



INNOVATE2018

ONLINE CONFERENCE



Serving Business Insights with the AWS Analytics Suite (Level 200)

Ekta Parashar, Solution Architect Manager

What to Expect from this Session

- AWS Toolkit for advanced analytics
- Understand Stakeholders
- Services for Stakeholders
- Q & A

Who Am I ?

- Solution Architect at AWS since 2015.
- Based in Mumbai, India.

Central Storage
*Secure, Cost Effective
Storage in S3*



S3

Data Ingestion

Get your data into S3
quickly and securely



Firehose



Direct Connect



Snowball



DMS



Central Storage

Secure, Cost Effective
Storage in S3



S3

Data Ingestion

Get your data into S3 quickly and securely



Firehose



Direct Connect



Snowball



DMS



Central Storage

Secure, Cost Effective Storage in S3



S3



Processing & Analytics

Use predictive and prescriptive analytics to gain better understanding



Athena



QuickSight



EMR



Redshift

Catalog & Search

Capture, Access, and Search Metadata



Glue



Macie



DynamoDB



Amazon ES

Data Ingestion

Get your data into S3 quickly and securely



Firehose



Direct Connect



Snowball



DMS

Central Storage

*Secure, Cost Effective
Storage in S3*



S3

Processing & Analytics

*Use predictive and prescriptive
analytics to gain better understanding*



Athena



QuickSight



EMR



Redshift

Catalog & Search

Capture, Access, and Search Metadata



Glue



Macie



DynamoDB



Amazon ES

Access & User Interface

Give your users easy & secure access



API Gateway



IAM



Cognito

Data Ingestion

Get your data into S3 quickly and securely



Firehose



Direct Connect



Snowball



DMS

Central Storage

