aws

**AWS PARTNER CERTIFICATION READINESS**

# Content Review Session

Week 4 – Domain 3

Data Operations and Support

1

# *OPTIONAL* AWS Skill Builder Subscription

The Skill Builder subscription provides access to official AWS Certification practice exams, self-paced digital training content including open-ended challenges, self-paced labs, and game-based learning. **Please note, the Skill Builder subscription is not required for this Accelerator program.**

## Free digital training
### *LINK HERE*

**Special features include:**

- 600+ digital courses
- Learning plans
- 10 Practice Question Sets
- *AWS Cloud Quest (Foundational)*

## Individual subscription
### *LINK HERE*

**Everything in free digital training, plus:**

- AWS Cloud Quest (Intermediate - Advanced)
- AWS Certification Official Practice Exams
- Enhanced Exam Prep Courses
- Unlimited access to 1000+ hands-on labs
- AWS Jam Journeys (lab-based challenges)
- AWS Digital Classroom (Annual only)

Individual subscriptions are priced **at $29 USD per month** (*Flexibility to cancel anytime*) or **$449 USD per year**.

Access **65** Data Engineer - Associate Practice Exam Questions with feedback on your answer choices

# Today's Learning Outcomes

During this session, we will cover:

- Automate data processing by using AWS services

- Analyze data by using AWS services

- Maintain and monitor data pipelines

- Ensure data quality

# 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

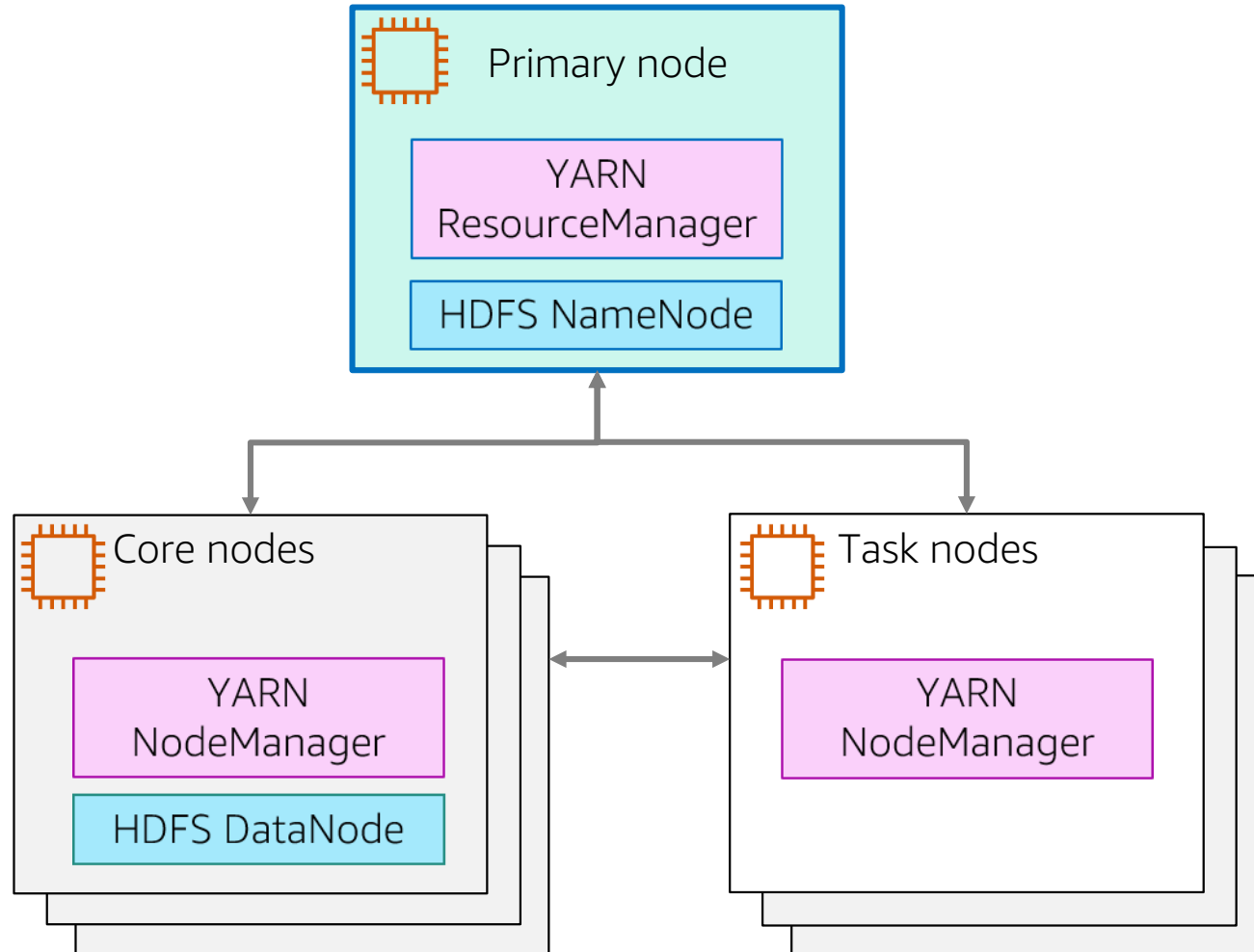- Which services accept scripting (for example, Amazon EMR, Amazon Redshift, AWS Glue)

