



AWS Solutions Architect Associate

Session 302

Analytics & Database: Glue,
Athena & Redshift Spectrum

July/2024

SERVERLESS!!!



S3 Select

Using S3 Select Feature from S3, you can query S3 File without downloading or download file, increase efficiency.
Don't need compute service. Run on S3 on its conditions, encryption.
Very Simple.

Input: CSV, JSON or Parquet.

Input file coding: UTF-8

Input Compression: None, GZIP, BZIP2. Parquet doesn't apply.

Output: CSV, JSON

Not works for pure Big Bata because it's only work on one file.

Adv: Not parsing, apply on one file, works "on the fly", without compute unit.

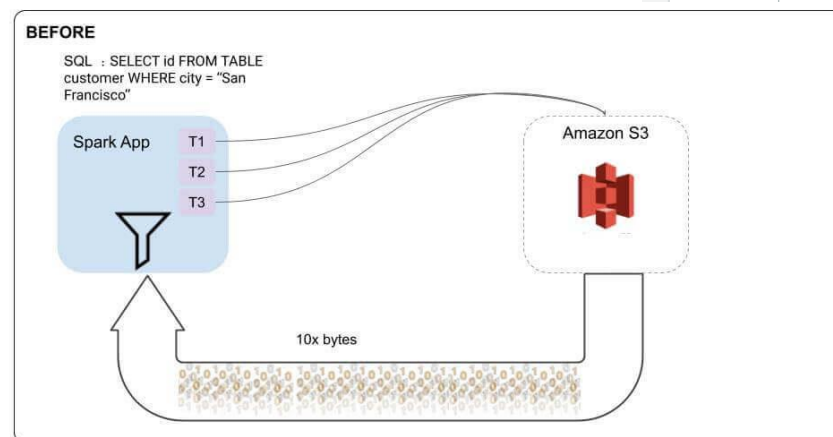
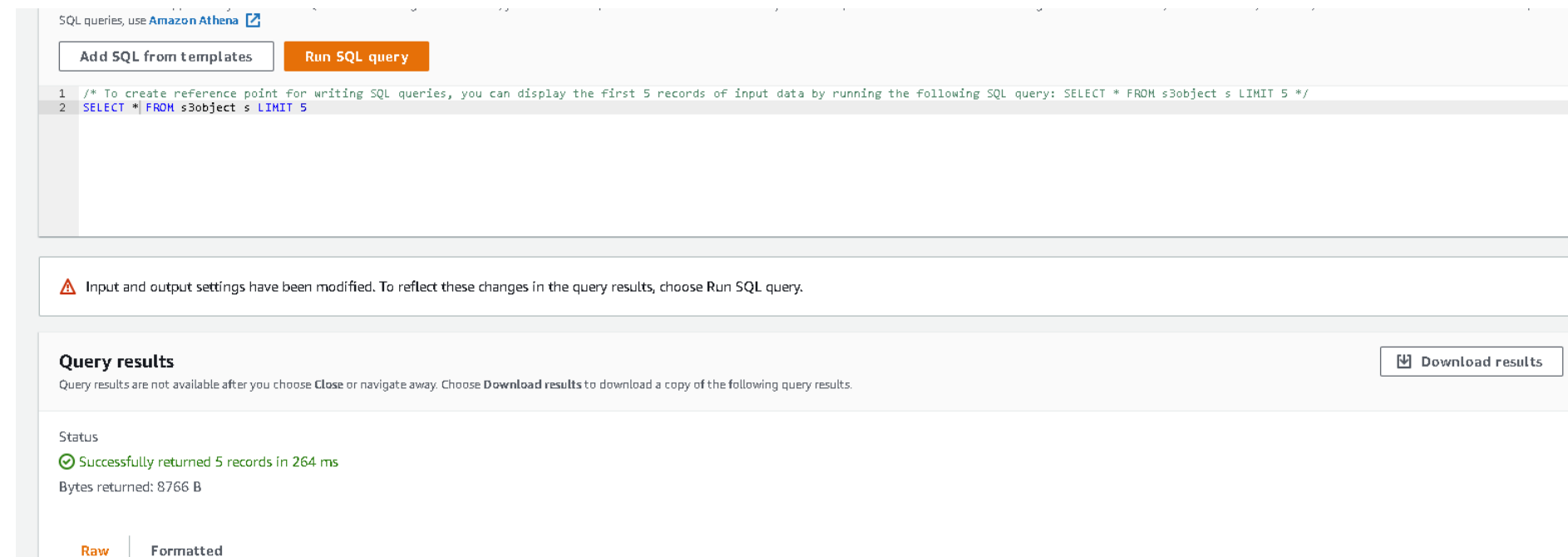
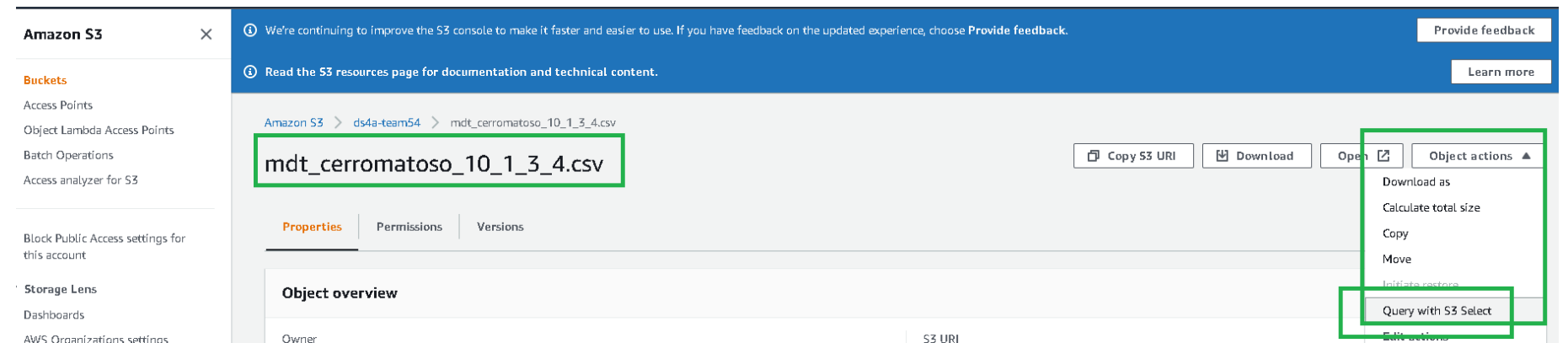