*Secure, Cost Effective
Storage in S3*



S3

Processing & Analytics

*Use predictive and prescriptive
analytics to gain better understanding*



Athena



QuickSight



EMR



Redshift



Catalog & Search

Capture, Access, and Search Metadata



Glue



Macie



DynamoDB



Amazon ES

Access & User Interface

Give your users easy & secure access



API Gateway



IAM



Cognito

Data Ingestion

Get your data into S3 quickly and securely



Firehose



Direct Connect



Snowball



DMS

Central Storage

Secure, Cost Effective Storage in S3



S3

Processing & Analytics

Use predictive and prescriptive analytics to gain better understanding



Athena



Quickstart



EMR



Redshift

Protect & Secure

Use entitlements to ensure data is secure and users identities are verified



Security Token Service



Cloudwatch



Cloudtrail

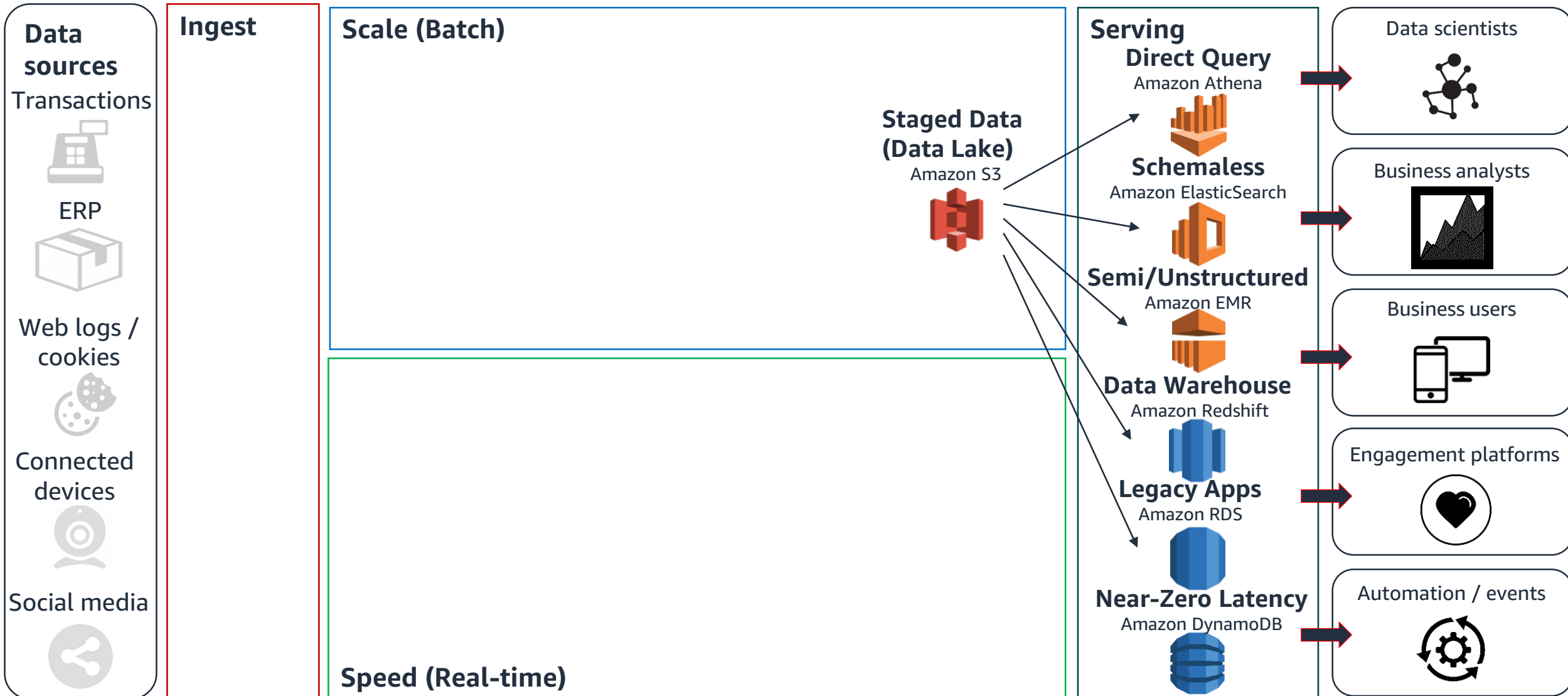


KMS



Modern data architecture

Insights to enhance business applications, new digital services

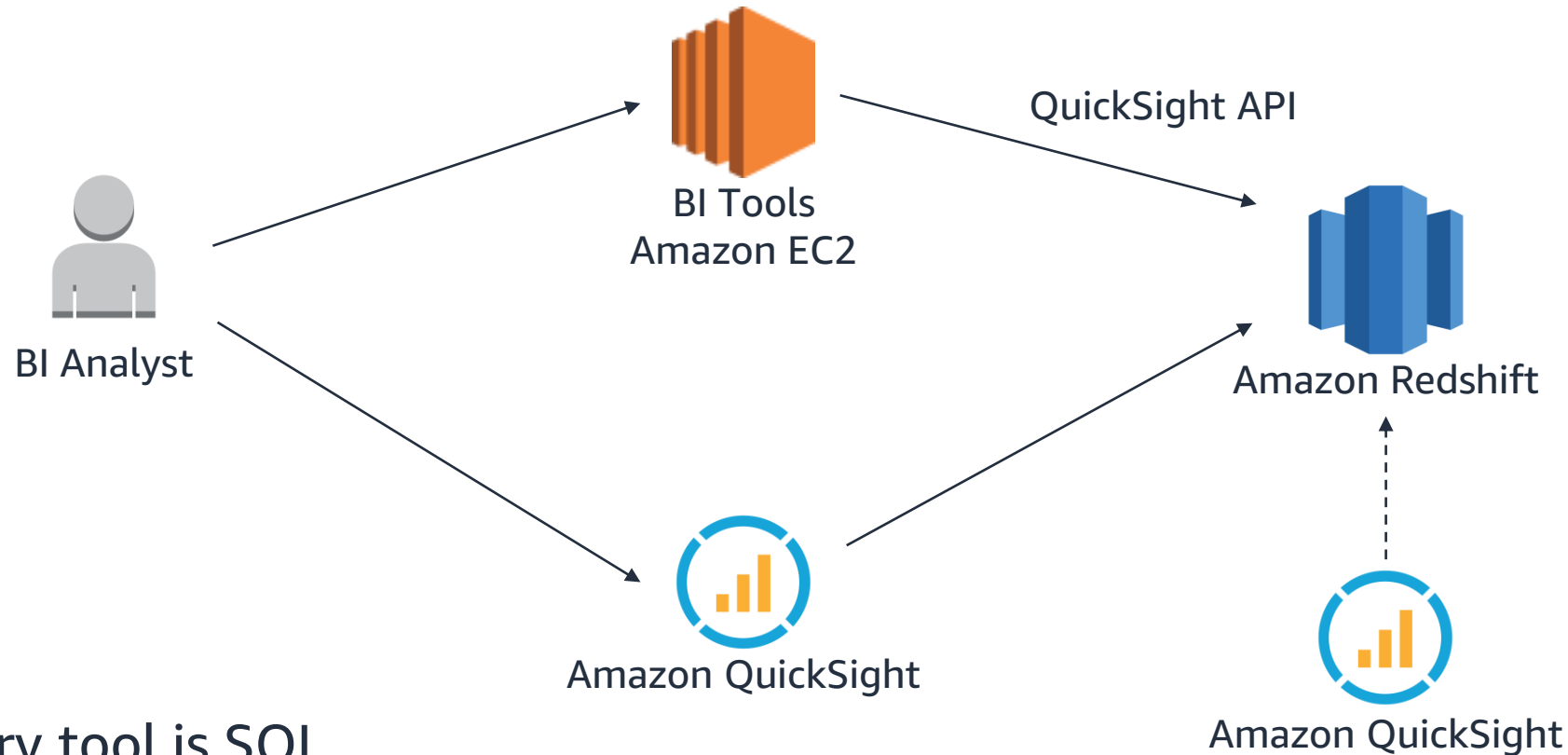


Match toolset to right persona

- **Business Intelligence (BI) Analyst**
 - Primary tool is SQL
 - Historical data resides in data warehouse such as Redshift, Spectrum, QuickSight
- **Data Scientist** – Uses programmatic languages such as R or python, explore data using Athena
- **Data Engineer**
 - Familiar with Hadoop and Spark
 - Wants to access data in Redshift in diverse ways

BI Analyst

BI Analyst with existing BI Tools



- Primary tool is SQL
- Data is largely structured with well known data sources
- Primary concern is fast, consistent performance
- Need to extend SQL with custom functions

