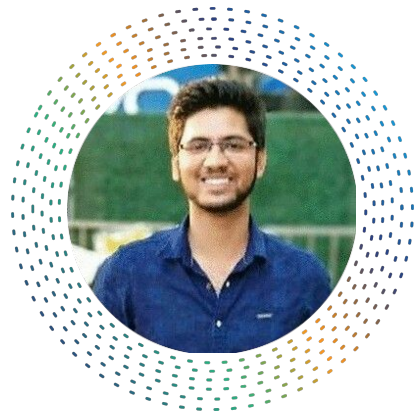


Data Engineering Workshop

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

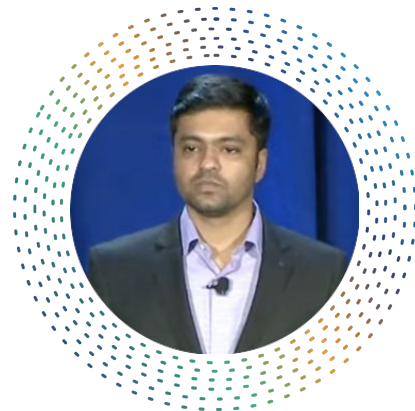


Speakers



Rajas Walavalkar

Solution Architect - AWS at Quantiphi
AWS Community Builder

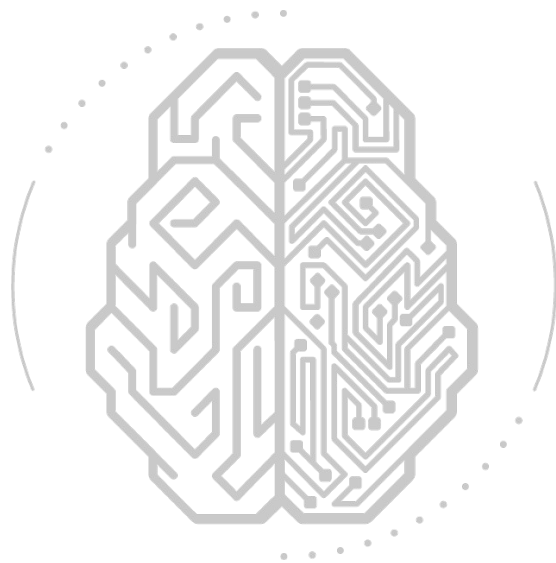


Sanchit Jain

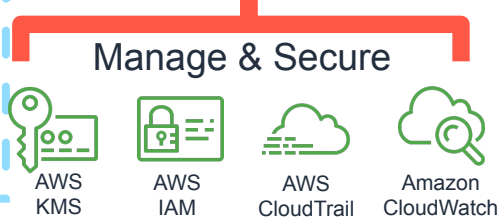
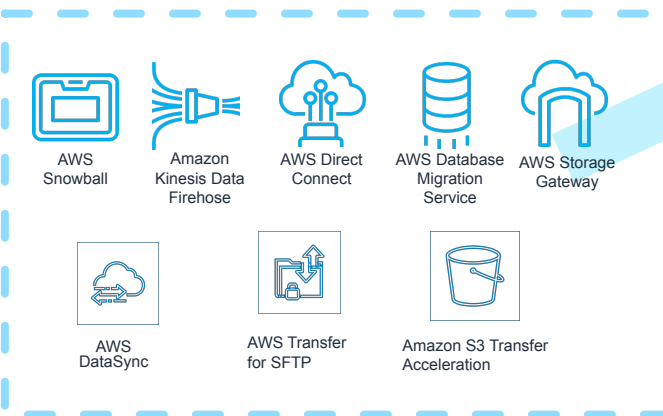
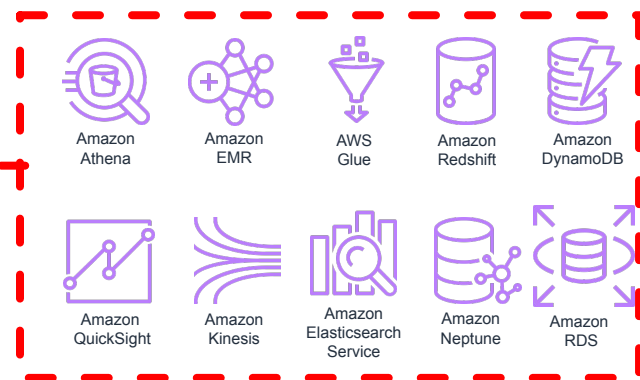
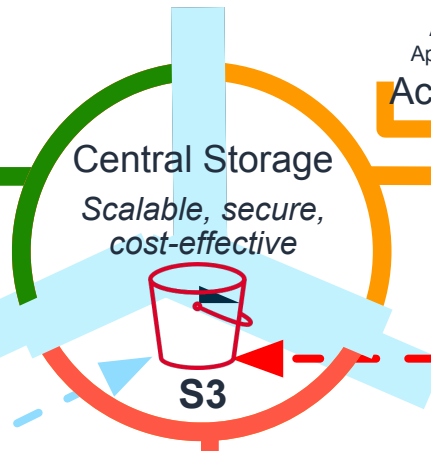
DNA & Cloud Practice Lead - AWS at Quantiphi
AWS Hero & AWS Ambassador