Figure 1

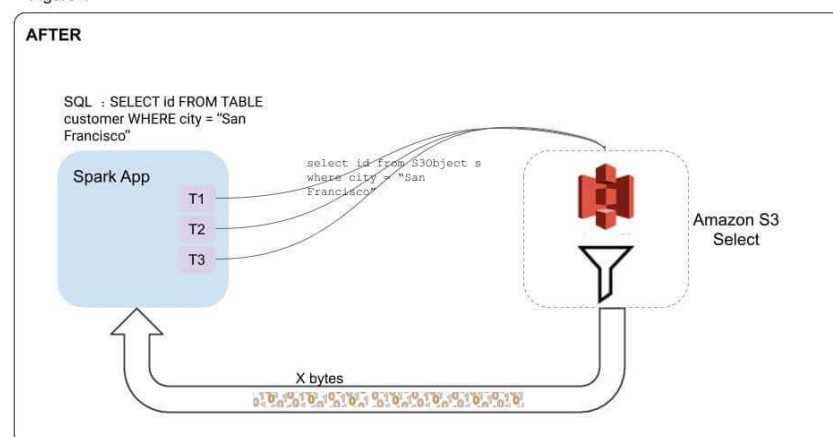


Figure 2

lst,total_power,mean_LCC,stddev_LCC,totalcurrent,period_day,totalmagnitude,lon_CC,lat_CC,n_storms,mean_dist_WC,std_dist_WC,mean_power_WC,std_power_WC,mean_mindist_WC,std_

L7:15:00,2018-05-02 17:25:00,145.0,0.0,6.140960326675138,3.004653087526722,30.982056245460225,26.012694767932594,0.7999999999999998,2.0,13.8,-74.84808824894765,7.18916749

L7:16:00,2018-05-02 17:26:00,142.0,0.0,6.140960326675138,2.634984601492083,33.88795131270982,25.556742532236147,0.7999999999999998,2.0,13.8,-74.84808824894765,7.189167498



ETL – Extract Transform and Load >>> **DATA INTEGRATION**

Data Sources: S3 (**semistructured**), DB: RDS, NoSQL: DynamoDB.