Amazon Redshift system architecture

Leader node

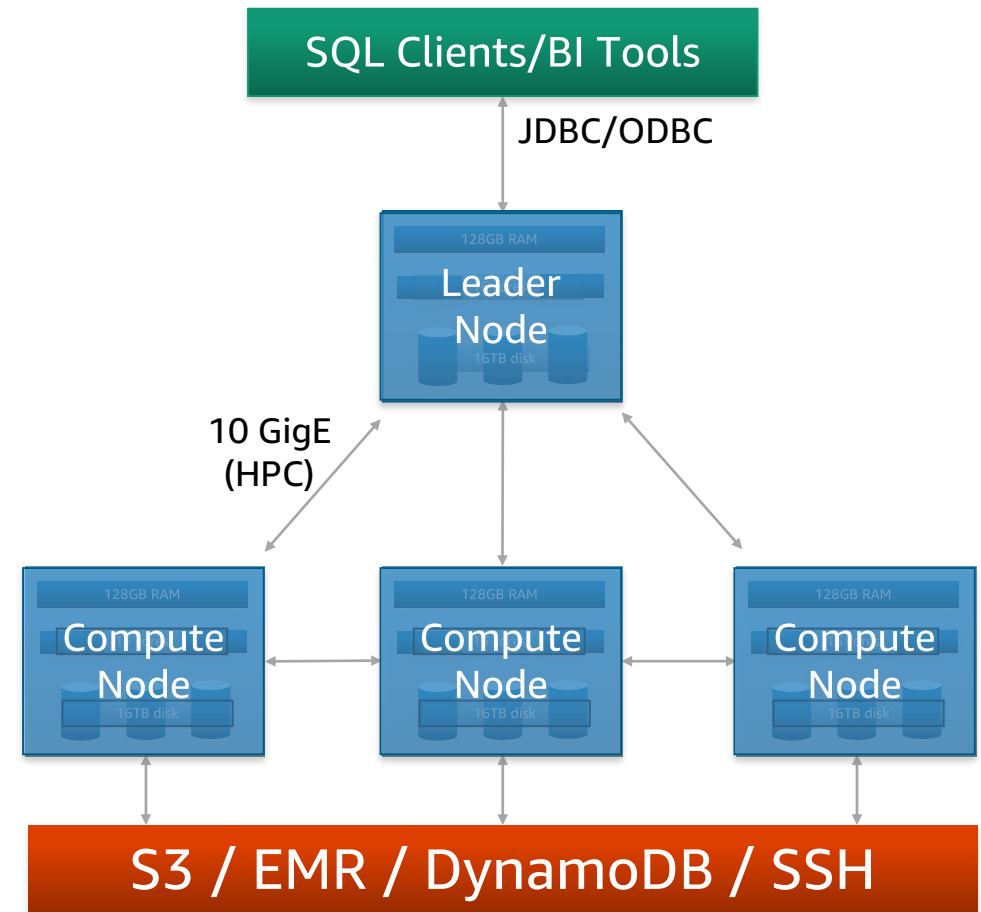
- SQL endpoint
- Stores metadata
- Coordinates query execution

Compute nodes

- Local, columnar storage
- Execute queries in parallel
- Load, backup, restore via Amazon S3; load from Amazon DynamoDB, Amazon EMR, or SSH

Two hardware platforms

- Optimized for data processing
- DS2: HDD; scale from 2TB to 2PB
- DC1: SSD; scale from 160GB to 356TB

