, and integration with other AWS services.

Q. What are some best practices for designing DynamoDB tables?

Some best practices for designing DynamoDB tables include:

- Choose the right partition key to evenly distribute data and avoid hot partitions.
- Use sparse indexes to reduce storage costs and improve query performance.
- Use the right data types to minimize storage space and optimize query performance.
- Implement DynamoDB streams for data change capture and real-time processing.

- Leverage secondary indexes and materialized views to support different access patterns.

Q. How does DynamoDB encryption work? DynamoDB offers encryption at rest to protect your data stored in tables. It uses AWS Key Management Service (KMS) to manage encryption keys. When you enable encryption for a DynamoDB table, all data stored in the table, including backups and replicas, is encrypted using an AWS-managed encryption key or a customer-managed encryption key.

AWS Step functions

AWS Step Functions is a fully managed serverless orchestration service provided by AWS. It allows you to coordinate and orchestrate multiple AWS services and Lambda functions into serverless workflows using visual workflows.

Q. How does AWS Step Functions work? AWS Step Functions works by defining state machines that represent the workflow logic using JSON-based Amazon States Language. Each state in the state machine represents a specific task or action to be executed. Step Functions manages the execution and coordination of tasks, handles retries and error handling, and provides visibility into the state of the workflow.

Q. What are the main components of AWS Step Functions? The main components of AWS Step Functions include:

- State machine: Represents the workflow logic defined using Amazon States Language.

- States: Individual tasks or actions within the state machine.
- Execution: An instance of a state machine execution.
- Input and output: Data passed between states during execution.
- Execution history: A log of state transitions and events during the execution of a state machine.

Q. What are the benefits of using AWS Step Functions?

Some benefits of using AWS Step Functions include:

- Simplicity: Allows you to build and visualize complex workflows using a graphical interface.
- Scalability: Automatically scales to handle high volumes of workflow executions.
- Reliability: Handles retries, error handling, and state management automatically.
- Integration: Integrates seamlessly with other AWS services, including Lambda, ECS, SNS, SQS, and more.
- Visibility: Provides detailed execution logs and monitoring metrics for workflows.

Q. How can you trigger an AWS Step Functions workflow?

AWS Step Functions workflows can be triggered in several ways, including:

- Direct invocation: You can invoke a state machine directly using the Step Functions API.
- Event-based invocation: You can trigger a state machine execution in response to events from other AWS services, such as S3 events, SNS notifications, or CloudWatch Events.
- Scheduled invocation: You can schedule state machine executions to run at specified intervals using CloudWatch Events.

Q. What are the different types of states supported by AWS Step Functions? AWS Step Functions supports several types of states, including:

- Task states: Represent a single unit of work to be performed, such as invoking a Lambda function or calling an AWS service API.
- Choice states: Define conditional branching logic based on the result of a previous state.
- Parallel states: Enable concurrent execution of multiple branches of execution.

- Wait states: Pause the execution of the state machine for a specified duration or until a specific event occurs.
- Pass states: Simply pass input to output without performing any work.

Q. How does error handling work in AWS Step Functions?

AWS Step Functions provides built-in error handling capabilities, including retries and catch clauses. You can specify retry policies for individual states or the entire state machine to handle transient errors. Catch clauses allow you to define error handling logic and transitions to handle specific error conditions.

Q. Can you integrate AWS Step Functions with other AWS services? Yes, AWS Step Functions integrates seamlessly with other AWS services, including:

- AWS Lambda: Invoke Lambda functions as tasks within Step Functions workflows.
- Amazon ECS: Coordinate container-based workflows using ECS tasks.
- Amazon SNS and Amazon SQS: Send notifications and messages between states in a state machine.
- AWS Glue: Orchestrate ETL workflows using Glue jobs within Step Functions.

Q. How does AWS Step Functions handle state transitions and retries? AWS Step Functions manages state transitions and retries automatically based on the defined workflow logic. If a state fails or encounters an error, Step Functions retries the state based on the configured retry policy. If the maximum number of retries is exceeded, Step Functions transitions to an error state or a catch clause based on the error handling logic defined in the state machine.

Q. What are some use cases for AWS Step Functions? Some common use cases for AWS Step Functions include:

- Workflow automation: Orchestrating multi-step business processes and workflows.
- Microservices orchestration: Coordinating interactions between microservices in distributed systems.
- Data processing pipelines: Managing complex ETL (Extract, Transform, Load) workflows for data processing and analysis.
- Application workflows: Implementing stateful workflows for application logic and business processes.

AWS Identity and Access Management (IAM)

AWS Identity and Access Management (IAM) is a web service provided by AWS that enables you to securely control access to AWS services and resources. IAM allows you to manage users, groups, roles, and permissions to securely delegate access to AWS resources.

Q. What are the main components of AWS IAM

The main components of AWS IAM include:

- Users: Represent individuals or entities that interact with AWS services and resources.
- Groups: Collections of users with similar permissions, making it easier to manage access.
- Roles: Define a set of permissions that can be assumed by users or AWS services.
- Policies: Documents that define permissions and are attached to users, groups, or roles.
- Access keys: Used for programmatic access to AWS services through APIs.

Q. What is an IAM policy?

An IAM policy is a document that defines permissions for actions and resources in AWS. Policies are written in JSON format and can be attached to IAM identities (users, groups, roles) or resources to specify who has access to what. Policies can grant or deny permissions, and they are evaluated when an IAM identity attempts to access a resource.

Q. What is the principle of least privilege?

The principle of least privilege is a security best practice that states that users, groups, and roles should only be granted the minimum level of access necessary to perform their tasks. By following this principle, you reduce the risk of unauthorized access and potential security vulnerabilities.

Q. How do you grant access to an AWS resource using IAM? To grant access to an AWS resource using IAM, you create an IAM policy that specifies the desired permissions and attach the policy to the IAM identity (user, group, or role) that needs access to the resource. The policy defines which actions are allowed or denied and which resources the actions can be performed on.

Q. What is IAM role chaining? IAM role chaining is the process of assuming multiple IAM roles in a sequence to escalate privileges or access resources across multiple AWS accounts. It involves granting a role permission to assume another role, allowing the initial role to temporarily assume the permissions of the second role.

Q. How do you secure access to IAM resources? To secure access to IAM resources, you can implement the following best practices:

- Enable multi-factor authentication (MFA) for IAM users with elevated privileges.
- Regularly review and audit IAM policies and permissions to ensure they align with security requirements.
- Use IAM roles and temporary security credentials instead of long-term access keys whenever possible.
- Enable AWS CloudTrail to monitor and log API activity for IAM actions.
- Apply IAM password policies to enforce strong password requirements for IAM users.

Q. What is IAM Federation? IAM Federation is the process of linking your existing identity management system with AWS IAM to allow users to access AWS resources using their existing corporate credentials. This can be achieved through integration with identity providers (IdPs) such as Active Directory, LDAP, or SAML-based identity providers.

Q. How do you rotate IAM access keys? To rotate IAM access keys, you can follow these steps:

- ☐ Generate a new access key for the IAM user.
- ☐ Update any applications or scripts that use the old access key with the new access key.
- ☐ Test the updated applications or scripts to ensure they work correctly with the new access key.
- ☐ Delete the old access key from the IAM user once you confirm that the new access key is working properly.

Q. What is IAM policy evaluation logic? IAM policy evaluation logic follows a deny-by-default model, where all requests are denied by default unless explicitly allowed by an attached policy. When a request is made to AWS, IAM evaluates all policies attached to the identity and resource involved in the request. If any policy explicitly denies the requested action, the request is denied. Otherwise, if all policies either explicitly allow or do not mention the action, the request is allowed.

Amazon Simple Notification Service (SNS)

Amazon Simple Notification Service (SNS) is a fully managed messaging service provided by AWS. It enables you to send messages or notifications to distributed systems, mobile devices, or other AWS services in a highly reliable and scalable manner.

Q. How does Amazon SNS work? Amazon SNS works by allowing you to create topics, to which subscribers can subscribe. When a message is published to a topic, Amazon SNS delivers copies of the message to each subscribed endpoint, such as HTTP/S, email, SMS, Lambda, SQS, or mobile push notification.

Q. What are the main components of Amazon SNS? The main components of Amazon SNS include:

- Topics: Logical channels for publishing messages.
- Subscriptions: Endpoints that receive messages published to a topic.
- Messages: The content sent from publishers to subscribers via topics.
- Publishers: Entities that send messages to SNS topics.
- Subscribers: Entities that receive messages published to SNS topics.

Q. What are the different delivery protocols supported by Amazon SNS?

Amazon SNS supports various delivery protocols, including:

- HTTP/S
- Email
- SMS

- Mobile push notifications (Apple Push Notification Service, Google Cloud Messaging, etc.)

- Amazon SQS

- AWS Lambda

Q. How does message filtering work in Amazon SNS? Message filtering in Amazon SNS allows subscribers to receive only the messages that are of interest to them. You can set filter policies on subscriptions to specify which messages are delivered based on message attributes or message structure.

Q. What are some use cases for Amazon SNS? Some common use cases for Amazon SNS include:

- Push notifications to mobile devices or web browsers.

- Event-driven communication between microservices in distributed systems.

- Sending alerts and notifications for system monitoring and alarms.