Store on target Data Stores.

Discover/Infer the data schema using classifiers.

Tool to Extract: AWS Glue Crawler.

Data Sources

AWS Glue supports the following data sources:

- Data stores
 - Amazon S3
 - Amazon Relational Database Service (Amazon RDS)
 - Third-party JDBC-accessible databases
 - Amazon DynamoDB
 - MongoDB and Amazon DocumentDB (with MongoDB compatibility)
- Data streams
 - Amazon Kinesis Data Streams
 - Apache Kafka

Data Targets

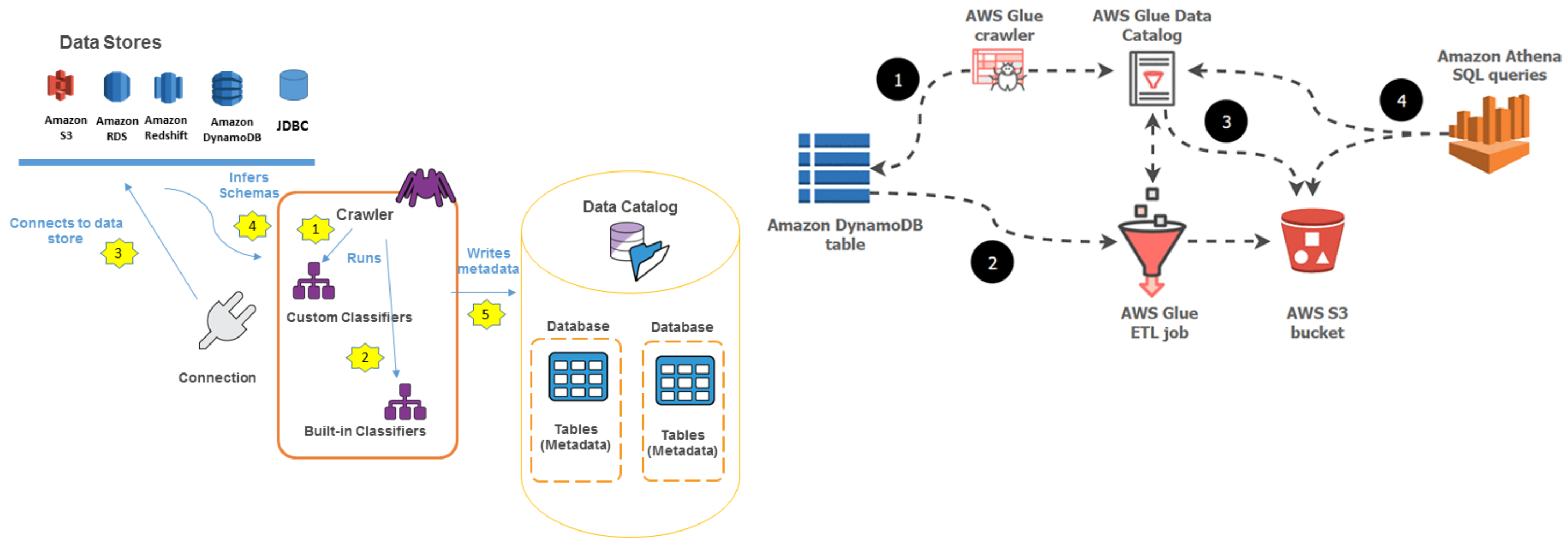
AWS Glue supports the following data targets:

- Amazon S3
- Amazon Relational Database Service (Amazon RDS)
- Third-party JDBC-accessible databases
- MongoDB and Amazon DocumentDB (with MongoDB compatibility)