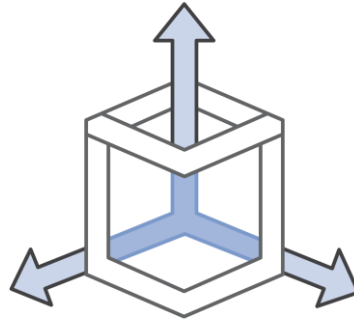


Amazon Redshift Spectrum

Run SQL queries directly against data in S3 using thousands of nodes



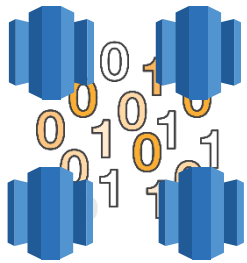
Fast @ Exabyte scale



Elastic & highly available



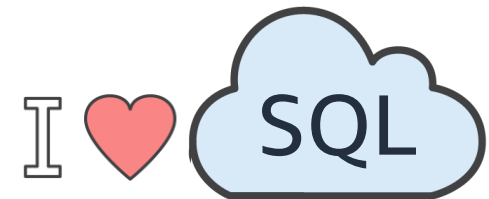
On-demand, pay-per-query



High concurrency: Multiple clusters access same data

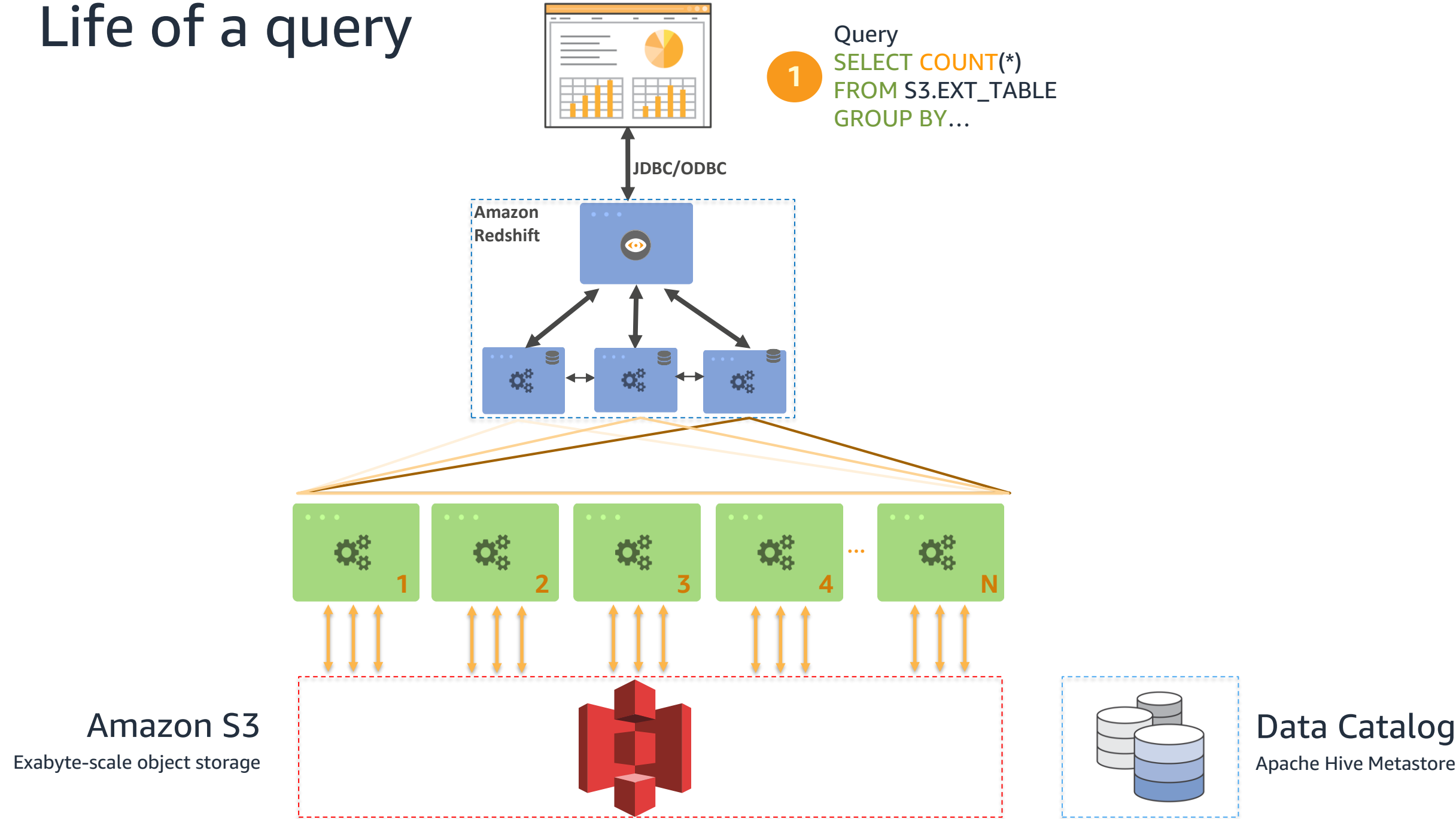


No ETL: Query data in-place using open file formats

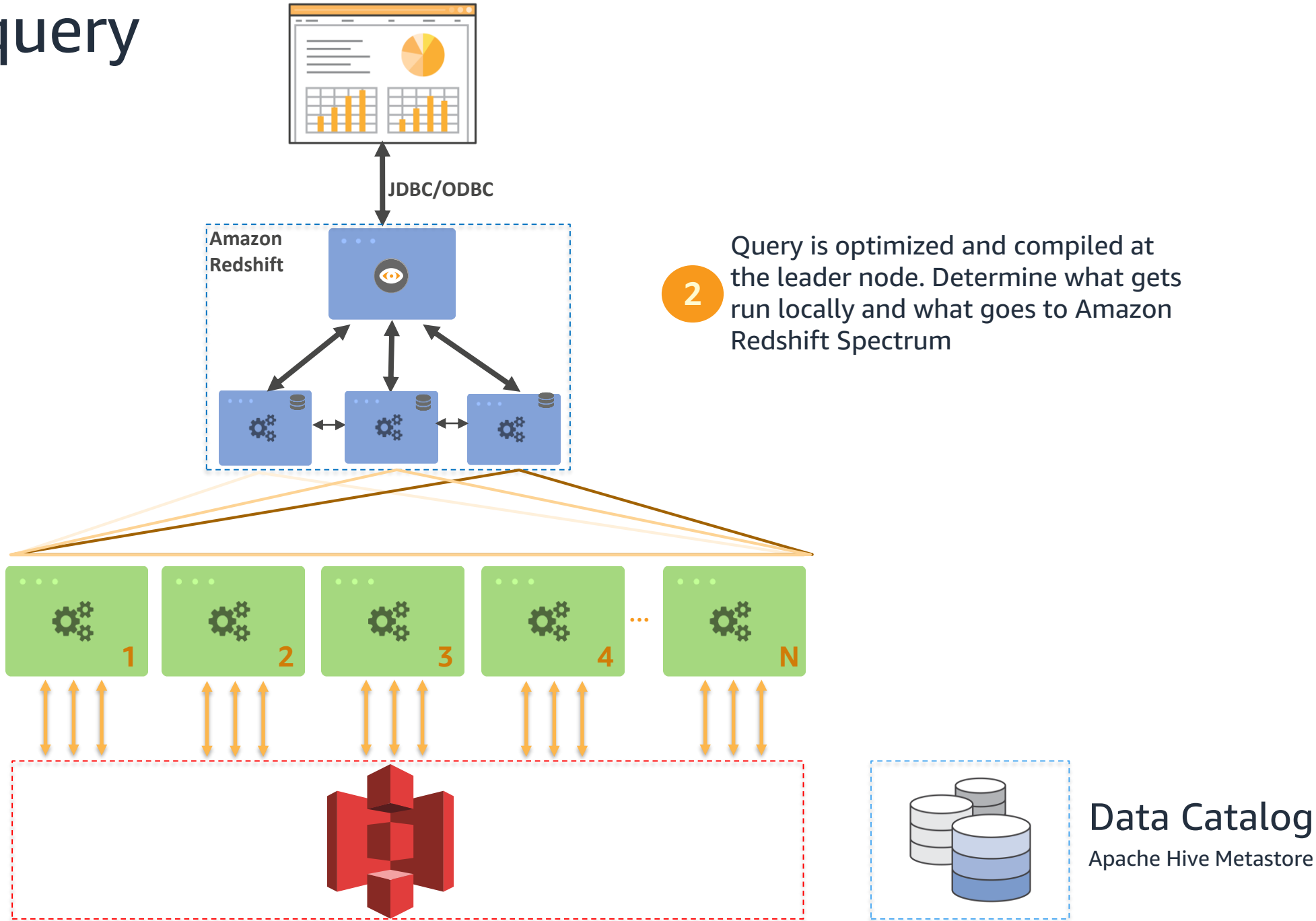