- Decoupling communication between application components using publish/subscribe messaging.

Q. How does message ordering work in Amazon SNS? By default, Amazon SNS does not guarantee the order in which messages are delivered. However, you can use message attributes or sequencing information to implement message ordering and ensure that messages are processed in the desired order by subscribers.

Q. What are message attributes in Amazon SNS? Message attributes are key-value pairs associated with messages published to Amazon SNS topics. They allow you to provide additional metadata or context to messages and enable advanced message filtering and processing based on attribute values.

Q. How does message delivery retry and error handling work in Amazon SNS?

Amazon SNS automatically retries message delivery for transient failures, such as network errors or service throttling. If delivery attempts fail repeatedly, Amazon SNS sends the message to a dead-letter queue (DLQ) or triggers a delivery failure notification to the publisher.

Q. How can you secure access to Amazon SNS? We can secure access to Amazon SNS by using AWS Identity and Access Management (IAM) to control user permissions, encrypting messages in transit using HTTPS, and enabling server-side encryption for messages stored in SNS topics.

Amazon Simple Queue Service (SQS)

Amazon Simple Queue Service (SQS) is a fully managed message queuing service provided by AWS. It enables you to decouple and scale microservices, distributed

systems, and serverless applications by reliably passing messages between components.

Q. How does Amazon SQS work? Amazon SQS works by allowing you to create message queues, where messages are temporarily stored until they are processed by a consumer. Producers send messages to the queue, and consumers retrieve messages from the queue to process them asynchronously.

Q. What are the main components of Amazon SQS?

The main components of Amazon SQS include:

- Message: The content sent from producers to consumers via queues.
- Queue: A container that holds messages until they are processed by consumers.
- Producer: Entities that send messages to SQS queues.
- Consumer: Entities that retrieve messages from SQS queues for processing.

Q. What are the types of queues supported by Amazon SQS?

Amazon SQS supports two types of queues:

- Standard queues: Provide best-effort ordering and at-least-once delivery of messages. They offer nearly unlimited throughput and can scale horizontally to handle high volumes of messages.
- FIFO (First-In-First-Out) queues: Guarantee that messages are delivered exactly once and in the order in which they are sent. They are suitable for applications that require strict message ordering and deduplication.

Q. How does message visibility timeout work in Amazon SQS? Message visibility timeout in Amazon SQS determines how long a message remains invisible to other consumers after being retrieved by a consumer for processing. If the message processing exceeds the visibility timeout, the message becomes visible again in the queue and can be retried by other consumers.

Q. What is the purpose of dead-letter queues (DLQs) in Amazon SQS?

Dead-letter queues (DLQs) in Amazon SQS are used to capture and store messages that cannot be processed successfully after a certain number of retries. DLQs enable you to isolate and investigate failed messages separately from the main processing flow.

Q. How does message deduplication work in Amazon SQS FIFO queues?

Message deduplication in Amazon SQS FIFO queues prevents duplicate messages from being sent or processed. When a message with a specific message deduplication ID is sent to a FIFO queue within a specified deduplication interval, SQS checks if a message with the same ID has been sent recently. If a matching message is found, the new message is discarded.

Q. What are some use cases for Amazon SQS?

Some common use cases for Amazon SQS include:

- Decoupling components in distributed systems to improve scalability and reliability.
- Implementing message-driven architectures for event-driven processing and asynchronous communication.
- Offloading time-consuming or background tasks from synchronous applications to asynchronous queues.
- Handling bursts of traffic and load spikes in web applications and microservices.

Q. How does Amazon SQS handle message delivery retries and error handling?

Amazon SQS automatically retries message delivery for transient failures, such as network errors or service throttling. If message processing fails repeatedly, SQS moves the message to a dead-letter queue (DLQ) or triggers a visibility timeout to allow for retry attempts by other consumers.

Q. How can you secure access to Amazon SQS? You can secure access to Amazon SQS by using AWS Identity and Access Management (IAM) to control user permissions, encrypting messages in transit using HTTPS, and restricting access to queues using resource-based policies.

Essential Amazon Web Services (AWS) Tools used by Data Engineers

(ALL SERIES IN SINGLE DOCUMENT)



OVERVIEW

👤 Data engineers have a vast toolkit when working with AWS, covering all aspects of data ingestion, processing, storage, and analysis. Here's a breakdown of some key AWS services data engineers rely on:

→ Data Ingestion & Management:

- 🔗 [Amazon S3](#): Object storage for raw data, scalable and cost-effective.
- 🔗 [Amazon Kinesis](#): Streaming platform for real-time data ingestion.
- 🔗 [Amazon SQS](#): Queueing service for buffering and managing data flow.
- 🔗 [AWS Data Pipeline](#): Orchestration service for building and managing data pipelines.
- 🔗 [AWS Glue](#): Serverless ETL service for data extraction, transformation, and loading.

→ Data Storage & Processing:

- 🔗 [Amazon Redshift](#): Scalable data warehouse for large-scale data analysis.
- 🔗 [Amazon DynamoDB](#): NoSQL database for fast and flexible data storage.
- 🔗 [Amazon Elasticsearch Service](#): Highly scalable search and analytics platform.
- 🔗 [Amazon EMR](#): Hadoop and Spark cluster service for big data processing.
- 🔗 [Amazon Lambda](#): Serverless compute service for running code without provisioning or managing servers.
- 🔗 [Amazon SageMaker](#): Machine learning platform for building, training, and deploying models.

→ Data Analysis & Visualization:

- 🔗 [Amazon Athena](#): Serverless interactive query service for analyzing data in S3.
- 🔗 [Amazon QuickSight](#): Cloud-based business intelligence service for data visualization and reporting.
- 🔗 [Amazon OpenSearch Service](#): Open-source search and analytics platform with near real-time capabilities.
- 🔗 [Amazon CloudWatch](#): Monitoring and observability service for tracking resource utilization and performance.

→ Additional Services:

- 🔗 [IAM](#): Identity and access management for controlling access to AWS resources.
- 🔗 [CloudFormation](#): Infrastructure as code service for automating the provisioning and management of AWS resources.
- 🔗 [CloudTrail](#): Logging service for tracking API calls made to AWS services.

→ Bonus Tip: Remember, the specific services used will depend on your specific data engineering

needs and project requirements. Don't hesitate to explore and experiment with different services to find the best fit for your workflow.

Amazon S3

📁 Amazon S3 (Simple Storage Service) is a scalable and widely used object storage service provided by Amazon Web Services (AWS). It allows data engineers to store and retrieve any amount of data at any time, making it a popular choice for building data lakes, storing backups, hosting static websites, and supporting various data-intensive applications.

→ How Data Engineers Use S3:

📁 Data Lake: S3 serves as a central data lake, a repository for all the raw data ingested from various sources. It's the first stop for most data pipelines, where data engineers can stage, organize, and pre-process it before feeding it into downstream applications.

📁 Backup and Archiving: S3 offers cost-effective storage for backups and archives of valuable data. Its durability and scalability make it perfect for long-term data retention, ensuring your information is safe and accessible even years later.

📁 Static Content Hosting: S3 can host static websites and web applications with ease. With its high availability and global distribution, it delivers content quickly and reliably to users worldwide.

📁 Data Sharing and Collaboration: S3 simplifies data sharing within your team or with external collaborators. You can set granular access permissions for different users and groups, ensuring secure and controlled data access.

📁 Big Data Processing: S3 seamlessly integrates with various big data processing frameworks like [Hadoop](#) and [Spark](#). Data engineers can directly read and write data from S3 within these platforms, eliminating the need for separate data transfer processes.

→ Examples:

📁 Example 1: A data engineer working on a customer recommendation system stores all user interaction data (clicks, purchases, etc.) in S3. This data lake serves as the source for training and running machine learning models that personalize recommendations for each user.

📁 Example 2: A team developing a real-time analytics platform uses S3 as a buffer for streaming data from sensors and devices. By pre-processing and storing the data in S3, they can enable near real-time analysis and reporting without overwhelming their processing systems.

📖 Example 3: A company needs to archive old financial records for regulatory compliance. S3's cost-effective storage and long-term data retention capabilities make it the perfect solution for storing and accessing these files without breaking the bank.



Amazon Kinesis

📖 Kinesis is a managed streaming service that handles the high-velocity ingestion and processing of data in real-time. It can handle data from various sources, like sensor readings, social media feeds, clickstream data, and more. Unlike traditional batch processing, Kinesis doesn't wait for data to accumulate – it analyzes it as it arrives, enabling immediate reactions and decisions.

→ How Data Engineers Use Kinesis:

📖 Fraud Detection: By analyzing real-time transaction streams, data engineers can identify suspicious activity and flag potential fraud attempts as they happen.

📖 IoT Platform Analytics: Kinesis streams sensor data from connected devices, allowing engineers to monitor performance, predict maintenance needs, and optimize operations in real-time.

📖 Personalization and Recommendations: Analyzing real-time user behavior through website clicks, app interactions, and purchase logs, Kinesis helps personalize content and recommendations, improving user engagement and conversion rates.

📖 Live Dashboards and Reporting: By processing data streams in real-time, Kinesis enables data engineers to build dynamic dashboards and reports that reflect the latest information, providing insightful snapshots into current trends and activities.