Access type that crawler uses	Data stores
Native client	<ul style="list-style-type: none">• Amazon Simple Storage Service (Amazon S3)• Amazon DynamoDB
JDBC	Within Amazon Relational Database Service (Amazon RDS) or external to Amazon RDS: <ul style="list-style-type: none">• Amazon Aurora• MariaDB• Microsoft SQL Server• MySQL• Oracle• PostgreSQL
MongoDB client	<ul style="list-style-type: none">• MongoDB• Amazon DocumentDB (with MongoDB compatibility)



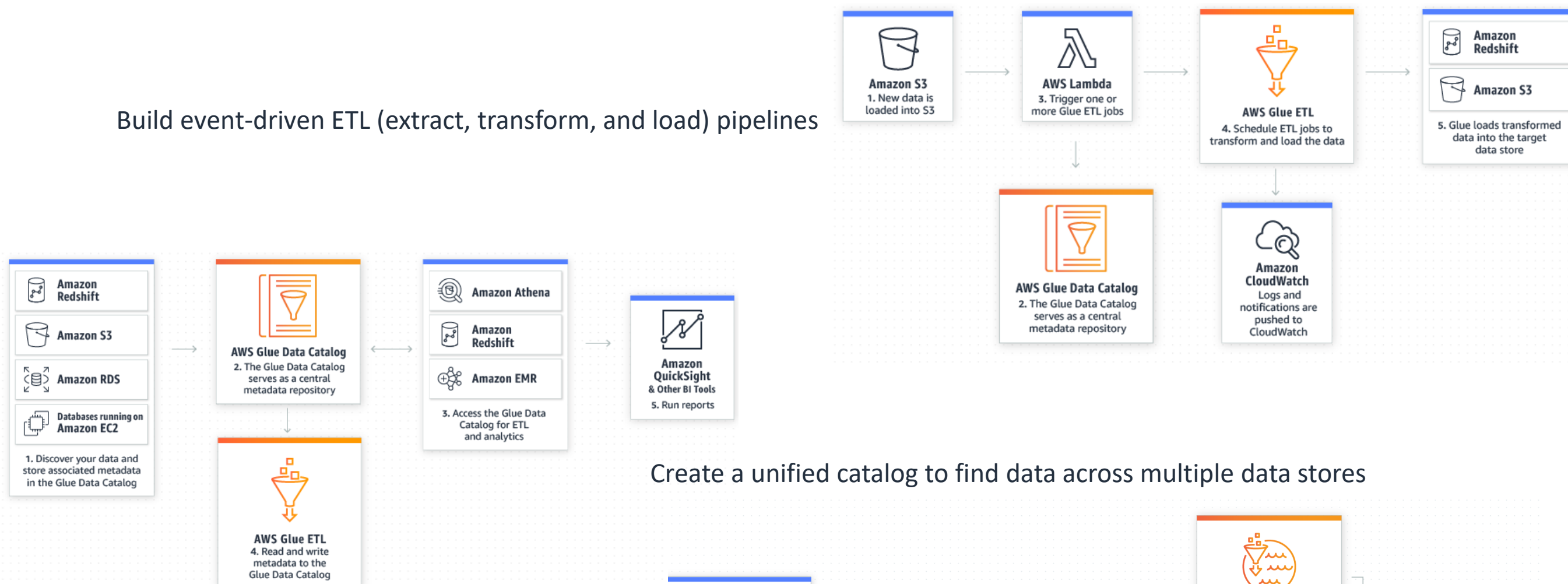
Glue – Main Use Cases





Glue – Use Cases

Build event-driven ETL (extract, transform, and load) pipelines



Create a unified catalog to find data across multiple data stores

Create, run, and monitor ETL jobs without coding



Other 2 common use cases: Explore data with self-service visual data preparation (AWS Glue DataBrew), Build materialized views to combine and replicate data (Glue Elastic Views in Beta).