What we will learn?



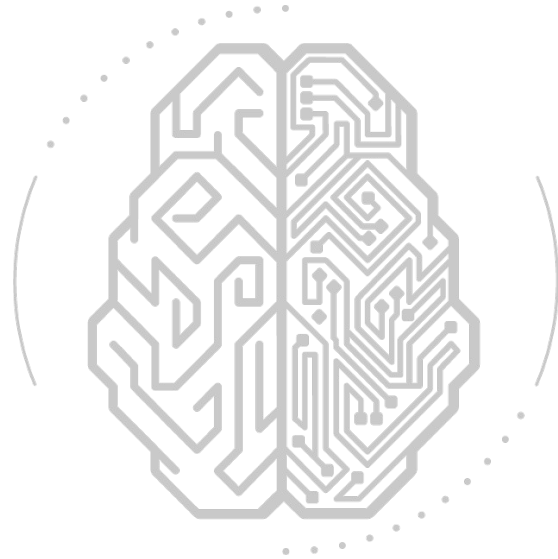
Session's Focus – Query The Data Lake



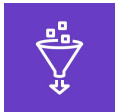
Data Ingestion

Analytics & Serving

AWS Glue - Features & Benefits

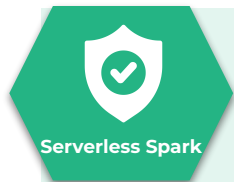


AWS Glue Overview



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



Cost Effective

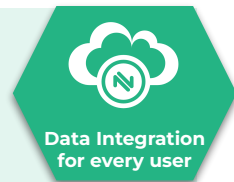
All-in-one pricing model includes infrastructure and is 55% cheaper than other cloud data integration options



No Lock-In

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

Development environments catered to different skill sets



Data Integration for every user

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



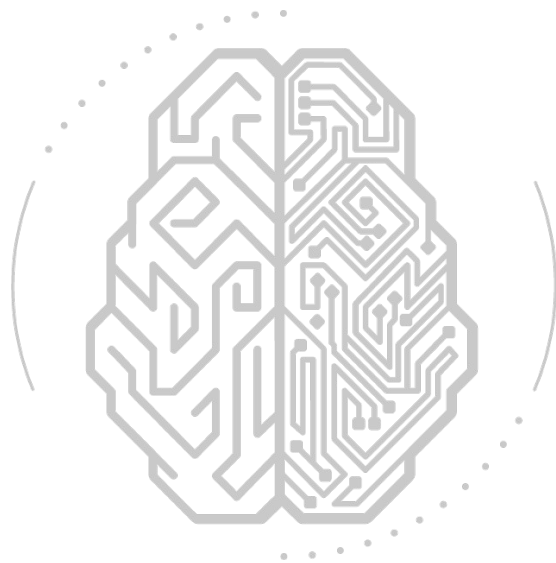
Handles complex workloads

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



More Power

AWS Glue - Components



AWS Glue: Components



Data Catalog

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



Job Authoring

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



Job Execution

- Runs jobs on a serverless Spark platform
- 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