📖 Event-Driven Architectures: Kinesis integrates seamlessly with other AWS services and triggers actions based on real-time data events. For example, a new social media comment might trigger a sentiment analysis or automated customer outreach.

→ Examples:

📖 Example 1: A stock trading platform uses Kinesis to stream real-time market data feeds. Data engineers can analyze these streams to identify price fluctuations, trigger trading alerts, and inform investment decisions instantly.

📁 Example 2: A fitness wearable app utilizes Kinesis to capture user activity data in real-time. This data feeds personalized workout recommendations, tracks progress towards goals, and provides instant feedback during exercise sessions.

📁 Example 3: A news aggregator platform employs Kinesis to collect and analyze news articles as they are published. This enables the platform to present users with curated and trending news content based on their real-time interests.



Amazon SQS

📁 SQS is a managed message queuing service that allows applications to communicate asynchronously. Instead of directly sending data from one application to another, you can send it to an SQS queue. This acts as a temporary buffer, decoupling the sender and receiver. The receiving application can then pull messages from the queue at its own pace, ensuring smooth processing even if it's slower than the sender.

→ How Data Engineers Use SQS:

📁 Decoupling Applications: SQS breaks the tight coupling between data producers and consumers, ensuring each component operates independently and scales easily. This improves overall system resilience and flexibility.

📁 Buffering Spikes: When data arrives in bursts, SQS acts as a buffer, preventing downstream systems from getting overwhelmed. It queues messages efficiently until the receiving application is ready to process them, smoothing out data flow and preventing bottlenecks.

📁 Micro-tasking and Workflows: SQS enables breaking down large tasks into smaller, independent messages. These messages can then be processed by multiple workers in parallel, accelerating overall processing and improving efficiency.

📁 Event-Driven Architectures: SQS integrates seamlessly with other AWS services and triggers actions based on messages in the queue. This allows building flexible and reactive applications that respond to events in real-time.

📁 Retry and Error Handling: SQS offers built-in mechanisms for retrying failed messages and handling errors gracefully. This ensures robust data processing and prevents data loss even in case of temporary failures.

→ Examples:

📁 Example 1: A data pipeline ingests large datasets from various sources. Instead of directly pushing data to a data warehouse, the pipeline sends messages to an SQS queue. This allows the warehouse to process data at its own pace without slowing down the ingestion process.

📁 Example 2: A video encoding service receives video files for processing. These files are uploaded to S3 and a message is sent to an SQS queue. Worker applications read the queue and encode the videos independently, scaling automatically based on the workload.

📁 Example 3: A customer support system uses SQS to handle high volumes of incoming support tickets. Tickets are added to an SQS queue, and support agents can pick them up at their own pace, ensuring efficient ticket management even during peak times.



AWS Data Pipeline

📁 Data Pipeline is a managed service that allows you to define and automate data transformations and movements within AWS. You can easily build data pipelines that extract, transform, and load (ETL) or extract, load, and transform (ELT) data from various sources like databases, files, and AWS services, delivering it to your desired destination, like data warehouses, analytics platforms, or applications.

→ How Data Engineers Use AWS Data Pipeline:

📁 Automating Data Flow: Instead of manually moving data, data engineers can build pipelines that run automatically with scheduled triggers or event-driven responses, ensuring consistent and reliable data processing.

📁 Orchestrating Complex Workflows: Data Pipeline allows chaining multiple data processing activities together, from simple data extraction to complex transformations and loading, creating multi-step data journeys with ease.

📁 Handling Different Data Formats: Data Pipeline supports various data formats like CSV, JSON,

XML, and more, making it versatile for integrating data from diverse sources.

📁 **Error Handling and Logging:** Data Pipeline offers built-in mechanisms for handling errors, retrying failed tasks, and logging activities, ensuring transparency and smooth pipeline execution.

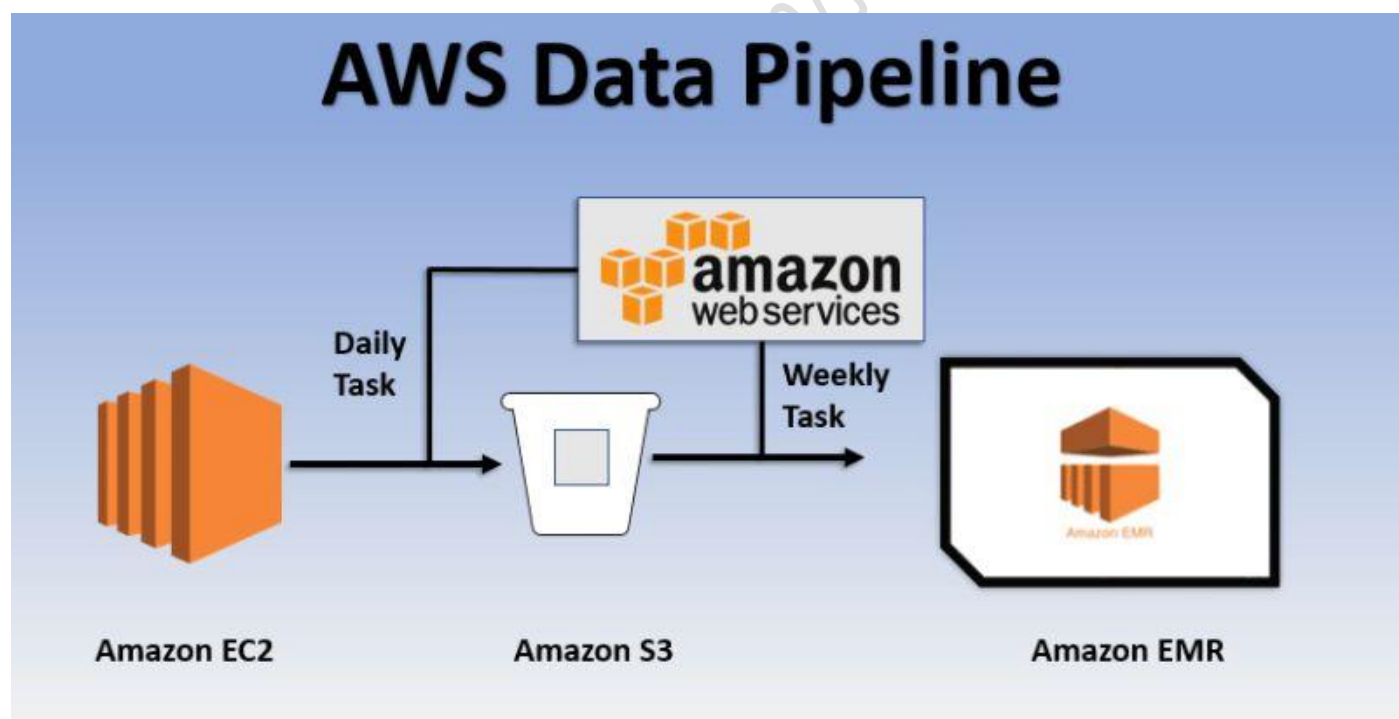
📁 **Cost-effective Orchestration:** You only pay for the resources your pipelines use, making it a cost-efficient solution for running data workflows of any size.

→ Examples:

📁 **Example 1:** A data engineer builds a pipeline that extracts sales data from different stores every hour, transforms it into a standardized format, and loads it into a Redshift data warehouse for daily sales analysis.

📁 **Example 2:** A pipeline automatically ingests sensor data from IoT devices in real-time, performs basic filtering and aggregation, and feeds it into a machine learning model for predictive maintenance insights.

📁 **Example 3:** A data pipeline retrieves customer reviews from different platforms, analyzes sentiment using a natural language processing service, and sends summarized feedback reports to marketing and product teams.



AWS Glue

📁 **Glue** simplifies data integration by automatically discovering and cataloging data sources within your AWS environment. It connects to various databases, data lakes (like [S3](#)), and other AWS services, providing a centralized view of your data assets. Additionally, Glue offers ETL (Extract, Transform, Load) capabilities, letting you cleanse, transform, and prepare your data for downstream applications.

→ How Data Engineers Use AWS Glue:

📁 Data Discovery and Cataloging: Glue scans and catalogs data sources, automatically building a metadata repository that describes your data assets, their structure, and location. This makes it easier to find, understand, and utilize your data.

📁 ETL Job Creation and Management: Build and schedule ETL jobs for transforming and preparing your data for analysis. Glue offers a visual interface for constructing data pipelines, simplifying even complex ETL workflows.

📁 Prebuilt Transformers and Scripts: Glue provides a library of prebuilt transformers and user-defined scripts for common data manipulations. This cuts down development time and ensures consistent data preparation across your projects.

📁 Scalable and Serverless: Glue runs on a serverless architecture, eliminating infrastructure management and scaling transparently based on your data processing needs. You only pay for the resources you use, making it cost-effective for any project size.

📁 Real-time and Batch Processing: Glue supports both real-time and batch processing, allowing you to handle both streaming and static data sources efficiently.

→ Examples:

📁 Example 1: A data engineer uses Glue to catalog data from various customer databases and S3 buckets. They then build ETL jobs to cleanse and format the data, preparing it for customer segmentation analysis in Amazon Redshift.

📁 Example 2: A team analyzes web logs using Glue. They leverage the service to discover log data in S3, transform it into structured format, and feed it into a machine learning model for predicting user behavior and enhancing website personalization.

