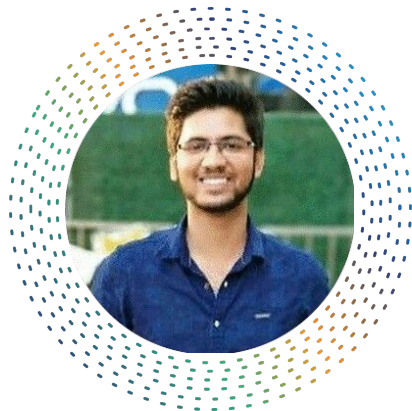# Data Engineering Workshop

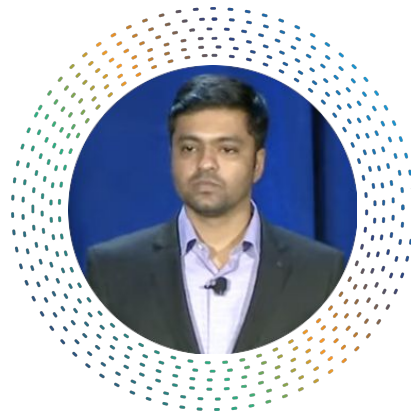**Rajas Walavalkar & Sanchit Jain**
**8th Dec, Thursday**
**2:00 PM to 4:00 PM**

# Speakers

**Rajas Walavalkar**

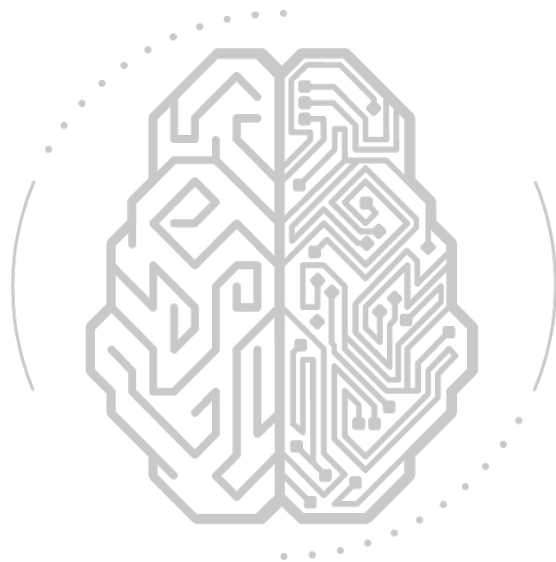Solution Architect - AWS at Quantiphi
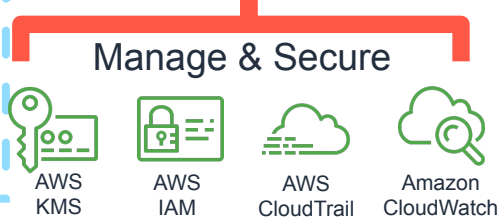
AWS Community Builder
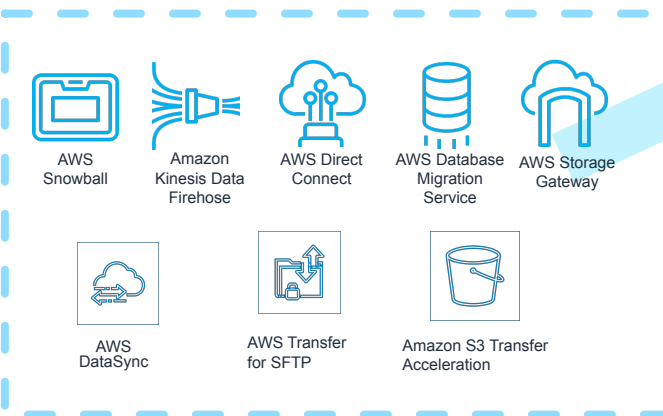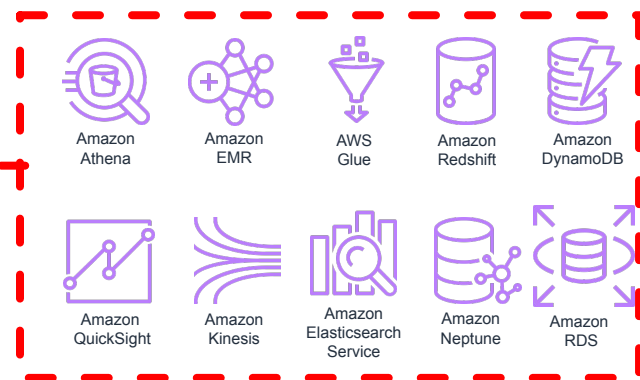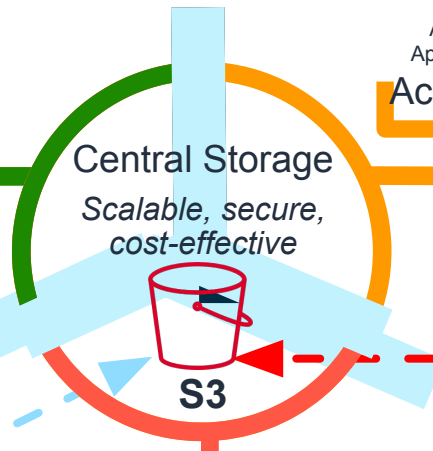
FOLLOW ME