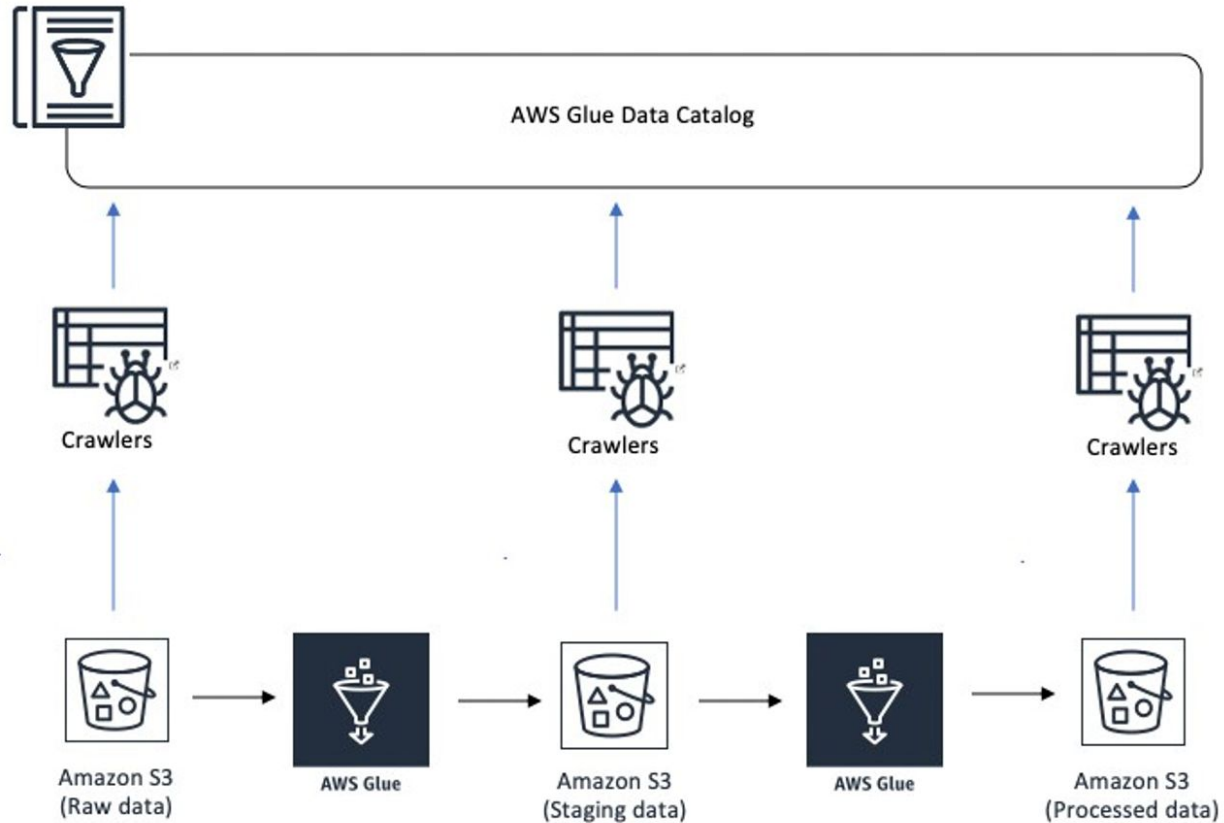
Populate using Hive DDL, bulk import, or automatically through **crawlers**

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

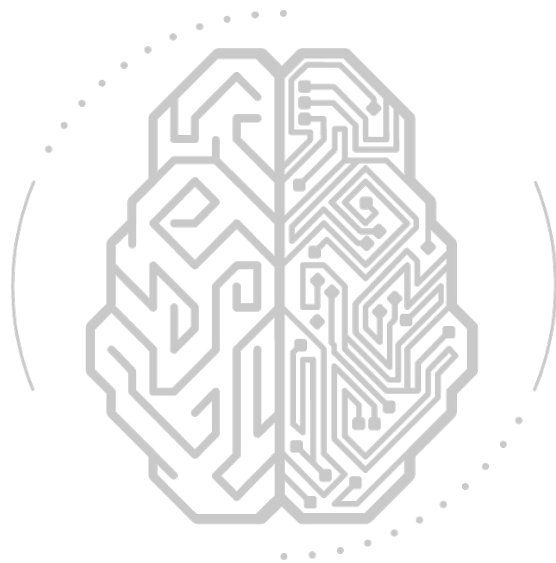
Expensive

- 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





Amazon Athena

Start querying data instantly. Get results in seconds. Pay only for the queries you run.