📁 Example 3: A company needs to integrate real-time sensor data from IoT devices with historical sensor data stored in S3. Glue's real-time capabilities enable seamless data ingestion and transformation, allowing for continuous monitoring and predictive maintenance insights.



Amazon Redshift

🔗 Redshift is a fully managed data warehouse designed for large-scale data analysis. It excels at storing and analyzing massive datasets from various sources, like relational databases, data lakes, and application logs. Unlike traditional databases, Redshift utilizes a parallel processing architecture, meaning it distributes your data across multiple nodes, enabling lightning-fast queries and analytical tasks even on datasets spanning billions of rows.

→ How Data Engineers Use Amazon Redshift:

🔗 Large-scale Data Analysis: Redshift empowers data engineers to analyze massive datasets for trends, patterns, and correlations that might be invisible in smaller samples. This unlocks opportunities for data-driven decision making, business optimization, and scientific discovery.

🔗 Building Data Marts: Redshift allows creating smaller, targeted data marts from the main data warehouse, focusing on specific business functions or departments. This improves query performance and data accessibility for various teams within an organization.

🔗 Data Visualization and Reporting: Redshift integrates seamlessly with popular business intelligence tools, enabling data engineers to create interactive dashboards and reports for visualizing complex data insights in a user-friendly manner.

🔗 Machine Learning Integration: Redshift plays a crucial role in preparing data for machine learning models. By pre-processing, cleansing, and joining data in Redshift, data engineers provide models with high-quality fuel for accurate predictions and intelligent applications.

🔗 Cost-effective Scalability: Redshift scales effortlessly with your data needs, allowing you to pay only for the resources you use. This makes it cost-effective for both small and large-scale data analytics projects.

→ Examples:

🔗 Example 1: A retail company leverages Redshift to analyze millions of customer transactions. They perform complex queries to identify buying patterns, predict customer churn, and personalize marketing campaigns for increased sales.

🔗 Example 2: A biopharmaceutical company uses Redshift to analyze clinical trial data from thousands of patients. By querying vast datasets, researchers discover potential drug interactions, identify promising treatment candidates, and accelerate medical advancements.

🔗 Example 3: A financial institution utilizes Redshift to combat fraud. They analyze real-time transaction data to detect suspicious activity, prevent financial losses, and ensure customer security.



amazon
REDSHIFT

Amazon DynamoDB

🔗 DynamoDB is a fully managed [NoSQL](#) database that stores data in key-value pairs. Unlike traditional relational databases, it doesn't rely on schema and can handle diverse data structures effortlessly. This flexibility makes it ideal for applications requiring high availability, fast performance, and the ability to scale instantly with changing data demands.

→ How Data Engineers Use [Amazon](#) DynamoDB:

🔗 Building Scalable Applications: DynamoDB's inherent scalability allows data engineers to build applications that seamlessly handle increasing data volumes without infrastructure management or performance bottlenecks.

🔗 Mobile and IoT Applications: Its fast response times and flexible data model make DynamoDB perfect for mobile and IoT applications, where real-time data access and dynamic schema adaptations are crucial.

🔗 Session Management and Caching: Data engineers leverage DynamoDB to store user sessions, application configurations, and frequently accessed data as a high-performance cache, ensuring smooth user experiences and efficient data access.

🔗 NoSQL Data Storage for Unstructured Data: When dealing with unstructured data like sensor readings, social media posts, or logs, DynamoDB's flexible structure eliminates schema limitations and simplifies data storage and retrieval.

🔗 Global Scale and Geographically Distributed Data: DynamoDB seamlessly replicates data across multiple AWS regions, ensuring high availability and low latency for globally distributed applications and geographically dispersed data access needs.

→ Examples:

🔗 Example 1: A social media platform utilizes DynamoDB to store user profiles, posts, and connections. Its scalability handles the ever-growing user base, and its flexibility accommodates evolving data structures as new features are introduced.

🔗 Example 2: A ride-hailing app uses DynamoDB to store driver and passenger locations, booking information, and real-time traffic data. Its fast response times ensure quick ride matching and

efficient route planning, while its scalability adapts to peak demand periods.

📖 Example 3: A fitness tracker app stores sensor data from wearables in DynamoDB. The flexible data model adapts to different sensor types, and the high availability ensures continuous data access for users to track their health and fitness goals.



Amazon Elasticsearch Service

📖 ES is a managed service based on the open-source Elasticsearch search engine. It allows you to store, search, and analyze large volumes of text and semi-structured data with incredible speed and accuracy. Its capabilities extend beyond simple keyword searches, offering advanced features like faceted navigation, geospatial search, and real-time analytics, making it a versatile tool for diverse data exploration and analysis needs.

→ How Data Engineers Use Amazon Elasticsearch Service:

📖 Building Search-based Applications: Integrate ES into your applications to offer users intuitive and powerful search experiences. From e-commerce sites to internal documentation platforms, ES delivers lightning-fast and relevant results, enhancing user interactions and information discovery.

📖 Log Analysis and Monitoring: Analyze application logs, security logs, and website traffic data in real-time using ES. Its ability to ingest and analyze high volumes of data makes it perfect for identifying trends, detecting anomalies, and troubleshooting issues in your systems.

📖 Content and Media Search: Power websites and applications with intelligent search for text, images, videos, and other multimedia content. ES's faceted search and relevance ranking features ensure users find exactly what they need, enhancing engagement and content discoverability.

📖 Personalization and Recommendation Engines: Implement AI-powered personalization and recommendation systems based on user search history, preferences, and past interactions. ES's analytics capabilities allow you to extract insights from user data and deliver personalized experiences that drive engagement and conversions.

→ Examples:

🔗 Example 1: A travel booking platform leverages ES to provide users with instant search results for hotels, flights, and experiences based on various criteria like location, budget, and preferences. The platform uses ES's faceted search to refine results and personalize recommendations, improving user experience and conversion rates.

🔗 Example 2: A news website integrates ES to power its search engine, allowing users to quickly find relevant articles based on keywords, topics, and even entities mentioned in the text. ES's near-real-time indexing ensures users access the freshest information with lightning-fast search speeds.



Amazon EMR

🔗 EMR stands for Elastic MapReduce, a managed cluster service for processing big data on AWS. It simplifies running popular open-source big data frameworks like [Apache Spark](#), Apache [Hadoop](#), and Apache [Hive](#), eliminating the need for manual infrastructure management and configuration. By providing pre-configured clusters and scaling them automatically, EMR empowers data engineers to focus on writing and executing their big data processing logic efficiently.

→ How Data Engineers Use Amazon EMR:

🔗 Large-scale Data Processing: Analyze petabytes of data efficiently using various frameworks like Spark for real-time analytics, Hadoop for batch processing, and Hive for data warehousing needs. EMR allows them to leverage the processing power of distributed clusters without infrastructure complexities.

🔗 ETL Workflows: Build and automate complex Extract, Transform, and Load (ETL) pipelines that ingest data from various sources, clean and transform it, and load it into data warehouses or analytics platforms for further analysis.

Machine Learning Model Training: Train large-scale machine learning models on vast datasets using frameworks like Spark MLlib or [TensorFlow](#) on EMR clusters, accelerating model development and deployment.

📁 Log Analysis and Security: Analyze massive log files and security data in real-time using EMR to identify trends, detect anomalies, and improve security posture.

📁 Data Lake Processing: Explore and analyze data stored in data lakes (like Amazon [S3](#)) using EMR, extracting valuable insights and transforming raw data into actionable knowledge.

→ Examples:

📁 Example 1: A company analyzes website clickstream data from millions of users using Spark on EMR. They gain insights into user behavior, personalize recommendations, and improve website conversion rates.

📁 Example 2: A research institute analyzes genomic data from thousands of patients using Hadoop on EMR. They discover genetic markers associated with diseases and accelerate medical research.

📁 Example 3: A financial services company builds ETL pipelines on EMR to ingest and process financial transactions from various sources. They gain insights into customer behavior and optimize their offerings.



Amazon Lambda

📁 Lambda lets you run code without provisioning or managing servers. Simply upload your code (functions), and Lambda takes care of everything else – allocating resources, running the code, and scaling automatically based on demand. This serverless approach means you focus on writing the code, not infrastructure complexities.

→ How Data Engineers Use Amazon Lambda:

📁 Serverless Data Processing: Build serverless data pipelines that trigger automatically based on events like new data arriving in [S3](#), changes in databases, or API calls. This eliminates the need for

dedicated servers and simplifies data processing workflows.

👉 **Microservices Architecture:** Break down complex data processing tasks into smaller, independent Lambda functions, enabling modularity, scalability, and easier code management.

👉 **Real-time Data Processing:** Process data streams in real-time using Lambda's near-instantaneous trigger capabilities. This allows for applications like fraud detection, log analysis, and real-time analytics dashboards.

👉 **API Gateway Integration:** Create serverless APIs using Lambda functions, allowing you to build RESTful APIs without managing servers or scaling infrastructure.