**Sanchit Jain**

DNA & Cloud Practice Lead  - AWS at Quantiphi

AWS Hero & AWS Ambassador

FOLLOW ME

# What we will learn?

# Session's Focus – Query The Data Lake
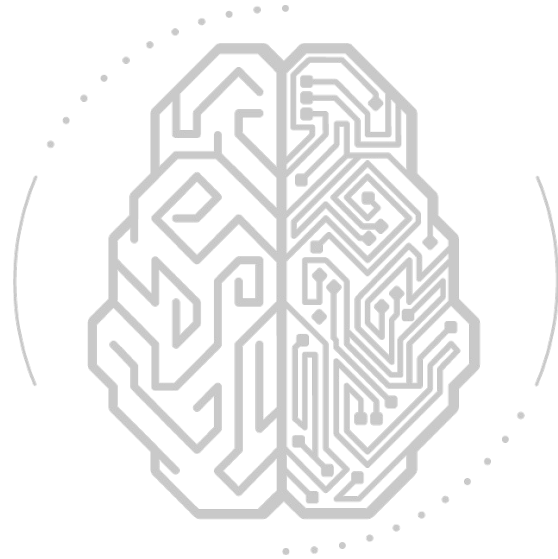
# AWS Simple Storage Service - S3

# Amazon S3 Introduction

- Amazon Simple Storage Service (Amazon S3) is storage for the Internet. We can use Amazon S3 to store and retrieve any amount of data at any time, from anywhere on the web.

- Amazon S3 is one of the main building blocks of AWS - It's advertised as "infinitely scaling" storage. Many websites use Amazon S3 as a backbone, and even many AWS services use Amazon S3 as an integration as well

- S3 Buckets
    - Buckets must have a globally unique name (across all regions all accounts)
    - Buckets are defined at the region level
    - S3 looks like a global service but buckets are created in a region

- S3 Objects
    - Objects (files) have a Key, and the key is the FULL path: s3://my-bucket/my_folder1/another_folder/my_file.txt

- S3 Durability:
    - High durability (99.999999999%, 11 9's) of objects across multiple AZ
    - If you store 10,000,000 objects with Amazon S3, you can on average expect to incur a loss of a single object once every 10,000 years
    - Same for all storage classes

- Availability
    - S3 standard has 99.99% availability, which means it will not be available 53 minutes a year
    - Varies depending on storage class
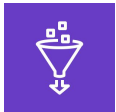
# Amazon S3 Features

- Amazon S3 - Static Websites
    - S3 can host static websites and have them accessible on the www

- Amazon S3 - Versioning
    - You can version your files in Amazon S3
    - It is enabled at the bucket level
    - It is best practice to version your buckets
        - Protect against unintended deletes (ability to restore a version)
        - Easy roll back to previous version