## 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)
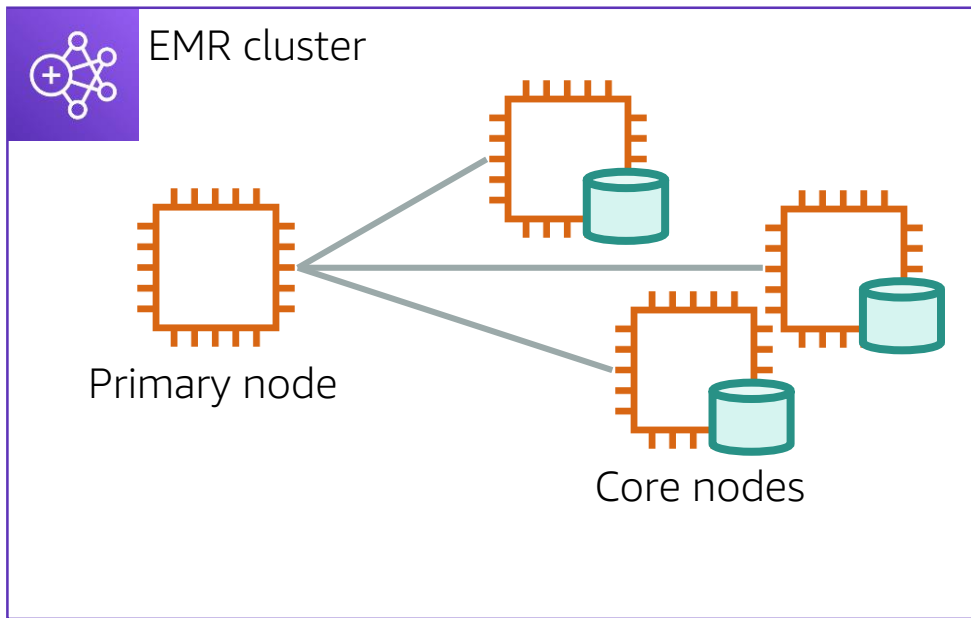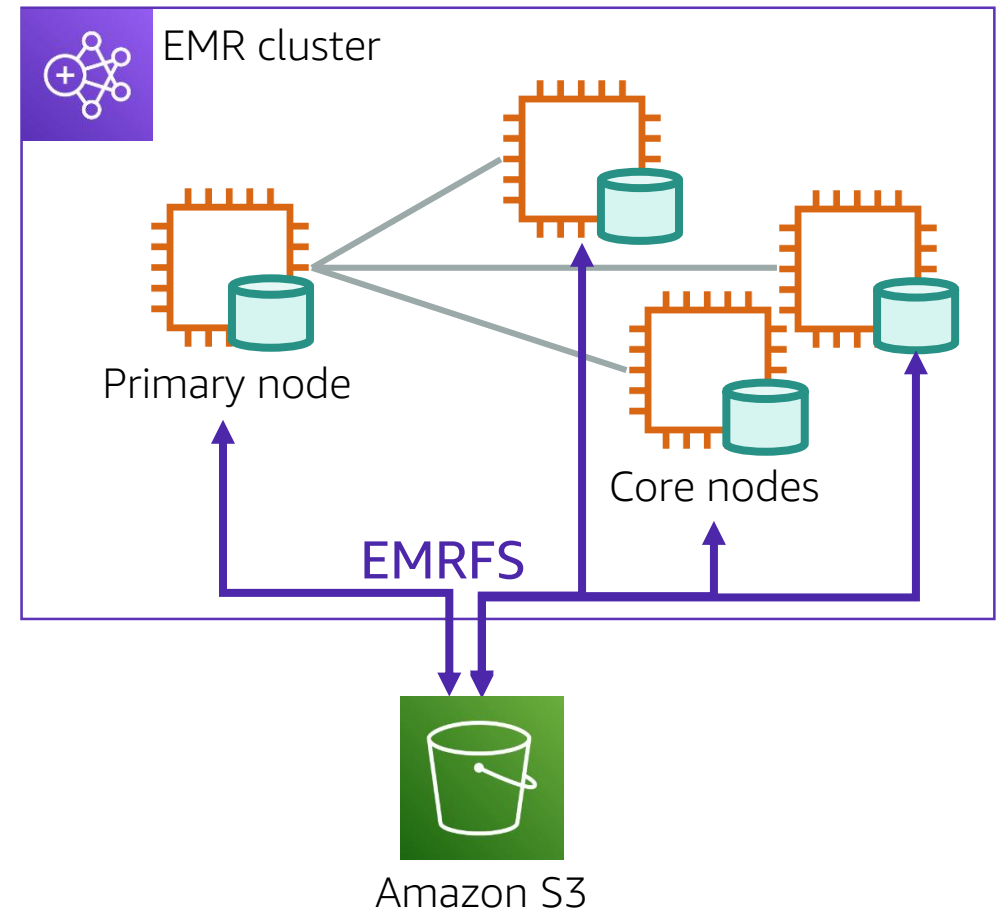
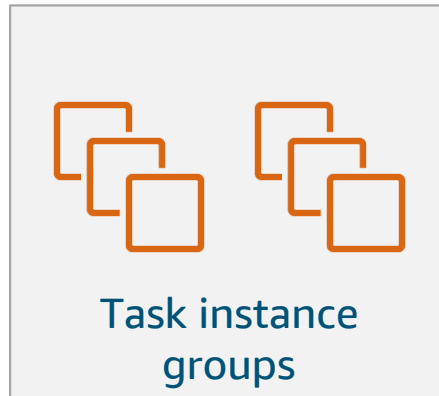# Amazon Elastic Map Reduce (EMR)