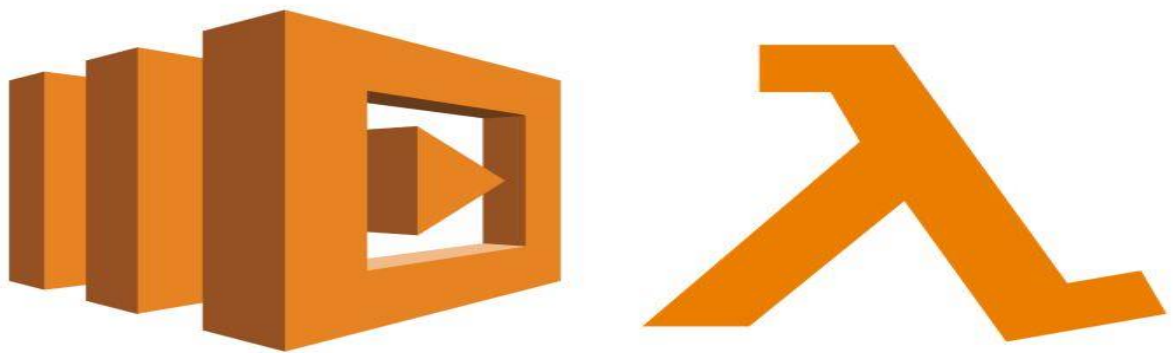
👉 **Data Transformation and Cleaning:** Cleanse, transform, and prepare data for further analysis or storage using Lambda functions triggered by data arrival events.

→ Examples:

👉 **Example 1:** A company builds a serverless [ETL](#) pipeline using Lambda functions triggered by new data uploads to S3. The functions automatically clean, transform, and load the data into a data warehouse.

👉 **Example 2:** A retail website uses Lambda functions to trigger personalized product recommendations based on user actions. This improves customer experience and increases sales.

👉 **Example 3:** A financial services company detects fraudulent transactions in real-time using Lambda functions triggered by payment events. This helps prevent financial losses and protects customers.



AWS Lambda

Amazon SageMaker

👉 SageMaker removes the complexities of machine learning development by offering a managed environment with various tools and services. You can quickly launch [Jupyter](#) notebooks for

experimentation, leverage pre-built algorithms for common tasks, or build custom models with your preferred frameworks. SageMaker handles infrastructure management, scaling resources, and automating repetitive tasks, allowing you to focus on the core ML tasks.

→ How Data Engineers Use [Amazon](#) SageMaker:

- 📁 Experimentation and Prototyping: Use Jupyter notebooks with pre-built SageMaker libraries to rapidly experiment with different algorithms and data sets, accelerating model development.
- 📁 Model Training and Tuning: Choose from a wide range of pre-built algorithms or bring your own custom models. SageMaker automates hyperparameter tuning and resource allocation, optimizing model performance.
- 📁 Model Deployment and Management: Easily deploy trained models into production with a few clicks. SageMaker manages infrastructure scaling and provides real-time monitoring for model performance.
- 📁 Machine Learning Pipelines: Build and automate end-to-end ML pipelines that involve data preparation, training, deployment, and monitoring, streamlining the ML workflow.
- 📁 Collaboration and Reproducibility: Share notebooks and models within teams, ensuring project consistency and reproducible results.

→ Examples:

- 📁 Example 1: A retail company uses SageMaker to build a recommendation engine that analyzes customer purchase history and predicts future purchases. This leads to personalized product recommendations and increased sales.
- 📁 Example 2: A financial services company trains a fraud detection model on SageMaker using historical transaction data. This helps identify fraudulent transactions in real-time and prevent financial losses.
- 📁 Example 3: A healthcare organization builds a medical imaging analysis model on SageMaker to analyze X-rays and identify potential abnormalities. This assists doctors in diagnoses and improves patient care.



Amazon SageMaker

Amazon Athena

👉 Athena is a serverless query engine that allows you to run interactive SQL queries directly on data stored in S3. It eliminates the need for complex data warehousing solutions or managing infrastructure, making it a cost-effective and scalable option for big data analysis.

→ How Data Engineers Use Amazon Athena:

👉 Exploring Data Lakes: Easily query and analyze data stored in data lakes without the need for data transformation or ETL pipelines. This allows for quick exploration and identification of trends and patterns.

👉 Ad-hoc Analysis: Empower business analysts and data scientists to perform ad-hoc analysis on large datasets without relying on data engineers to prepare the data. This fosters collaboration and accelerates decision-making.

👉 Cost-effective Analytics: Pay only for the queries you run, making Athena ideal for occasional or exploratory analysis where traditional data warehouses might be cost-prohibitive.

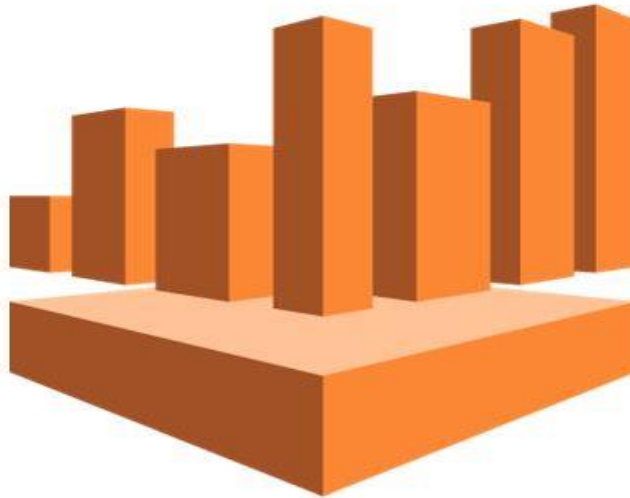
👉 Integration with Existing Tools: Seamlessly integrate Athena with other AWS services like [Amazon Redshift](#), [Amazon QuickSight](#), and Amazon EMR for a comprehensive data analytics ecosystem.

→ Examples:

👉 Example 1: A marketing team uses Athena to analyze website clickstream data stored in S3 to understand user behavior and optimize marketing campaigns.

👉 Example 2: A financial services company uses Athena to query historical transaction data to identify fraudulent activities and improve security measures.

📖 Example 3: A research institute uses Athena to analyze large genomic datasets stored in S3 to identify genetic markers associated with diseases and accelerate medical research.



Amazon Athena

Amazon QuickSight

📖 QuickSight is a serverless, cloud-based BI platform that simplifies data visualization creation and sharing. Data engineers can connect to various data sources (databases, data lakes, etc.), prepare and transform data, and build interactive dashboards with drag-and-drop simplicity. QuickSight empowers users to explore data, identify trends, and gain actionable insights through intuitive visualizations like charts, graphs, and maps.

→ How Data Engineers Use Amazon QuickSight:

📖 **Building Interactive Dashboards:** Create visually appealing and informative dashboards tailored to specific audiences (management, marketing, sales). Embed filters, drill-downs, and actions to enable deep-dive analysis within the dashboard.

📖 **Data Storytelling:** Craft compelling narratives by combining visualizations with text, images, and annotations to guide users through key insights and data exploration.

📖 **Empowering Business Users:** Facilitate self-service analytics by enabling business users to create and customize their own reports and dashboards, reducing reliance on data engineers for basic analysis.

📖 **Collaboration and Sharing:** Securely share dashboards and reports with colleagues and stakeholders, fostering data-driven decision-making across teams.

→ Examples:

📁 Example 1: A retail company builds a QuickSight dashboard that displays real-time sales trends, inventory levels, and customer demographics. This empowers store managers to optimize product placement and promotions.

📁 Example 2: A financial services company creates interactive dashboards for risk analysts to visualize market trends, identify potential risks, and make informed investment decisions.

📁 Example 3: A healthcare organization builds a QuickSight dashboard for medical professionals to analyze patient data, track treatment progress, and gain insights into patient populations.



Amazon OpenSearch Service

📁 OpenSearch Service is a fully managed offering based on the popular OpenSearch open-source search engine. It simplifies deploying, scaling, and managing OpenSearch clusters, allowing data engineers to focus on building search and analytics applications without infrastructure complexities. It offers diverse capabilities like:

Full-text Search: Locate relevant information within text-heavy documents with ease.

Structured Search: Search and filter data based on specific fields and attributes.

Geospatial Search: Find data associated with geographical locations.

Real-time Analytics: Analyze data streams and gain insights as they arrive.

Aggregation and Visualization: Summarize and visualize search results for deeper understanding.

→ How Data Engineers Use Amazon OpenSearch Service:

👉 Log Analysis: Gain insights into system logs, application logs, and security logs to identify errors, detect anomalies, and improve performance.

👉 E-commerce Search: Build powerful product search experiences for online stores, enabling customers to find what they need quickly and easily.

👉 Website Search: Implement website search functionality to improve user experience and navigation.

👉 Enterprise Search: Facilitate search across various enterprise data sources (documents, code, emails) for knowledge discovery and improved decision-making.

👉 Fraud Detection: Analyze transaction data to identify suspicious activity and prevent fraudulent transactions.

→ Examples:

👉 Example 1: A gaming company uses OpenSearch Service to analyze player logs in real-time, identifying bugs, balancing gameplay, and personalizing game experiences.

👉 Example 2: A media streaming platform utilizes OpenSearch Service to power its content search, enabling users to discover movies, shows, and music efficiently.

👉 Example 3: A research institute leverages OpenSearch Service to analyze large datasets of scientific publications, accelerating research and discovery.