- S3 Access Logs
    - Very helpful to come down to the root cause of an issue, or audit usage, view suspicious patterns, etc
    - Any request made to S3, from any account, authorized or denied, will be logged into another S3 bucket

- S3 Replication
    - Must enable versioning in source and destination, and Buckets can be in different accounts
    - Replication is Asynchronous in nature - Copying is asynchronous
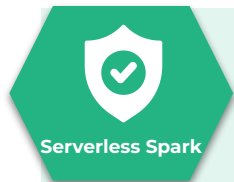
# AWS Glue - Features & Benefits

# AWS Glue Overview

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 all the capabilities needed for data integration so that you can start analyzing your data and putting it to use in minutes instead of months.

## FEATURES

**Serverless Spark**

There is no infrastructure to maintain. Allocate needed compute power and run jobs. Job starts in few seconds and can run at petabyte scale

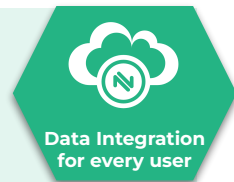**Data Integration for every user**

Development environments catered to different skill sets

**Cost Effective**

All-in-one pricing model includes infrastructure and is 55% cheaper than other cloud data integration options

**Handles complex workloads**
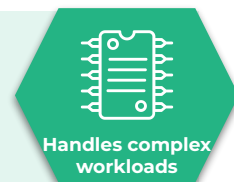
Glue connects to 60+ data sources, processes petabytes of data in real-time, batch and event driven modes

**No Lock-In**

Develop data integration pipelines in open source SparkSQL, PySpark and Scala

**More Power**

AWS Glue automates much of the effort spent in building, maintaining, and running ETL jobs

# AWS Glue - Components

# AWS Glue: Components

**Data Catalog**

- Hive metastore compatible with enhanced functionality
- Crawlers automatically extract metadata and create tables
- Integrated with Athena, Amazon Redshift Spectrum

**Job Authoring**

- Auto-generates ETL code
- Builds on open frameworks—Python and Spark
- Developer-centric—editing, debugging, sharing

**Job Execution**

- Runs jobs on a serverless Spark platform
- Provides flexible scheduling
- Handles dependency resolution, monitoring, and alerting

# AWS Glue Data Catalog

Manage table metadata through a Hive metastore API or Hive SQL. Supported by tools like Hive, Presto, Spark, etc. AWS added a few extensions:

- **Search** over metadata for data discovery

- **Connection info—**JDBC URLs, credentials

- **Classification** for identifying and parsing files

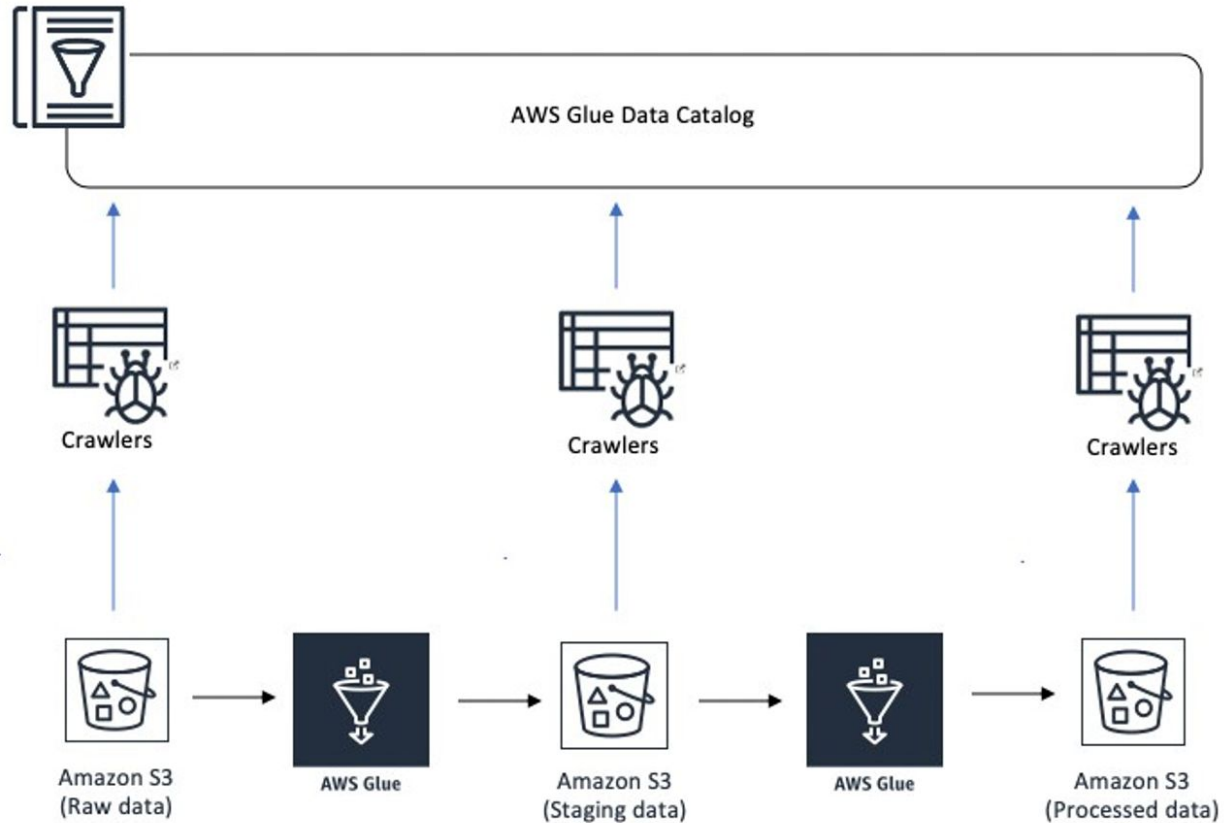- **Versioning** of table metadata as schemas evolve and other  metadata are updated

Populate using Hive DDL, bulk import, or automatically through **crawlers**

# AWS Glue Data Catalog: Crawlers

Crawlers automatically build your Data Catalog and keep it in sync

- Automatically discover new data, extract schema definitions

    - Detect schema changes and version tables

    - Detect Hive style partitions on Amazon S3

- Built-in classifiers for popular types; custom classifiers using Grok  expressions

- Run ad hoc or on a schedule; serverless—only pay when crawler runs

# AWS Glue In Action

# What is AWS Glue DataBrew ?

AWS Glue DataBrew is a visual data preparation tool that makes it easy for data analysts and data scientists to prepare data with an interactive, point-and-click visual interface without writing code.

## CAPABILITIES OF GLUE DATABREW

Profile

Clean and Normalize

Map Data Lineage

Automate

## NEED FOR DATABREW

*"Upto 80% of data analysis time is spent on preparing data"*

### Time Consuming
- Multi-step process to extract, clean, normalize & load data at scale
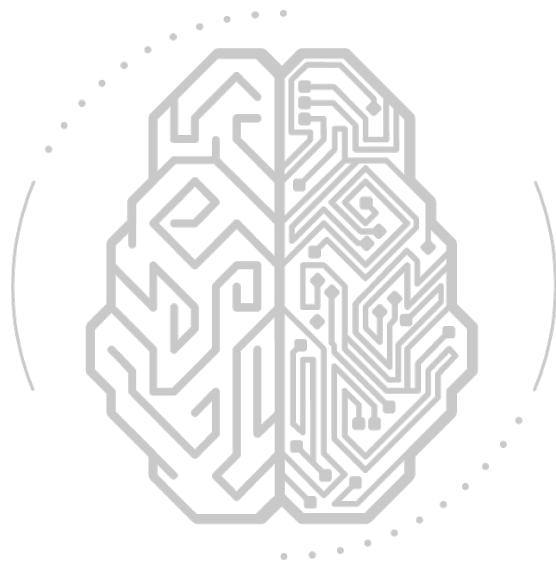- The right tools for the right persona must be integrated