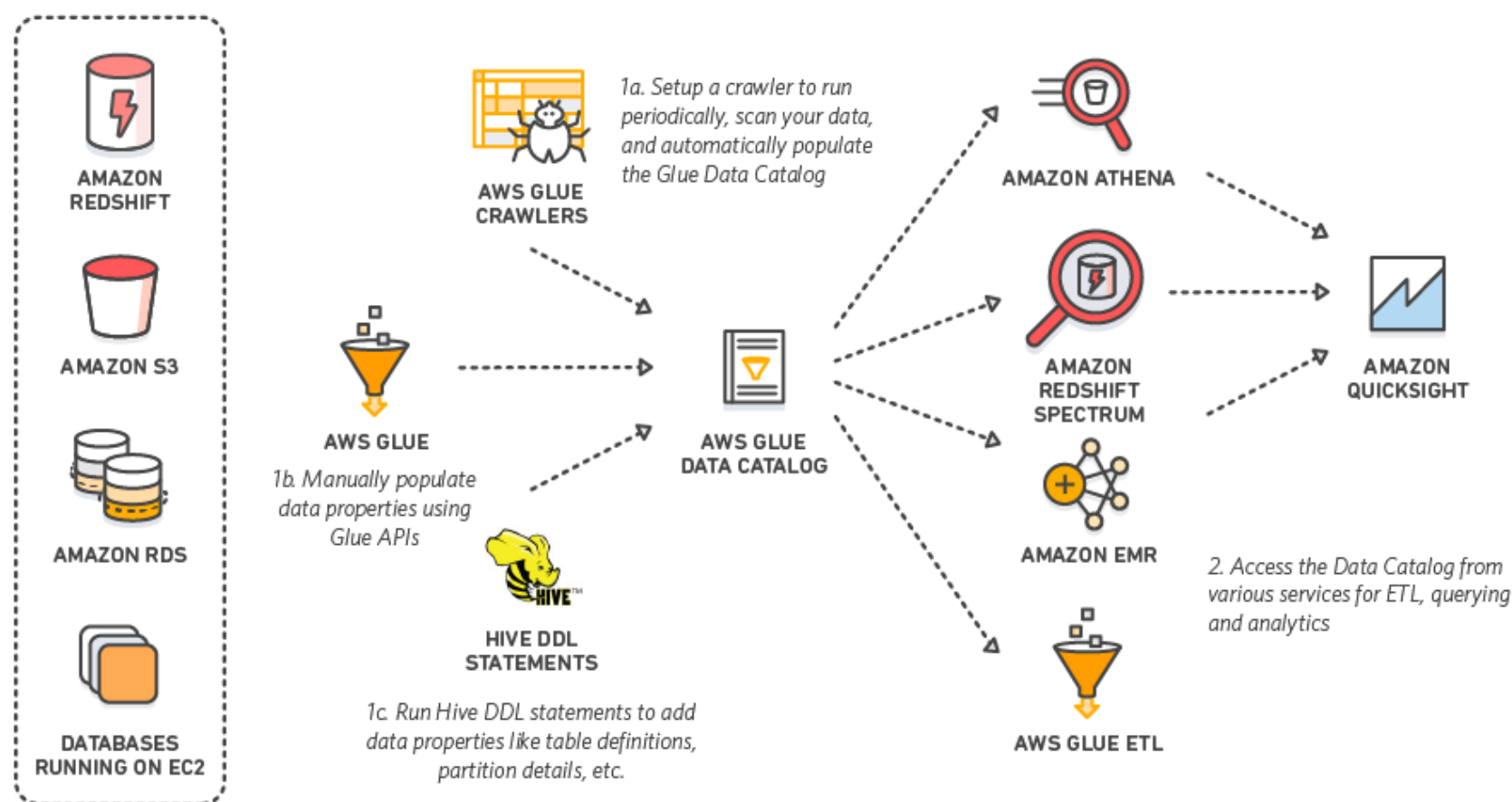


Query (ANSI SQL) service using AWS Glue Data Catalog as schema to several data sources. For native using S3 (Read and Write) and for Federated Query using Lambda Connectors: RDS, DynamoDB, Hbase, Cloudwatch Logs, etc.

Distributed query, based on Presto.

S3 Formats: CSV, JSON, Avro, Parquet, ORC.