Amazon OpenSearch Service

Amazon CloudWatch

👉 Amazon CloudWatch is a service used for monitoring and observing resources in real-time, built for DevOps engineers, developers, site reliability engineers (SREs), and IT managers. CloudWatch provides users with data and actionable insights to monitor their respective applications, stimulate system-wide performance changes, and optimize resource utilization. CloudWatch collects

monitoring and operational data in the form of logs, metrics, and events, providing its users with an aggregated view of AWS resources, applications, and services that run on AWS. The CloudWatch can also be used to detect anomalous behavior in the environments, set warnings and alarms, visualize logs and metrics side by side, take automated actions, and troubleshoot issues.

→ How Data Engineers Use Amazon CloudWatch:

👉 Integrations: Seamlessly integrate CloudWatch with other AWS services like [Lambda](#), [S3](#), and Step Functions to create comprehensive monitoring and observability workflows.

👉 APIs and SDKs: Leverage APIs and SDKs to programmatically interact with CloudWatch data, automating tasks and integrating it into your CI/CD pipelines.

👉 Anomaly Detection: Use machine learning-powered anomaly detection to identify unusual patterns in metrics and logs, proactively detecting potential issues before they become critical.

→ Examples:

👉 Examples 1: A data engineer sets up CloudWatch to monitor the performance of a large-scale batch processing job running on EMR. They track metrics like CPU utilization, memory usage, and job completion times to identify any bottlenecks and optimize the job configuration.

👉 Examples 2: A DevOps team uses CloudWatch to monitor the health and availability of their production web application. They receive alerts in real-time if any metrics exceed predefined thresholds, allowing them to quickly diagnose and resolve issues before they impact users.

👉 Examples 3: A security engineer uses CloudWatch to monitor logs from various AWS services for potential security threats. They set up custom filters and alerts to detect suspicious activity and take immediate action if needed.



Amazon IAM

🔑 AWS [Identity & Access Management \(IAM\)](#) and Access Management (IAM) is a web service that helps you securely control access to AWS resources. With IAM, you can centrally manage permissions that control which AWS resources users can access. You use IAM to control who is authenticated (signed in) and authorized (has permissions) to use resources.

→ How Data Engineers Use Amazon IAM:

🔑 Individual User Accounts: Create individual user accounts for each data engineer and assign specific IAM roles with appropriate permissions.

🔑 IAM Roles: Create roles for groups of users with similar access needs, simplifying permission management. Users can assume roles with the required permissions for specific tasks, enhancing flexibility and security.

🔑 Federated Access: Integrate with external identity providers ([ACTIVE](#) Directory, [Okta](#), etc.) for seamless user authentication and authorization, allowing data engineers to use their existing credentials to access AWS resources.

→ Examples:

🔑 Examples 1: A data engineer creates an IAM role with read-only access to a specific S3 bucket containing sensitive data. They assign this role to a data analyst who needs to analyze the data without granting write permissions.

🔑 Examples 2: A DevOps team creates IAM roles with specific permissions for deploying and managing different environments (development, staging, production). Data engineers can assume these roles based on the environment they need to work with.

🔑 Examples 3: A company integrates IAM with their Active Directory, allowing data engineers to use their existing credentials to access AWS resources, simplifying access management and reducing the need for separate AWS credentials.



Amazon CloudFormation

☞ Amazon Web Services (AWS) is the service offered by the AWS cloud it is mainly used to provision the service in the AWS like EC2, [S3](#), Autoscaling, load balancing and so on you can provision all the service automation with the Infrastructure as a code (IAC), instead of managing all of them manually you can manage with the help of AWS Cloudformation.

→ How Data Engineers Use Amazon CloudFormation:

☞ Provision Data Pipelines: Define your data pipelines (ETL, ELT) as CloudFormation templates, including resources like S3 buckets, [Redshift](#) clusters, and [Lambda](#) functions.

☞ Automated Data Lake Creation: Use CloudFormation to provision and configure data lakes with S3 buckets, IAM policies, and access controls, streamlining data storage and access.

☞ Versioning and Rollbacks: Easily roll back to previous deployments if issues arise, thanks to CloudFormation's versioning capabilities.

→ Examples:

☞ Examples 1: A data engineer creates a CloudFormation template to deploy an ETL pipeline consisting of an S3 bucket for data storage, a Lambda function for data transformation, and a Redshift cluster for data warehousing. This template can be reused and deployed across different environments.

☞ Examples 2: A DevOps team uses CloudFormation to provision a data lake with fine-grained access control for different user groups. The template defines S3 buckets, IAM roles, and policies, ensuring secure access to sensitive data.

☞ Examples 3: A large organization uses CloudFormation stacks to deploy their data infrastructure across multiple AWS accounts and regions, maintaining consistency and simplifying infrastructure management.



Amazon CloudTrail

☞ With AWS CloudTrail, you can monitor your AWS deployments in the cloud by getting a history of AWS API calls for your account, including API calls made by using the AWS Management Console, the AWS SDKs, the command line tools, and higher-level AWS services. You can also identify which users and accounts called AWS APIs for services that support CloudTrail, the source IP address from which the calls were made, and when the calls occurred. You can integrate CloudTrail into applications using the API, automate trail creation for your organization, check the status of your trails, and control how administrators turn CloudTrail logging on and off.

→ How Data Engineers Use Amazon CloudTrail:

☞ Detect Unauthorized Access: Identify potential security threats and suspicious activity by analyzing CloudTrail logs for unusual patterns or unauthorized API calls.

☞ Compliance Reporting: Simplify compliance audits by generating reports on user activity, resource changes, and API calls, demonstrating adherence to regulatory requirements.

☞ Forensic Analysis: Investigate security incidents or troubleshoot issues by analyzing CloudTrail logs for specific events or activities.

→ Examples:

☞ Examples 1: A data engineer sets up CloudTrail to track all API calls related to S3 buckets used for storing sensitive data. This helps them detect any unauthorized access attempts and ensure data security.

📁 Examples 2: A DevOps team uses CloudTrail logs to generate compliance reports for PCI-DSS, demonstrating their adherence to security regulations for their cloud infrastructure.

📁 Examples 3: A security analyst investigates a potential security incident by analyzing CloudTrail logs for suspicious API calls and user activity within a specific timeframe.



AWS Certified Data Engineer - Associate (DEA-C01) Exam Guide

Introduction

The AWS Certified Data Engineer - Associate (DEA-C01) exam validates a candidate's ability to implement data pipelines and to monitor, troubleshoot, and optimize cost and performance issues in accordance with best practices.

The exam also validates a candidate's ability to complete the following tasks:

- Ingest and transform data, and orchestrate data pipelines while applying programming concepts.
- Choose an optimal data store, design data models, catalog data schemas, and manage data lifecycles.
- Operationalize, maintain, and monitor data pipelines. Analyze data and ensure data quality.
- Implement appropriate authentication, authorization, data encryption, privacy, and governance. Enable logging.

Target candidate description

The target candidate should have the equivalent of 2–3 years of experience in data engineering. The target candidate should understand the effects of volume, variety, and velocity on data ingestion, transformation, modeling, security, governance, privacy, schema design, and optimal data store design. Additionally, the target candidate should have at least 1–2 years of hands-on experience with AWS services.

Recommended general IT knowledge

The target candidate should have the following general IT knowledge:

- Setup and maintenance of extract, transform, and load (ETL) pipelines from ingestion to destination
- Application of high-level but language-agnostic programming concepts as required by the pipeline
- How to use Git commands for source control
- How to use data lakes to store data
- General concepts for networking, storage, and compute

Recommended AWS knowledge

The target candidate should have the following AWS knowledge:

- How to use AWS services to accomplish the tasks listed in the Introduction section of this exam guide
- An understanding of the AWS services for encryption, governance, protection, and logging of all data that is part of data pipelines
- The ability to compare AWS services to understand the cost, performance, and functional differences between services
- How to structure SQL queries and how to run SQL queries on AWS services
- An understanding of how to analyze data, verify data quality, and ensure data consistency by using AWS services

Job tasks that are out of scope for the target candidate

The following list contains job tasks that the target candidate is not expected to be able to perform. This list is non-exhaustive. These tasks are out of scope for the exam:

- Perform artificial intelligence and machine learning (AI/ML) tasks.
- Demonstrate knowledge of programming language-specific syntax.
- Draw business conclusions based on data.

Refer to the Appendix for a list of in-scope AWS services and features and a list of out-of-scope AWS services and features.

Exam content

Response types

There are two types of questions on the exam:

- **Multiple choice:** Has one correct response and three incorrect responses (distractors)
- **Multiple response:** Has two or more correct responses out of five or more response options

Select one or more responses that best complete the statement or answer the question. Distractors, or incorrect answers, are response options that a candidate with incomplete knowledge or skill might choose. Distractors are generally plausible responses that match the content area.

Unanswered questions are scored as incorrect; there is no penalty for guessing. The exam includes 50 questions that affect your score.

Unscored content

The exam includes 15 unscored questions that do not affect your score. AWS collects information about performance on these unscored questions to evaluate these questions for future use as scored questions. These unscored questions are not identified on the exam.

Exam results

The AWS Certified Data Engineer - Associate (DEA-C01) exam has a pass or fail designation. The exam is scored against a minimum standard established by AWS professionals who follow certification industry best practices and guidelines.

Your results for the exam are reported as a scaled score of 100–1,000. The minimum passing score is 720. Your score shows how you performed on the exam as a whole and whether you passed. Scaled scoring models help equate scores across multiple exam forms that might have slightly different difficulty levels.