### Expensive
- Costly user licenses & siloed tools that cause rework
- Often requires moving large amount of data into silos

### Manual
- Needs a lot of code-based heavy lifting to work at scale
- Hard to operationalize & build repeatable workflows

# Querying the Data Lake with Amazon Athena

An interactive query service that makes it easy to analyze data directly from Amazon S3 using Standard SQL

# Amazon Athena

- Query data in your Amazon S3 based data lake

- Analyze infrastructure, operation, and application logs

- Interactive analytics using popular BI tools

- Self-service data exploration for data scientists

- Embed analytics capabilities into your applications

# What does it look like?

# Athena is Serverless

- No Infrastructure or administration

- Zero Spin up time

- Transparent upgrades

# Familiar Technologies Under the Covers

**Used for SQL Queries**

In-memory distributed query engine
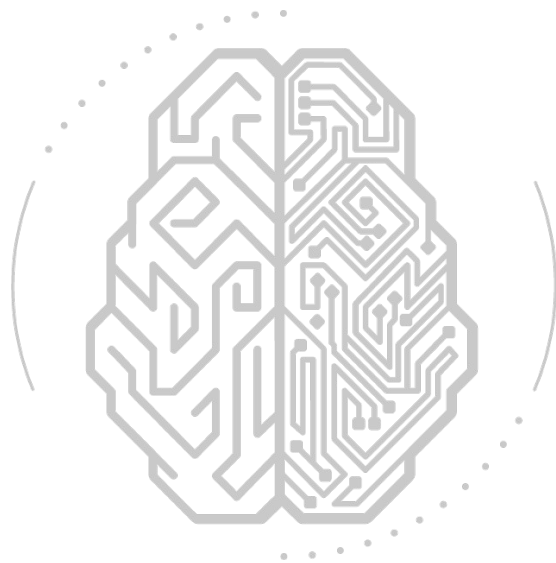
ANSI-SQL compatible with extensions

**Used for DDL functionality**

Complex data types

Multitude of formats

Supports data partitioning

# Amazon Athena is Cost Effective

- Pay per query

- $5 per TB scanned from S3

- DDL Queries and failed queries are free

- Save by using compression, columnar formats, partitions

# Athena Workgroups

# Athena Workgroups

Athena Workgroups are used to isolate queries between different teams, workloads or applications, and to set  on amount of data each query or the entire workgroup can process