# Amazon EMR data storage options
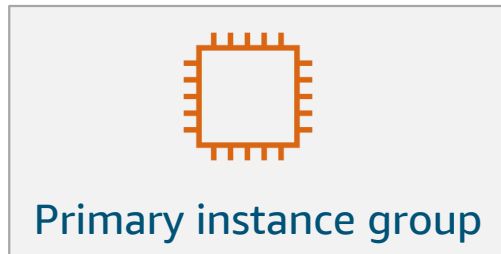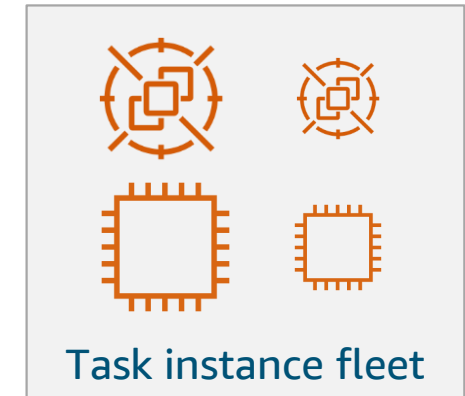
# Node configuration options

Instance groups

Primary instance group

Core instance group

Task instance groups

Instance fleets

Primary instance fleet

Core instance fleet

Task instance fleet

# Transforming data in Amazon EMR



**EMR cluster with Hadoop**

**Leader node**

Begin processing

**Step 1** — Formatting and compression

**Step 2** — Partitioning

**Step 3** — Cleansing and normalizing

**Step 4** — Aggregation

End processing

Core nodes

EMRFS — Raw data

EMRFS — Processed data

# Considerations in choosing Spark on Amazon EMR or AWS Glue ETL

| Consideration | Spark on Amazon EMR | AWS Glue ETL |
|---|---|---|
| Responsibility model | Fully managed | Serverless |
| Degree of control | More flexibility | More automation |
| Data model | Schema before load | Schema on read |
| Scalability | Policy-based or managed scaling | Specify max workers and concurrency |
| Pricing | Stable cost for consistent workloads | Pay for use based on needs |

aws

# Amazon EventBridge

- EventBridge is a serverless service for building event-driven, loosely-coupled applications by routing events between sources and targets

- Event buses receive events from many sources and deliver to multiple targets, with optional event transformation

# Redshift Data API



Amazon Elastic Compute Cloud (Amazon EC2)

Amazon Elastic Container Service (Amazon ECS)