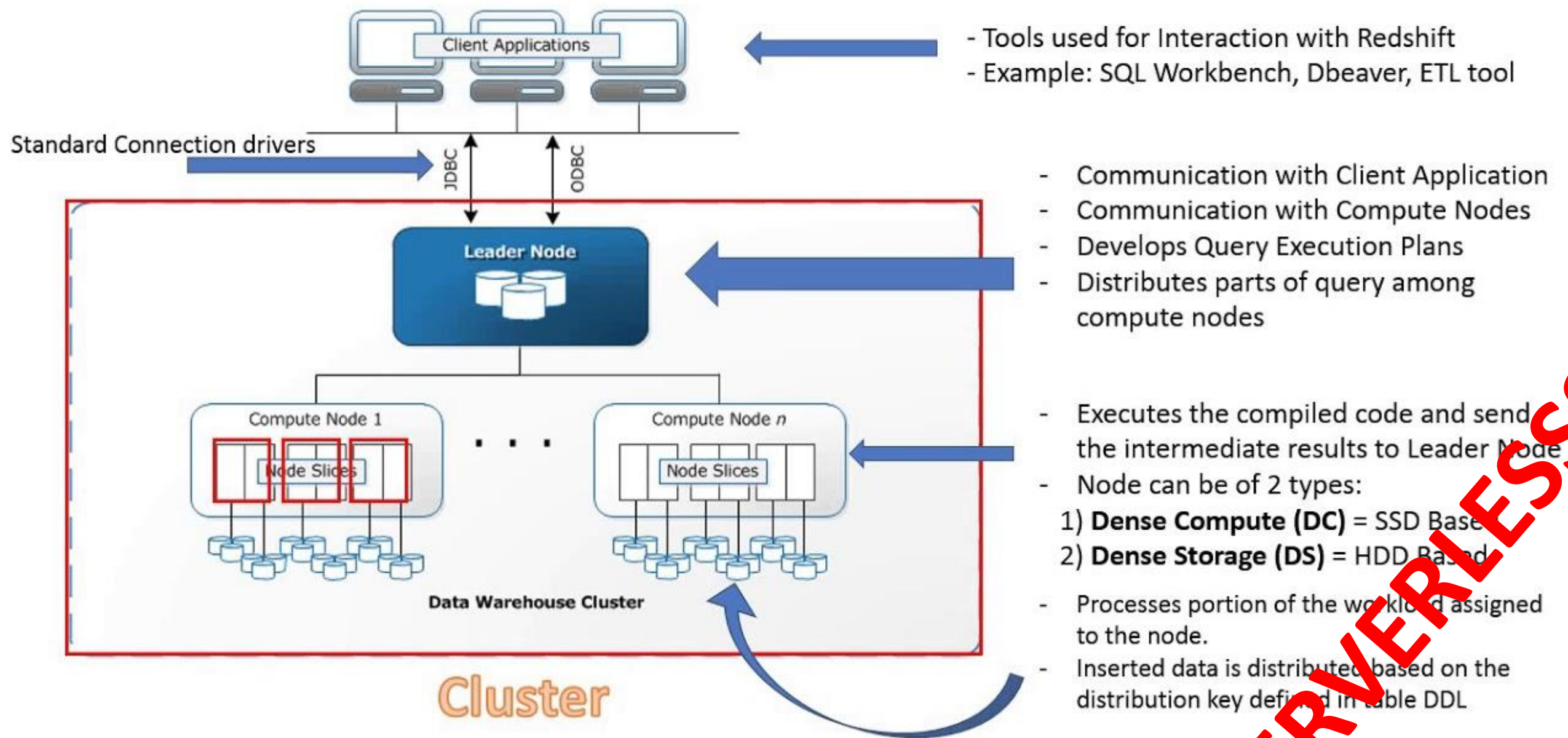
Serverless and Pay per Query Service (U\$5 per TB Scanned). Using compressed format is more cheaper.





Redshift is Data Warehouse (DWH) preferred by AWS. Exabyte Scale for multiples data sources: DB, DataLakes, etc.