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

Query Editor

Saved Queries

History

Catalog

Catalog

Sample Queries

DATABASE

default

Table Name

adsads

cloudfront_logs

cloudtrailtable

elb_logs

ncc

ncctest

prestoinstance

prestostat

taxinyc_csv

taxinyc_par

test3

test4

wikistats

language (string)

page_title (string)

hits (bigint)

retrived_size (bigint)

wikistats_parq

```
1 select sum(hits) as hits,language from default.wikistats
2 group by language
3 order by 1 desc
4 limit 10
```

Execute

Save As

or create a

New Query

Recent Queries

Query

Columns

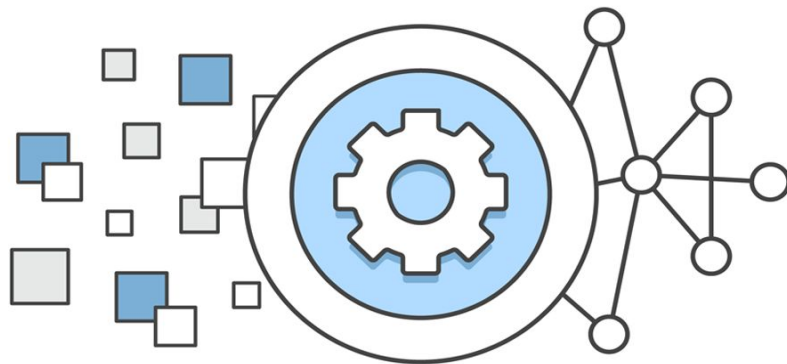
Results

Chart

	hits	language
1	213917076	en
2	43673225	ja
3	25593194	es
4	16637343	de
5	10579788	fr
6	7958810	commons.m
7	6664552	pt
8	6108102	ru
9	5857921	pl
10	5783183	it