AWS Lambda

AWS Tools and SDKs

Amazon EventBridge

AWS AppSync

Jupyter notebooks

Amazon Redshift Data API

Amazon Redshift

# Redshift Data API code example (Python)

Using the Redshift Data API to run SQL:

```python
query_str = "select  count(*) as record_count from stocksummary.stocks"

res = client_redshift.execute_statement(Database= db,
                                        SecretArn= secret_arn,
                                        Sql= query_str,
                                        ClusterIdentifier= cluster_id)
```
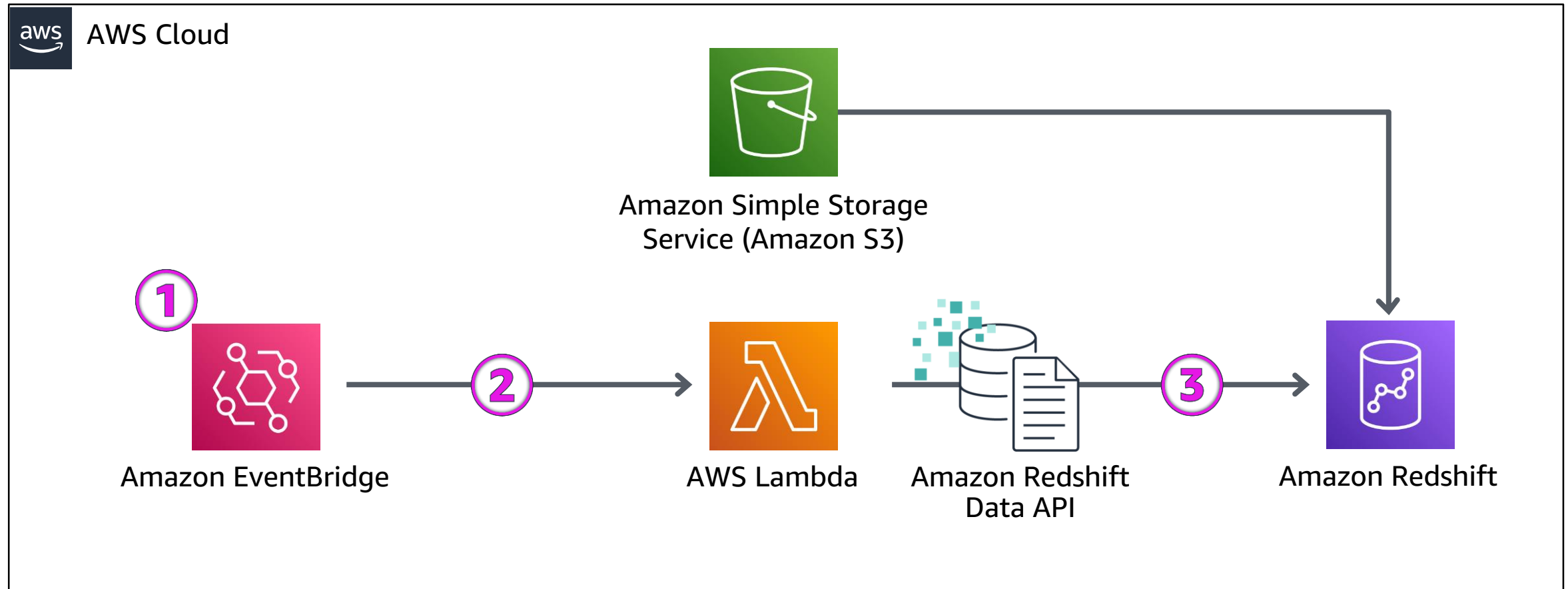