Your score report could contain a table of classifications of your performance at each section level. The exam uses a compensatory scoring model, which means that you do not need to achieve a passing score in each section. You need to pass only the overall exam.

Each section of the exam has a specific weighting, so some sections have more questions than other sections have. The table of classifications contains general information that highlights your strengths and weaknesses. Use caution when you interpret section-level feedback.

Content outline

This exam guide includes weightings, content domains, and task statements for the exam. This guide does not provide a comprehensive list of the content on the exam. However, additional context for each task statement is available to help you prepare for the exam.

The exam has the following content domains and weightings:

- Domain 1: Data Ingestion and Transformation (34% of scored content)
- Domain 2: Data Store Management (26% of scored content)
- Domain 3: Data Operations and Support (22% of scored content)
- Domain 4: Data Security and Governance (18% of scored content)

Domain 1: Data Ingestion and Transformation

Task Statement 1.1: Perform data ingestion.

Knowledge of:

- Throughput and latency characteristics for AWS services that ingest data
- Data ingestion patterns (for example, frequency and data history)
- Streaming data ingestion
- Batch data ingestion (for example, scheduled ingestion, event-driven ingestion)
- Replayability of data ingestion pipelines
- Stateful and stateless data transactions

Skills in:

- Reading data from streaming sources (for example, Amazon Kinesis, Amazon Managed Streaming for Apache Kafka [Amazon MSK], Amazon DynamoDB Streams, AWS Database Migration Service [AWS DMS], AWS Glue, Amazon Redshift)
- Reading data from batch sources (for example, Amazon S3, AWS Glue, Amazon EMR, AWS DMS, Amazon Redshift, AWS Lambda, Amazon AppFlow)
- Implementing appropriate configuration options for batch ingestion
- Consuming data APIs

- Setting up schedulers by using Amazon EventBridge, Apache Airflow, or time-based schedules for jobs and crawlers
- Setting up event triggers (for example, Amazon S3 Event Notifications, EventBridge)
- Calling a Lambda function from Amazon Kinesis
- Creating allowlists for IP addresses to allow connections to data sources
- Implementing throttling and overcoming rate limits (for example, DynamoDB, Amazon RDS, Kinesis)
- Managing fan-in and fan-out for streaming data distribution

Task Statement 1.2: Transform and process data.

Knowledge of:

- Creation of ETL pipelines based on business requirements
- Volume, velocity, and variety of data (for example, structured data, unstructured data)
- Cloud computing and distributed computing
- How to use Apache Spark to process data
- Intermediate data staging locations

Skills in:

- Optimizing container usage for performance needs (for example, Amazon Elastic Kubernetes Service [Amazon EKS], Amazon Elastic Container Service [Amazon ECS])
- Connecting to different data sources (for example, Java Database Connectivity [JDBC], Open Database Connectivity [ODBC])
- Integrating data from multiple sources
- Optimizing costs while processing data
- Implementing data transformation services based on requirements (for example, Amazon EMR, AWS Glue, Lambda, Amazon Redshift)
- Transforming data between formats (for example, from .csv to Apache Parquet)
- Troubleshooting and debugging common transformation failures and performance issues
- Creating data APIs to make data available to other systems by using AWS services

Task Statement 1.3: Orchestrate data pipelines.

Knowledge of:

- How to integrate various AWS services to create ETL pipelines
- Event-driven architecture
- How to configure AWS services for data pipelines based on schedules or dependencies
- Serverless workflows

Skills in:

- Using orchestration services to build workflows for data ETL pipelines (for example, Lambda, EventBridge, Amazon Managed Workflows for Apache Airflow [Amazon MWAA], AWS Step Functions, AWS Glue workflows)
- Building data pipelines for performance, availability, scalability, resiliency, and fault tolerance
- Implementing and maintaining serverless workflows
- Using notification services to send alerts (for example, Amazon Simple Notification Service [Amazon SNS], Amazon Simple Queue Service [Amazon SQS])

Task Statement 1.4: Apply programming concepts.

Knowledge of:

- Continuous integration and continuous delivery (CI/CD) (implementation, testing, and deployment of data pipelines)
- SQL queries (for data source queries and data transformations)
- Infrastructure as code (IaC) for repeatable deployments (for example, AWS Cloud Development Kit [AWS CDK], AWS CloudFormation)
- Distributed computing
- Data structures and algorithms (for example, graph data structures and tree data structures)
- SQL query optimization

Skills in:

- Optimizing code to reduce runtime for data ingestion and transformation
- Configuring Lambda functions to meet concurrency and performance needs
- Performing SQL queries to transform data (for example, Amazon Redshift stored procedures)
- Structuring SQL queries to meet data pipeline requirements
- Using Git commands to perform actions such as creating, updating, cloning, and branching repositories
- Using the AWS Serverless Application Model (AWS SAM) to package and deploy serverless data pipelines (for example, Lambda functions, Step Functions, DynamoDB tables)
- Using and mounting storage volumes from within Lambda functions

Domain 2: Data Store Management

Task Statement 2.1: Choose a data store.

Knowledge of:

- Storage platforms and their characteristics
- Storage services and configurations for specific performance demands
- Data storage formats (for example, .csv, .txt, Parquet)
- How to align data storage with data migration requirements
- How to determine the appropriate storage solution for specific access patterns
- How to manage locks to prevent access to data (for example, Amazon Redshift, Amazon RDS)

Skills in:

- Implementing the appropriate storage services for specific cost and performance requirements (for example, Amazon Redshift, Amazon EMR, AWS Lake Formation, Amazon RDS, DynamoDB, Amazon Kinesis Data Streams, Amazon MSK)
- Configuring the appropriate storage services for specific access patterns and requirements (for example, Amazon Redshift, Amazon EMR, Lake Formation, Amazon RDS, DynamoDB)

- Applying storage services to appropriate use cases (for example, Amazon S3)
- Integrating migration tools into data processing systems (for example, AWS Transfer Family)
- Implementing data migration or remote access methods (for example, Amazon Redshift federated queries, Amazon Redshift materialized views, Amazon Redshift Spectrum)

Task Statement 2.2: Understand data cataloging systems.

Knowledge of:

- How to create a data catalog
- Data classification based on requirements
- Components of metadata and data catalogs

Skills in:

- Using data catalogs to consume data from the data's source
- Building and referencing a data catalog (for example, AWS Glue Data Catalog, Apache Hive metastore)
- Discovering schemas and using AWS Glue crawlers to populate data catalogs
- Synchronizing partitions with a data catalog
- Creating new source or target connections for cataloging (for example, AWS Glue)

Task Statement 2.3: Manage the lifecycle of data.

Knowledge of:

- Appropriate storage solutions to address hot and cold data requirements
- How to optimize the cost of storage based on the data lifecycle
- How to delete data to meet business and legal requirements
- Data retention policies and archiving strategies
- How to protect data with appropriate resiliency and availability

Skills in:

- Performing load and unload operations to move data between Amazon S3 and Amazon Redshift
- Managing S3 Lifecycle policies to change the storage tier of S3 data
- Expiring data when it reaches a specific age by using S3 Lifecycle policies
- Managing S3 versioning and DynamoDB TTL

Task Statement 2.4: Design data models and schema evolution.

Knowledge of:

- Data modeling concepts
- How to ensure accuracy and trustworthiness of data by using data lineage
- Best practices for indexing, partitioning strategies, compression, and other data optimization techniques
- How to model structured, semi-structured, and unstructured data
- Schema evolution techniques

Skills in:

- Designing schemas for Amazon Redshift, DynamoDB, and Lake Formation
- Addressing changes to the characteristics of data
- Performing schema conversion (for example, by using the AWS Schema Conversion Tool [AWS SCT] and AWS DMS Schema Conversion)
- Establishing data lineage by using AWS tools (for example, Amazon SageMaker ML Lineage Tracking)

Domain 3: Data Operations and Support

Task Statement 3.1: Automate data processing by using AWS services.

Knowledge of:

- How to maintain and troubleshoot data processing for repeatable business outcomes
- API calls for data processing
- Which services accept scripting (for example, Amazon EMR, Amazon Redshift, AWS Glue)

Skills in:

- Orchestrating data pipelines (for example, Amazon MWAA, Step Functions)
- Troubleshooting Amazon managed workflows
- Calling SDKs to access Amazon features from code
- Using the features of AWS services to process data (for example, Amazon EMR, Amazon Redshift, AWS Glue)
- Consuming and maintaining data APIs
- Preparing data transformation (for example, AWS Glue DataBrew)
- Querying data (for example, Amazon Athena)
- Using Lambda to automate data processing
- Managing events and schedulers (for example, EventBridge)

Task Statement 3.2: Analyze data by using AWS services.

Knowledge of:

- Tradeoffs between provisioned services and serverless services
- SQL queries (for example, SELECT statements with multiple qualifiers or JOIN clauses)
- How to visualize data for analysis
- When and how to apply cleansing techniques
- Data aggregation, rolling average, grouping, and pivoting

Skills in:

- Visualizing data by using AWS services and tools (for example, AWS Glue DataBrew, Amazon QuickSight)
- Verifying and cleaning data (for example, Lambda, Athena, QuickSight, Jupyter Notebooks, Amazon SageMaker Data Wrangler)
- Using Athena to query data or to create views
- Using Athena notebooks that use Apache Spark to explore data

Task Statement 3.3: Maintain and monitor data pipelines.