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



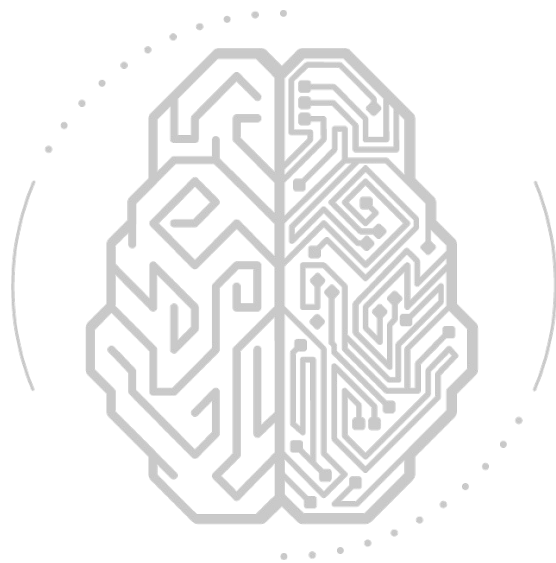
Used for DDL functionality

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

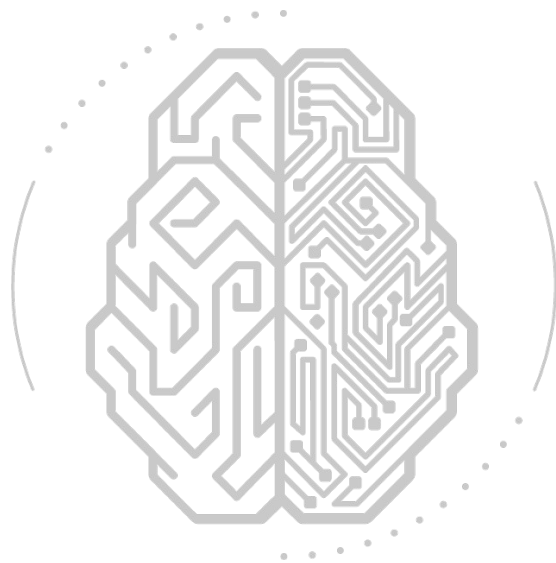
Workgroups – Cost Controls

- 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

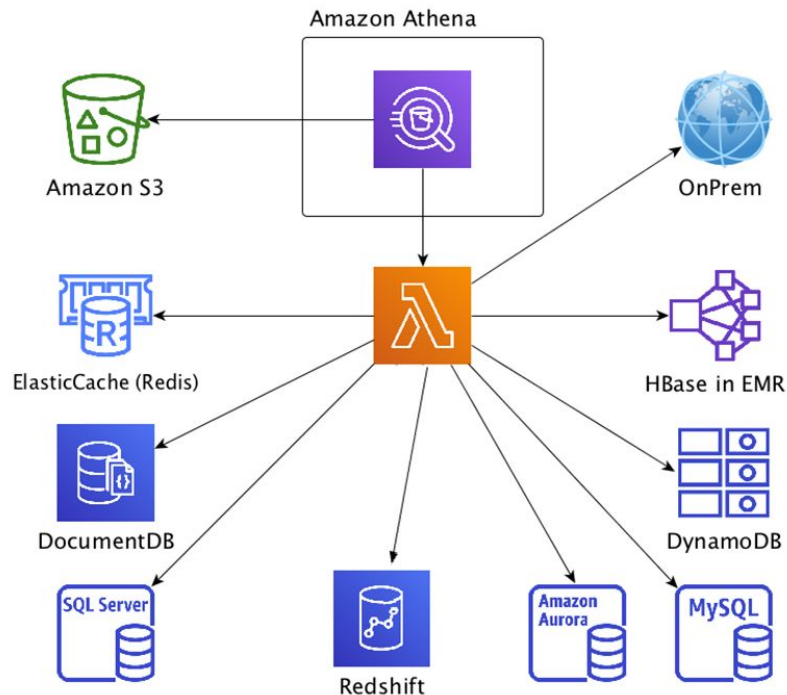
Data limit	Time period	Action
10 Gigabytes	Not applicable	Query will be cancelled.
1 Terabytes	24 hours	Send notification to topic : arn:aws:sns:us-east-1:9 9:AthenaAlarm
10 Gigabytes	1 hour	Send notification to topic : arn:aws:sns:us-east-1:9 9:AthenaAlarm

Athena federated query

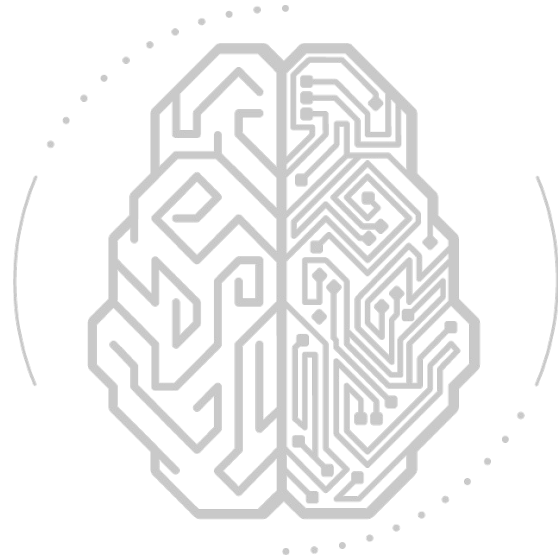


Athena federated query

- 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

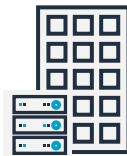


Connect to your data, wherever it is

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



Adobe Analytics



servicenow

Data Prep

Optional step when creating data sets:

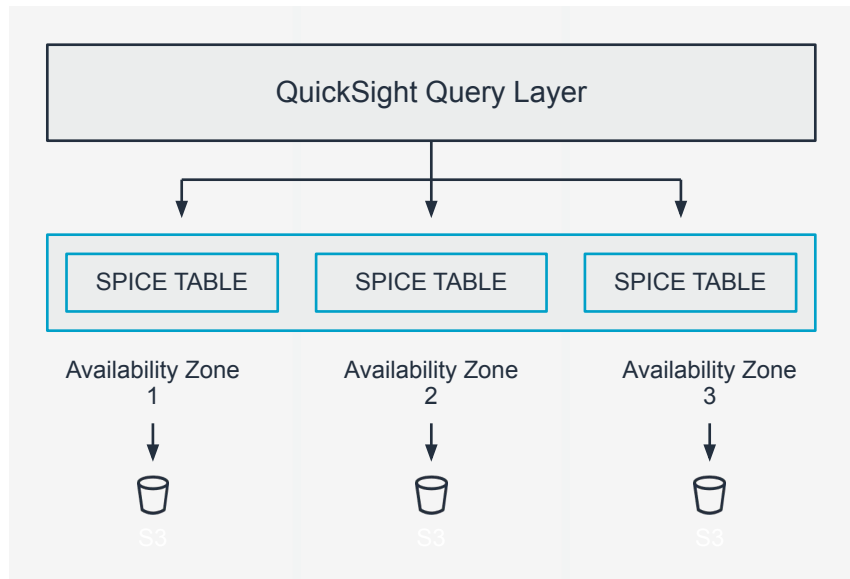
- Preview data
- Rename, remove fields, change data types
- Create new calculated fields
- Filter rows
- Issue direct query or ingest to SPICE
- Push down custom SQL queries
- Join across all data sources supported by QuickSight including file-to-file, file-to-database, and database-to-database joins

The screenshot shows the Qlik Sense interface with the following components:

- Top Bar:** Includes a search bar with "retail_sales", buttons for "Save & visualize", "Save", and "Cancel", and user information for "N. Virginia" and "5676796...".
- Data source:** Shows "retail_sales" as the selected data source. It includes a "Refresh now" button and a "SPICE" status indicator.
- Fields:** A list of fields is shown, including "customer id", "customer name", and "customer type".
- Data sources:** A table of data sources is displayed with columns: cust..., city, state, coun..., zip c..., region, and age r... The table contains 10 rows of data.
- Join types:** A section for configuring joins, showing the relationship between "all_flights" and "airport_id" as an "Inner" join.
- Configure join:** A dialog box is open, showing the relationship between "all_flights" and "airport_id" as an "Inner" join.

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

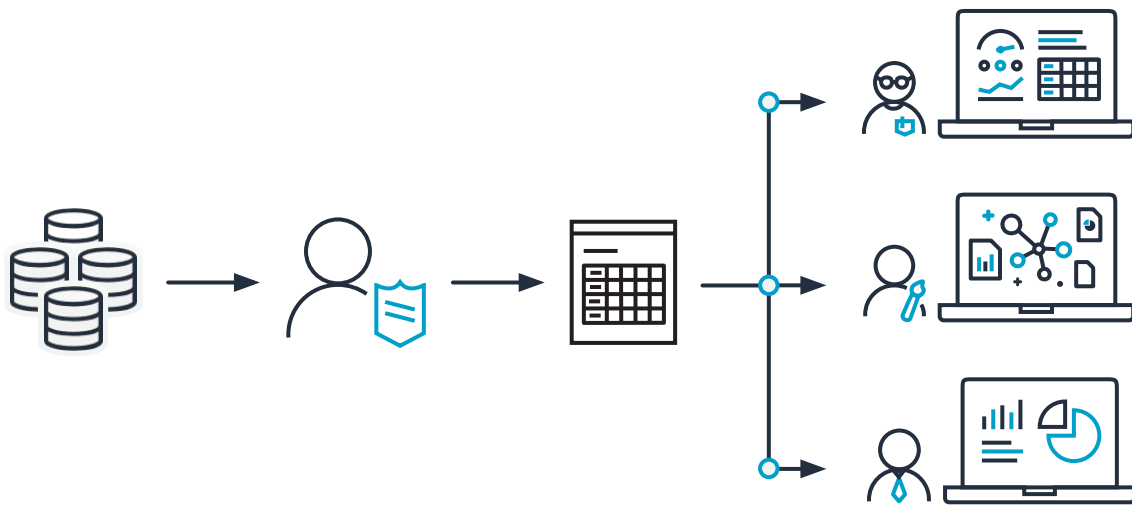
- Anyone
- Can be internal or external users (customers/partners)

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

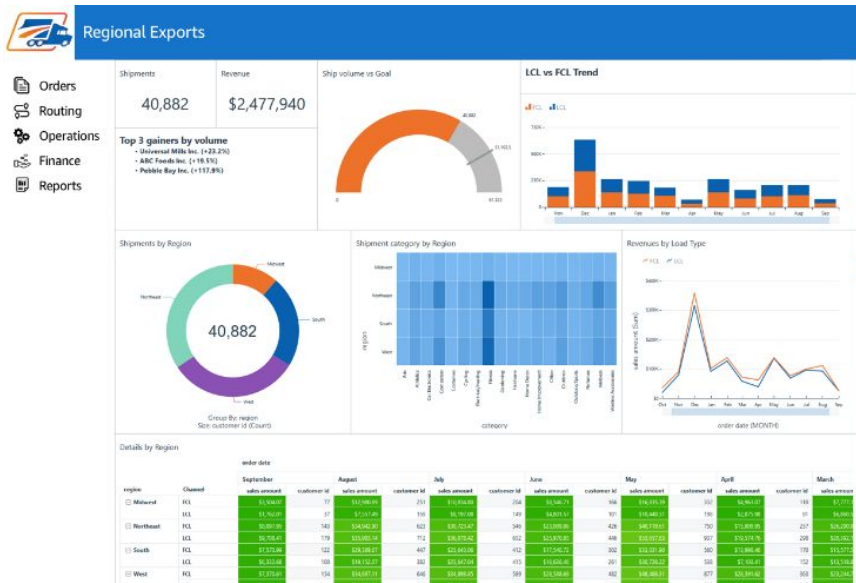
Discover unexpected trends and outliers against millions of business metrics



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