Using the COPY command and the Redshift Data API:

```python
query = "COPY stocksummary.stocks FROM '" + s3_data_path + "' IAM_ROLE '" + redshift_iam_role + "' CSV IGNOREHEADER 1;"

resp = client_redshift.execute_statement(Database= db, SecretArn= secret_arn, Sql= query, ClusterIdentifier= cluster_id)
```

# Architecting an event-driven architecture using the Amazon Redshift Data API

aws

AWS PARTNER CERTIFICATION READINESS

# Domain 3: Data Operations and Support

Analyze data by using AWS services
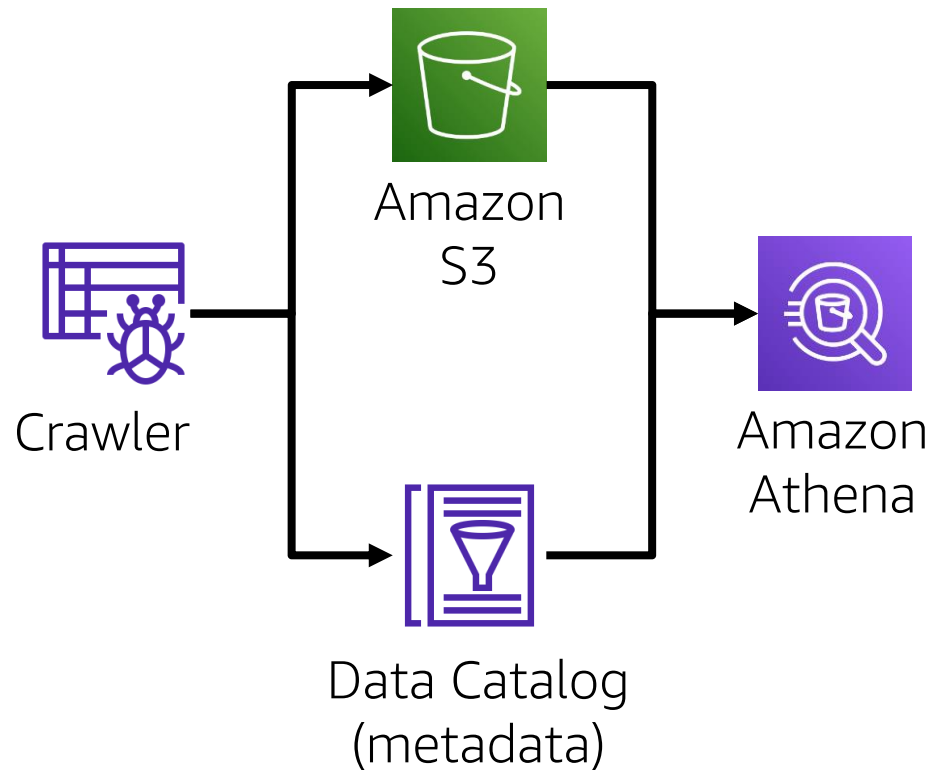
# Analyze data by using AWS services