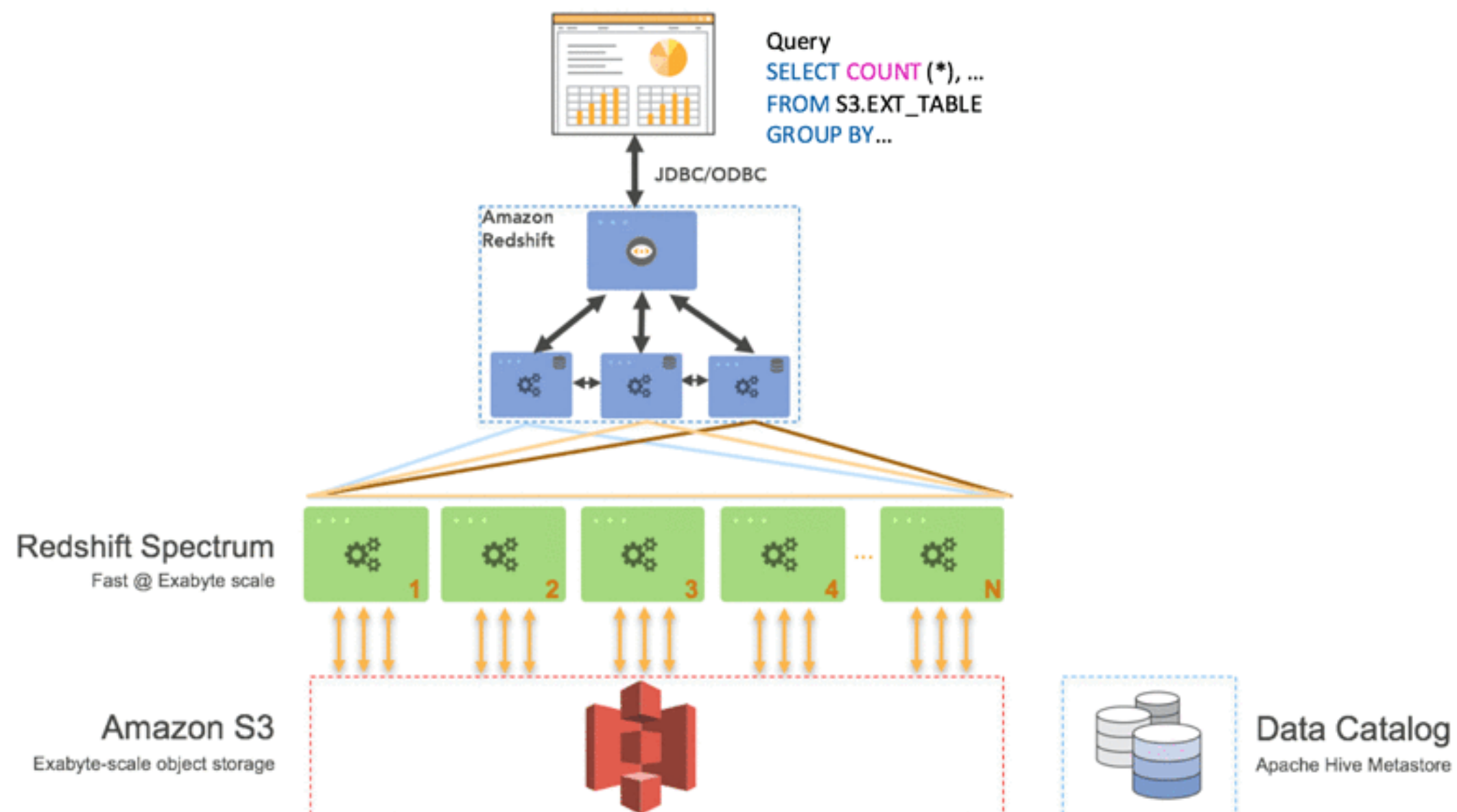
AWS Redshift - Architecture



NO SERVERLESS!!!



Redshift Spectrum





The Emerging Analytics Architecture

Storage



Amazon S3
Exabyte-scale Object Storage



AWS Glue Data Catalog
Hive-compatible Meta store

Serverless Compute



Amazon Kinesis Firehose
Real-Time Data Streaming



AWS Glue
ETL & Data Catalog



Amazon Redshift Spectrum
Fast @ Exabyte scale



AWS Lambda
Trigger-based Code Execution

Data Processing



Amazon EMR
Managed Hadoop Applications



Amazon Redshift
Petabyte-scale Data Warehousing



Amazon Athena
Interactive Query