Knowledge of:

- How to log application data
- Best practices for performance tuning
- How to log access to AWS services
- Amazon Macie, AWS CloudTrail, and Amazon CloudWatch

Skills in:

- Extracting logs for audits
- Deploying logging and monitoring solutions to facilitate auditing and traceability
- Using notifications during monitoring to send alerts
- Troubleshooting performance issues
- Using CloudTrail to track API calls
- Troubleshooting and maintaining pipelines (for example, AWS Glue, Amazon EMR)
- Using Amazon CloudWatch Logs to log application data (with a focus on configuration and automation)
- Analyzing logs with AWS services (for example, Athena, Amazon EMR, Amazon OpenSearch Service, CloudWatch Logs Insights, big data application logs)

Task Statement 3.4: Ensure data quality.

Knowledge of:

- Data sampling techniques
- How to implement data skew mechanisms
- Data validation (data completeness, consistency, accuracy, and integrity)
- Data profiling

Skills in:

- Running data quality checks while processing the data (for example, checking for empty fields)
- Defining data quality rules (for example, AWS Glue DataBrew)
- Investigating data consistency (for example, AWS Glue DataBrew)

Domain 4: Data Security and Governance

Task Statement 4.1: Apply authentication mechanisms.

Knowledge of:

- VPC security networking concepts
- Differences between managed services and unmanaged services
- Authentication methods (password-based, certificate-based, and role-based)
- Differences between AWS managed policies and customer managed policies

Skills in:

- Updating VPC security groups
- Creating and updating IAM groups, roles, endpoints, and services
- Creating and rotating credentials for password management (for example, AWS Secrets Manager)
- Setting up IAM roles for access (for example, Lambda, Amazon API Gateway, AWS CLI, CloudFormation)
- Applying IAM policies to roles, endpoints, and services (for example, S3 Access Points, AWS PrivateLink)

Task Statement 4.2: Apply authorization mechanisms.

Knowledge of:

- Authorization methods (role-based, policy-based, tag-based, and attribute-based)
- Principle of least privilege as it applies to AWS security
- Role-based access control and expected access patterns
- Methods to protect data from unauthorized access across services

Skills in:

- Creating custom IAM policies when a managed policy does not meet the needs
- Storing application and database credentials (for example, Secrets Manager, AWS Systems Manager Parameter Store)
- Providing database users, groups, and roles access and authority in a database (for example, for Amazon Redshift)
- Managing permissions through Lake Formation (for Amazon Redshift, Amazon EMR, Athena, and Amazon S3)

Task Statement 4.3: Ensure data encryption and masking.

Knowledge of:

- Data encryption options available in AWS analytics services (for example, Amazon Redshift, Amazon EMR, AWS Glue)
- Differences between client-side encryption and server-side encryption
- Protection of sensitive data
- Data anonymization, masking, and key salting

Skills in:

- Applying data masking and anonymization according to compliance laws or company policies
- Using encryption keys to encrypt or decrypt data (for example, AWS Key Management Service [AWS KMS])
- Configuring encryption across AWS account boundaries
- Enabling encryption in transit for data.

Task Statement 4.4: Prepare logs for audit.

Knowledge of:

- How to log application data
- How to log access to AWS services
- Centralized AWS logs

Skills in:

- Using CloudTrail to track API calls
- Using CloudWatch Logs to store application logs
- Using AWS CloudTrail Lake for centralized logging queries
- Analyzing logs by using AWS services (for example, Athena, CloudWatch Logs Insights, Amazon OpenSearch Service)
- Integrating various AWS services to perform logging (for example, Amazon EMR in cases of large volumes of log data)

Task Statement 4.5: Understand data privacy and governance.

Knowledge of:

- How to protect personally identifiable information (PII)
- Data sovereignty

Skills in:

- Granting permissions for data sharing (for example, data sharing for Amazon Redshift)
- Implementing PII identification (for example, Macie with Lake Formation)
- Implementing data privacy strategies to prevent backups or replications of data to disallowed AWS Regions
- Managing configuration changes that have occurred in an account (for example, AWS Config)

Appendix

In-scope AWS services and features

The following list contains AWS services and features that are in scope for the exam. This list is non-exhaustive and is subject to change. AWS offerings appear in categories that align with the offerings' primary functions:

Analytics:

- Amazon Athena
- Amazon EMR
- AWS Glue
- AWS Glue DataBrew
- AWS Lake Formation
- Amazon Kinesis Data Firehose
- Amazon Kinesis Data Streams
- Amazon Managed Service for Apache Flink
- Amazon Managed Streaming for Apache Kafka (Amazon MSK)
- Amazon OpenSearch Service
- Amazon QuickSight

Application Integration:

- Amazon AppFlow
- Amazon EventBridge
- Amazon Managed Workflows for Apache Airflow (Amazon MWAA)
- Amazon Simple Notification Service (Amazon SNS)
- Amazon Simple Queue Service (Amazon SQS)
- AWS Step Functions

Cloud Financial Management:

- AWS Budgets
- AWS Cost Explorer

Compute:

- AWS Batch
- Amazon EC2
- AWS Lambda
- AWS Serverless Application Model (AWS SAM)

Containers:

- Amazon Elastic Container Registry (Amazon ECR)
- Amazon Elastic Container Service (Amazon ECS)
- Amazon Elastic Kubernetes Service (Amazon EKS)

Database:

- Amazon DocumentDB (with MongoDB compatibility)
- Amazon DynamoDB
- Amazon Keyspaces (for Apache Cassandra)
- Amazon MemoryDB for Redis
- Amazon Neptune
- Amazon RDS
- Amazon Redshift

Developer Tools:

- AWS CLI
- AWS Cloud9
- AWS Cloud Development Kit (AWS CDK)
- AWS CodeBuild
- AWS CodeCommit
- AWS CodeDeploy
- AWS CodePipeline

Frontend Web and Mobile:

- Amazon API Gateway

Machine Learning:

- Amazon SageMaker

Management and Governance:

- AWS CloudFormation
- AWS CloudTrail
- Amazon CloudWatch
- Amazon CloudWatch Logs
- AWS Config
- Amazon Managed Grafana
- AWS Systems Manager
- AWS Well-Architected Tool

Migration and Transfer:

- AWS Application Discovery Service
- AWS Application Migration Service
- AWS Database Migration Service (AWS DMS)
- AWS DataSync
- AWS Schema Conversion Tool (AWS SCT)
- AWS Snow Family
- AWS Transfer Family

Networking and Content Delivery:

- Amazon CloudFront
- AWS PrivateLink
- Amazon Route 53
- Amazon VPC

Security, Identity, and Compliance:

- AWS Identity and Access Management (IAM)
- AWS Key Management Service (AWS KMS)
- Amazon Macie
- AWS Secrets Manager
- AWS Shield
- AWS WAF

Storage:

- AWS Backup
- Amazon Elastic Block Store (Amazon EBS)
- Amazon Elastic File System (Amazon EFS)
- Amazon S3
- Amazon S3 Glacier

Out-of-scope AWS services and features

The following list contains AWS services and features that are out of scope for the exam. This list is non-exhaustive and is subject to change. AWS offerings that are entirely unrelated to the target job roles for the exam are excluded from this list:

Analytics:

- Amazon FinSpace

Business Applications:

- Alexa for Business
- Amazon Chime
- Amazon Connect
- Amazon Honeycode
- AWS IQ
- Amazon WorkDocs
- Amazon WorkMail

Compute:

- AWS App Runner
- AWS Elastic Beanstalk
- Amazon Lightsail
- AWS Outposts
- AWS Serverless Application Repository

Containers:

- Red Hat OpenShift Service on AWS (ROSA)

Database:

- Amazon Timestream

Developer Tools:

- AWS Fault Injection Simulator (AWS FIS)
- AWS X-Ray

Frontend Web and Mobile:

- AWS Amplify
- AWS AppSync
- AWS Device Farm
- Amazon Location Service
- Amazon Pinpoint
- Amazon Simple Email Service (Amazon SES)

Internet of Things (IoT):

- FreeRTOS
- AWS IoT 1-Click
- AWS IoT Device Defender
- AWS IoT Device Management
- AWS IoT Events
- AWS IoT FleetWise
- AWS IoT RoboRunner
- AWS IoT SiteWise
- AWS IoT TwinMaker

Machine Learning:

- Amazon CodeWhisperer
- Amazon DevOps Guru

Management and Governance:

- AWS Activate
- AWS Managed Services (AMS)

Media Services:

- Amazon Elastic Transcoder
- AWS Elemental Appliances and Software
- AWS Elemental MediaConnect
- AWS Elemental MediaConvert
- AWS Elemental MediaLive
- AWS Elemental MediaPackage
- AWS Elemental MediaStore
- AWS Elemental MediaTailor
- Amazon Interactive Video Service (Amazon IVS)
- Amazon Nimble Studio

Migration and Transfer:

- AWS Mainframe Modernization
- AWS Migration Hub

Storage:

- EC2 Image Builder

Survey

How useful was this exam guide? Let us know by [taking our survey](#).