## 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

# Athena

## Serverless query engine



Crawler → Amazon S3 → Amazon Athena
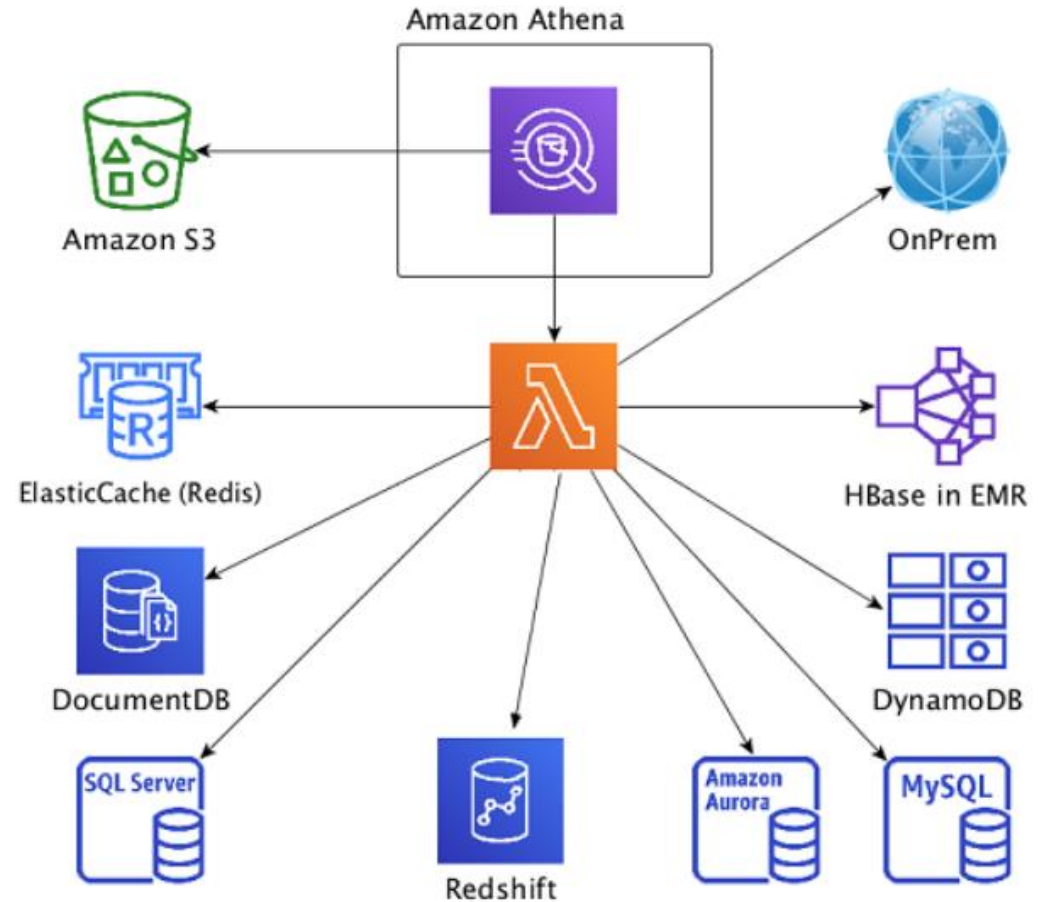
Data Catalog (metadata)