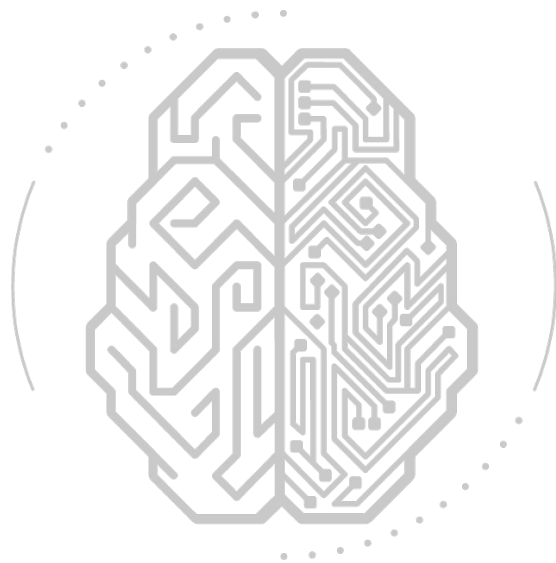Workload Isolation

Query Metrics

Cost Controls

# Workgroups – Cost Controls

- Per query data scanned threshold; exceeding, will cancel query

- Trigger alarms to notify of increasing usage and cost
- Disable Workgroup when all queries exceed a maximum threshold
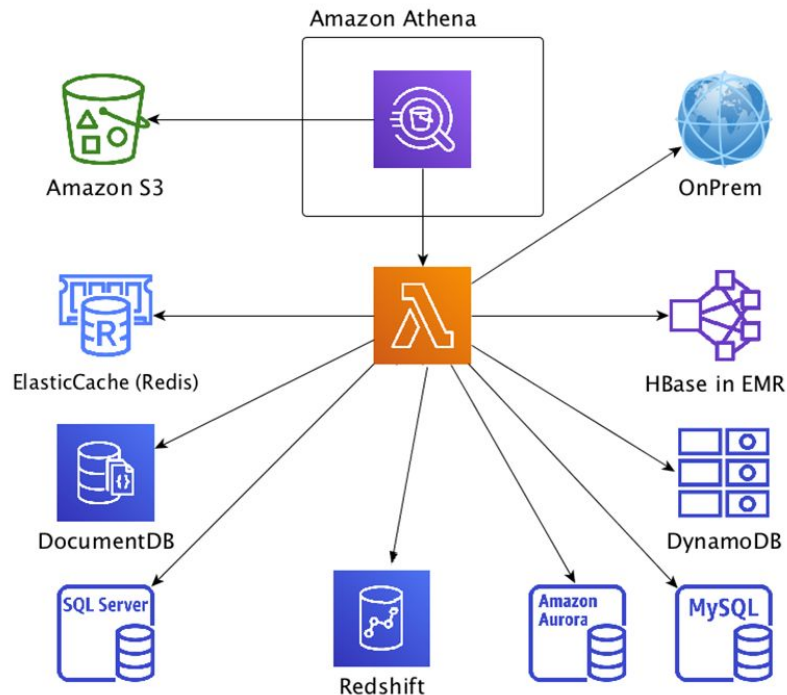
Any Athena metric

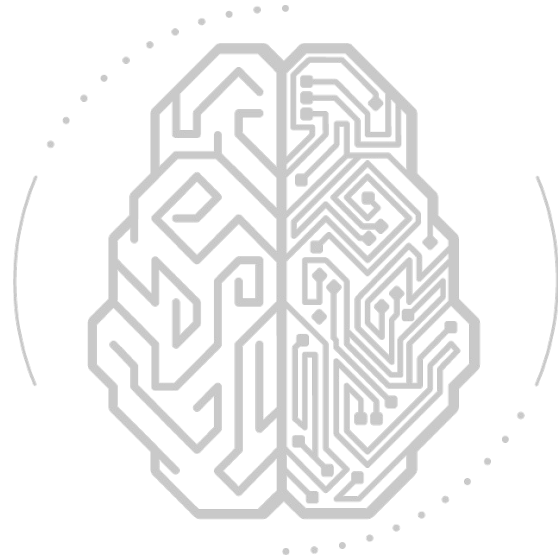| Data limit | Time period | Action | |
|---|---|---|---|
| 10 Gigabytes | Not applicable | Query will be cancelled. | |
| 1 Terabytes | 24 hours | Send notification to topic : arn:aws:sns:us-east-1:9 | 9:AthenaAlarm |
| 10 Gigabytes | 1 hour | Send notification to topic : arn:aws:sns:us-east-1:9 | 9:AthenaAlarm |

# Athena federated query

# Athena federated query

- Run query across relational, non-relational, object, or custom data sources

- Run query across On-Premises or cloud data sources

- Can be used for ad-hoc investigations, or complex pipelines, or applications

# Visualizing the Data Lake using Amazon QuickSight

# Why Amazon QuickSight?

## Cloud native = No servers = Auto-Scale
No servers or software to manage, maintain, deploy. Start with 10s of users and scale to 10s of 1000s

## Fully integrated with AWS
Build end-to-end analytics in AWS. Secure private VPC access, fine-grained access control, ML integrations

## Secure and global
End-to-end encryption. Native High Availability. 10 Global regions. HIPAA, PCI, ISO, SOC and FedRamp eligibility

## Easy to develop and maintain
Design with Amazon QuickSight, integrate with APIs. Secure data with row-level security and authenticate seamlessly via single sign-on

## Fast, consistent performance
Fast, predictable performance every time. Concurrent users or increased interactions do not slow down the system

## ML insights
Contextual, relevant insights with ML-powered anomaly detection, forecasting, alerts and customizable narratives