Full Amazon Redshift SQL support

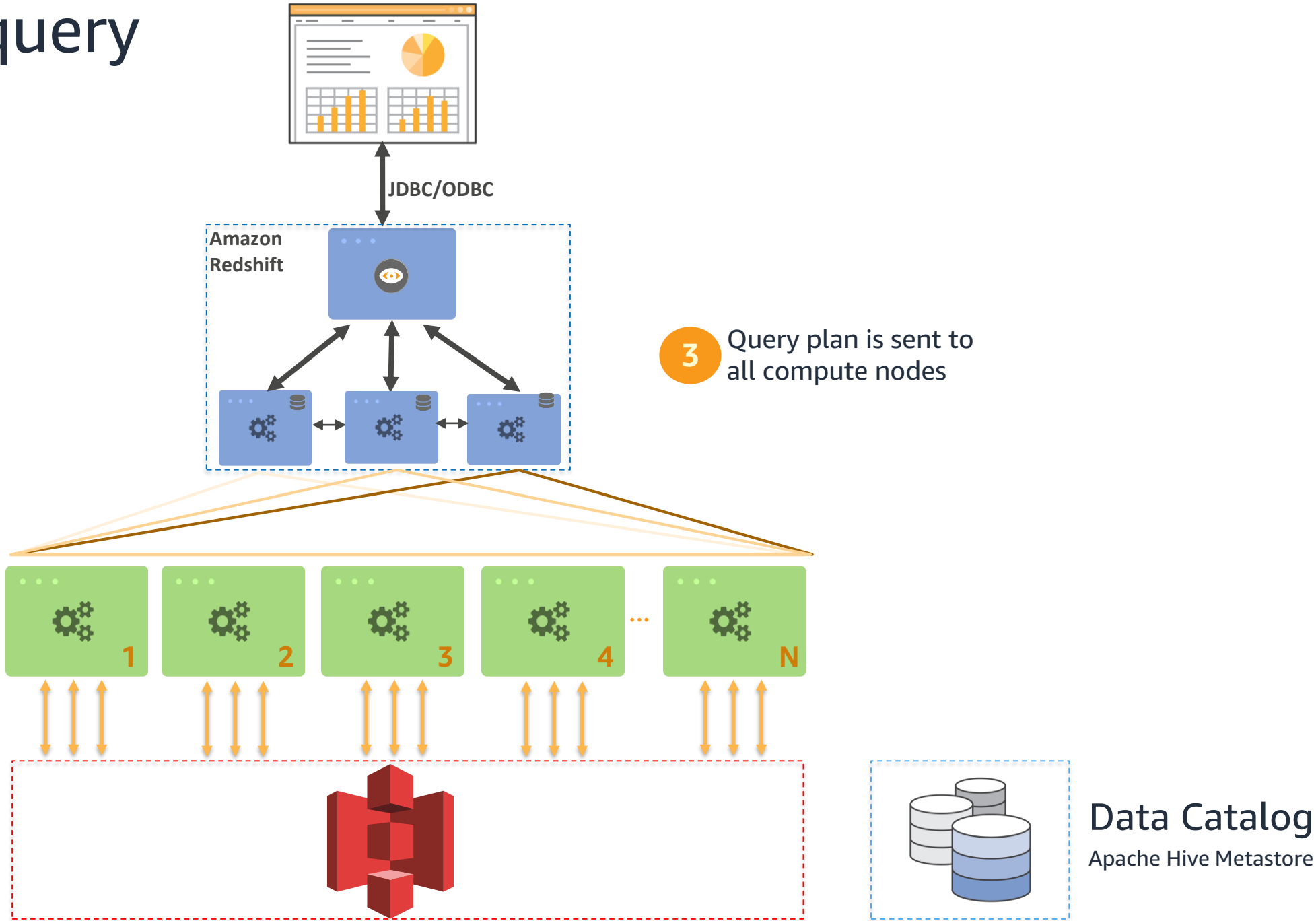
Life of a query



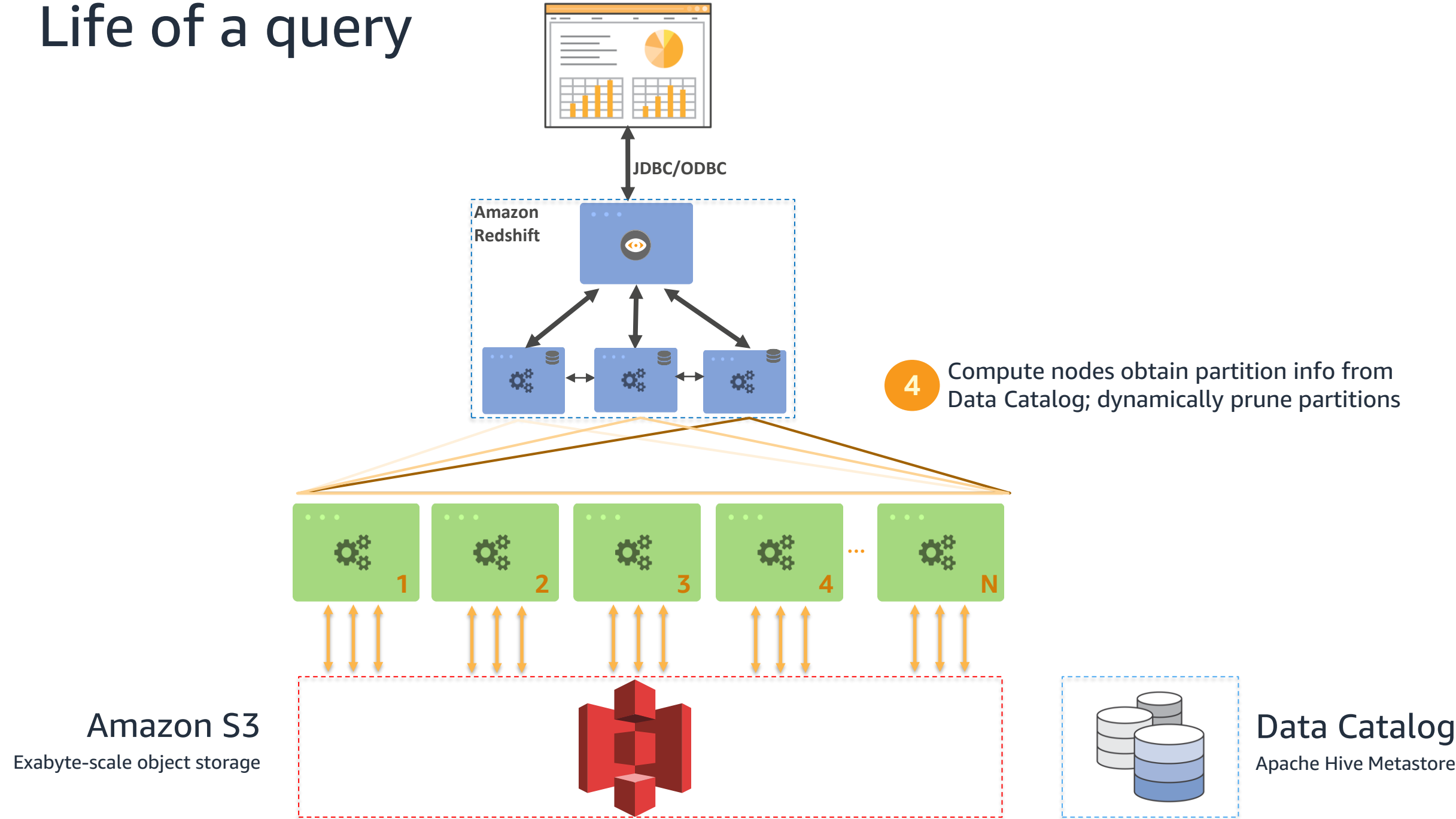
Life of a query



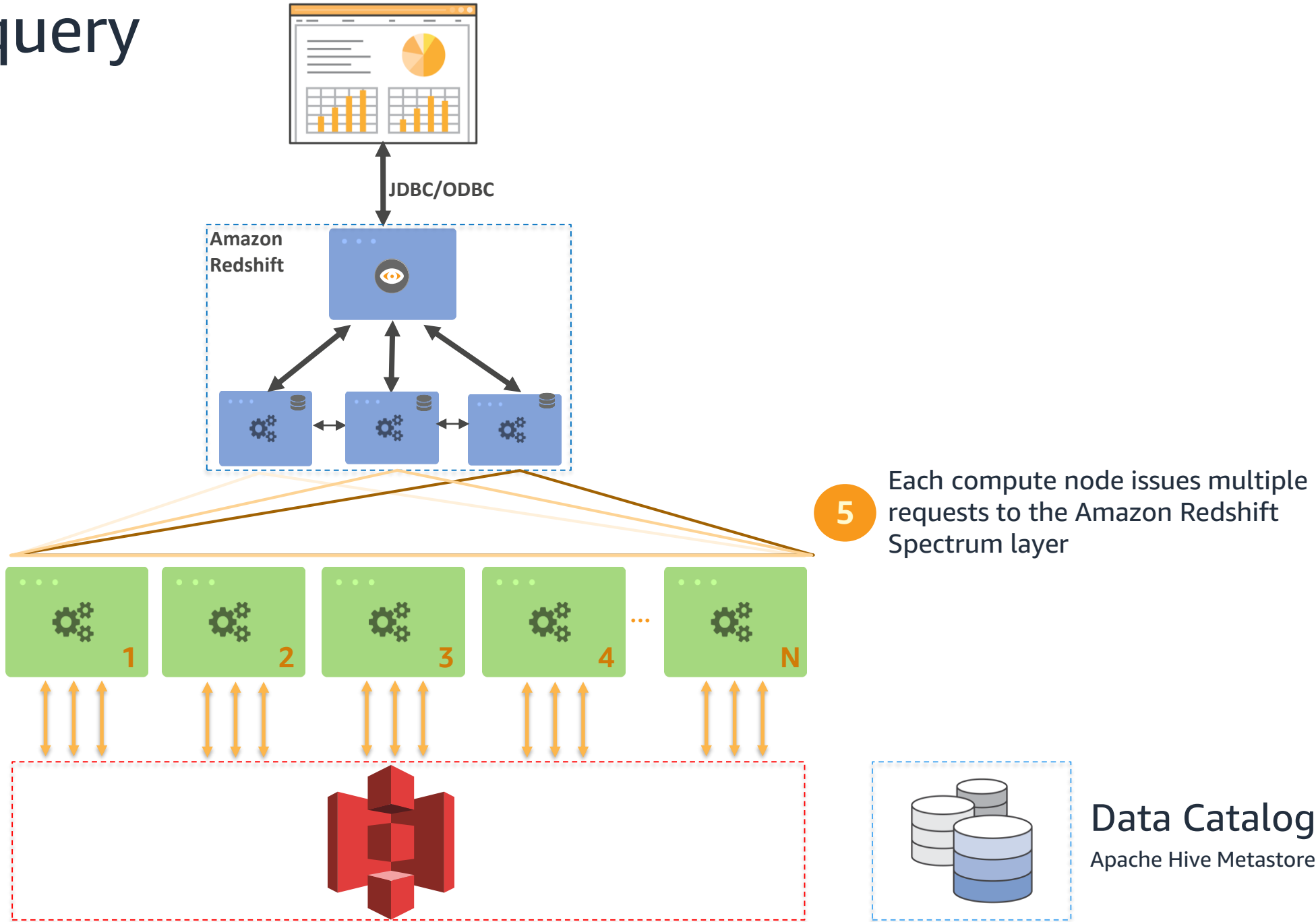
Life of a query



Life of a query



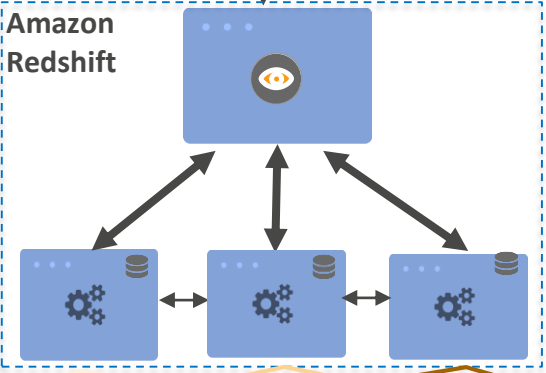
Life of a query



Life of a query



JDBC/ODBC



9 Result is sent back to client

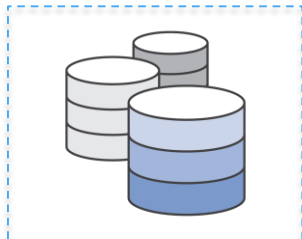
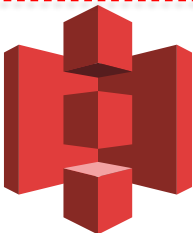
8 Final aggregations and joins with local Amazon Redshift tables done in-cluster