## Benefits

- Query Amazon S3 data directly or the AWS Glue Data Catalog

- Use SQL based on Presto

- Supports CSV, JSON, ORC, Avro, Parquet

- Can integrate with Amazon QuickSight

- Support for federated query

# Athena Federated Query

- Enables federated querying across multiple data sources with SQL

- Utilizes Lambda connectors for different data sources

- Allows building custom connectors for proprietary data sources

# Athena workgroups

- Isolate workloads.

- Control user access.

- Manage query usage, costs, and data limits

- A workgroup is an IAM resource managed by Athena.

```
{
    "Version": "2012-10-17",
    "Statement": [
        {
            "Effect": "Allow",
            "Action": [
                "athena:StartQueryExecution",
                "athena:StopQueryExecution"
            ],
            "Resource": [
                "arn:aws:athena:us-east-1:123456789012:workgroup/test_workgroup"
            ]
        }
    ]
}
```

# QuickSight

- Scalable, serverless business intelligence service
- Embeddable into your existing applications
- Dynamic, machine learning (ML) powered insights

# Accessing data with QuickSight

# SPICE

## Super-fast, Parallel, In-memory Calculation Engine



- Faster processing
- Reduced wait time vs. direct queries
- Reduced cost through reuse

# QuickSight embedded analytics

Embed visualizations in any application for desktop or mobile applications.

1. Create a dashboard
2. Apply permissions
3. Authenticate your app server
4. Embed via JavaScript SDK



Desktop or mobile →

# 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)

# Data pipelines on AWS

An efficient and well-designed data integration pipeline is critical for making the data available and trusted amongst the analytics consumers.

Here are some considerations to review when designing data pipelines:

| Factor | AWS Glue Workflow | AWS Step Function | Amazon Managed Workflow for Apache Airflow (MWAA) |
|---|---|---|---|
| **Use case** | Suitable when your pipeline consists of mostly AWS Glue jobs and crawlers. | Suitable when there is a need to integrate with different services, including AWS Lambda, SSM, and so on. | Compatible with open-source Airflow and suitable when you want to reuse existing Airflow assets. |
| **Infrastructure** | Serverless | Serverless | Managed service |
| **Building a data pipeline** | Build a data pipeline using an AWS Glue job written in Python or /Scala and crawlers. | Build a data pipeline using the Step Functions console. Possible to integrate with non-supported services using Lambda. | Workflows are created as DAGs, which are defined within a Python file that defines the DAG's structure as code. |

# AWS Glue Workflow

- Create/visualize complex ETL workflows in AWS Glue

- Workflow triggers:
  - Schedules
  - On-demand
  - EventBridge events

- Visual graph of execution progress and component dependencies

# Monitoring a Data Lake

Data sources → Ingestion → Data stores → Catalog and processing → Search and analytics → Visualization

← Security and monitoring →

Amazon CloudWatch

Logs  Rule  Events  Alarm

AWS CloudTrail

aws

# Domain 3: Data Operations and Support

Ensure data quality

# 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)

# Glue Databrew



Data sources → AWS Glue DataBrew → Amazon S3 → Analytics and machine learning

## AWS Glue DataBrew

- Visual data prep environment
- Cleanses and prepares data
- Over 250 built-in transformations
- Evaluates data quality
- Automates at scale

# Glue Databrew

# SageMaker

- Fully managed ML service

- End-to-end model development

- Built-in algorithms/models, auto-tuning, hosting, and more

# SageMaker Data Wrangler

- Prepare and featurize data with Data Wrangler flows involving minimal coding.

- Analyze data quality, visualize features, and perform quick modeling.

- Export workflows to integrate with SageMaker pipelines, Feature Store, or custom scripts.

aws

# Thank you for attending this session