## Insights for everyone
Provide access to all users, pay only for usage. No upfront costs, no charges for inactive users

## Customize and embed
Embed in applications and enable analytics in hours, not months or years. Use themes to match application/corporate branding
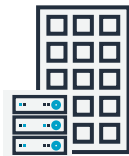
# Connect to your data, wherever it is

QuickSight is natively integrated with AWS data sources, as well as on-premises and hosted databases and third party business applications

## On-premises

Securely connect to on-premise databases and flat files like Excel and CSV

- Excel
- CSV
- Teradata
- MySQL
- SQL Server
- PostgreSQL

## In the cloud

Connect to hosted database, big data formats, and secure VPCs

- Presto
- Spark
- SQL Server
- Postgre SQL
- MariaDB
- Snowflake
- IoT Analytics

## Applications

Connect directly to third party business applications

- Salesforce
- Square
- Adobe Analytics
- Jira
- ServiceNow
- Twitter
- Github

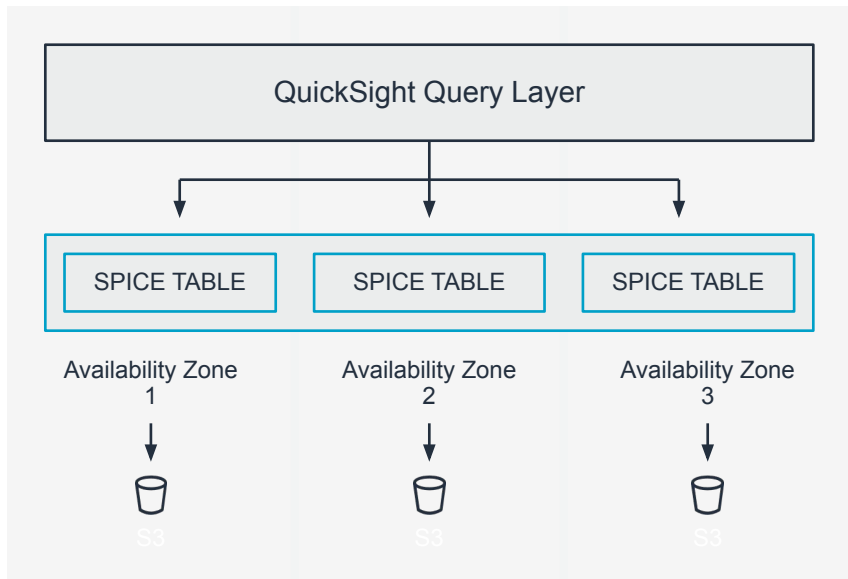# Data Prep

Optional step when creating data sets:

- Preview data

- Rename, remove fields, change data types

- Create new calculated fields

- Filter rows

- Issue direct query or ingest to SPICE

- Push down custom SQL queries

- Join across all data sources supported by QuickSight including file-to-file,file-to-database, and database-to-database joins

# SPICE

QuickSight is powered by SPICE, a super-fast calculation engine that delivers performance and scale, regardless of how many users are active.



Up to 10X faster (millisecond latency)

Fault-tolerant, self-healing

Support for high concurrency

Backed up in S3 (Write Ahead Log)

Instant failover with zero impact

# User Types / User Roles

## Admin
Manage Users
Manage SPICE Capacity
Manage VPC Connections
Manage Account Settings

## Author
Create Data Sets
Create Analyses
Create Dashboards

## Reader
Consume Dashboards

## QS Admin
Sometimes separate from Business Users, sometimes the same
Usually has AWS Console

## Analyst
Sometimes in IT, sometimes Business Users
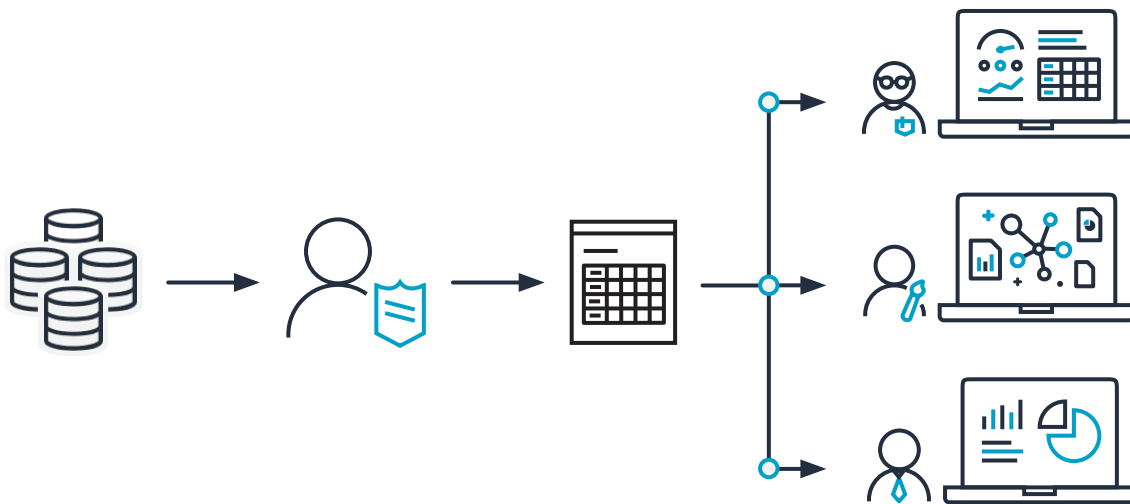'Data Analyst'
'Data Engineer'
'BI Engineer'

## Business User
Anyone
Can be internal or external users (customers/partners3rd

# Data governance

Create managed datasets that give power users and authors the flexibility to perform self-serve analytics on data that you control.

**Create datasets that:**

- Can be shared with any user
- Automatically refresh
- Have row level security
- Users cannot modify
- Dynamically update with changes

# Differentiate with natural language and ML

**Auto narratives**
Summarize your business metrics in plain language

**Forecasting**
Machine learning forecasting with point and click simplicity
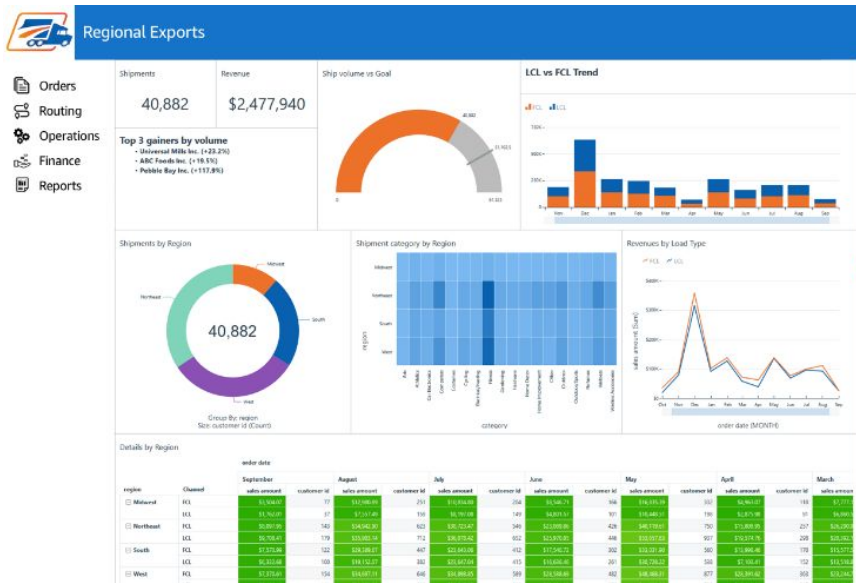
**Anomaly detection**
Discover unexpected trends and outliers against millions of business metrics

**ML predictions**
Visualize and build predictive dashboards with Amazon SageMaker models

# Embed Amazon QuickSight Dashboards



Fully interactive with drill down, filtering, & external links

Personalized views with row-level security

No servers to manage, no long-term commitments

Pay for usage with pay-per-session reader pricing

Seamless authentication

# Amazon QuickSight Examples

THANK YOU