7 Amazon Redshift Spectrum projects, filters, and aggregates



6 Amazon Redshift Spectrum nodes scan your S3 data

Amazon S3
Exabyte-scale object storage



Data Catalog
Apache Hive Metastore

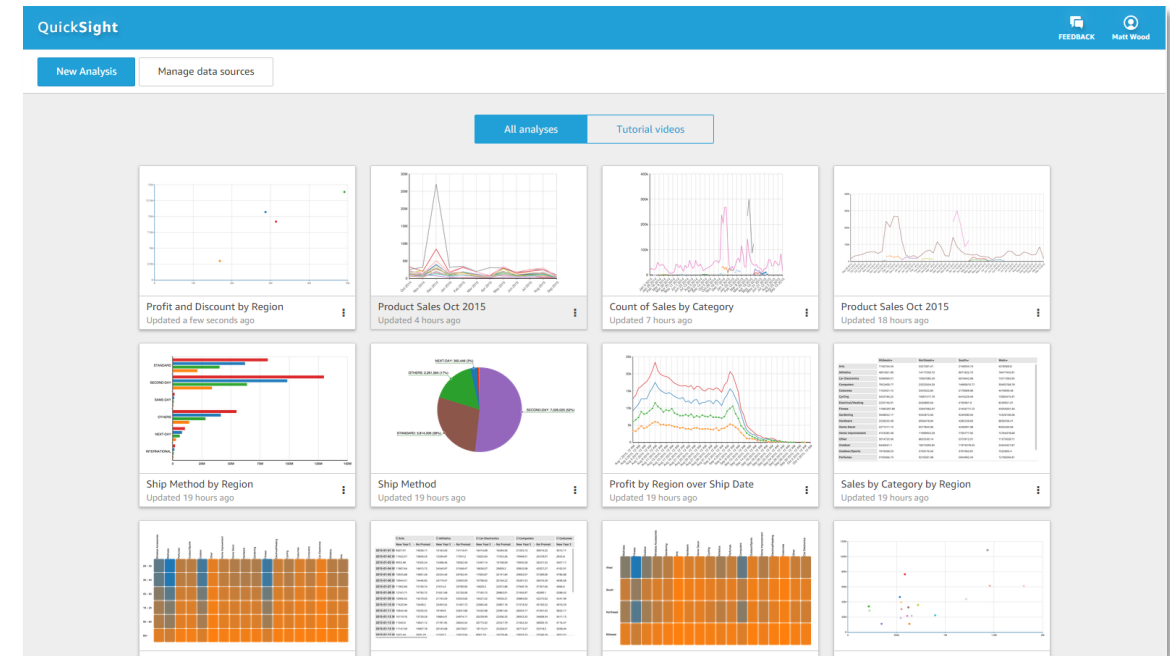
Amazon QuickSight

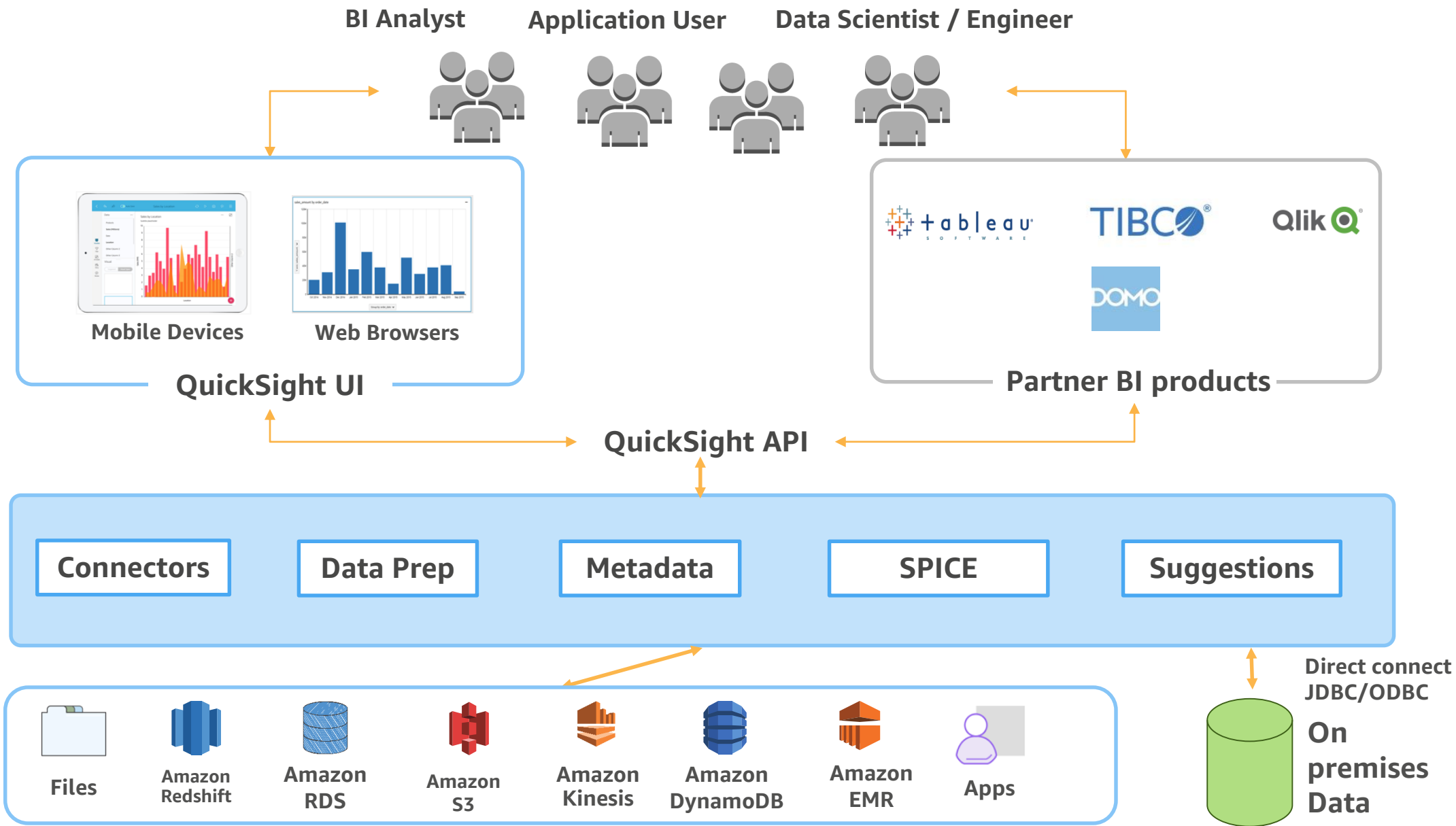
QuickSight – fast access to all your analytics



Business user

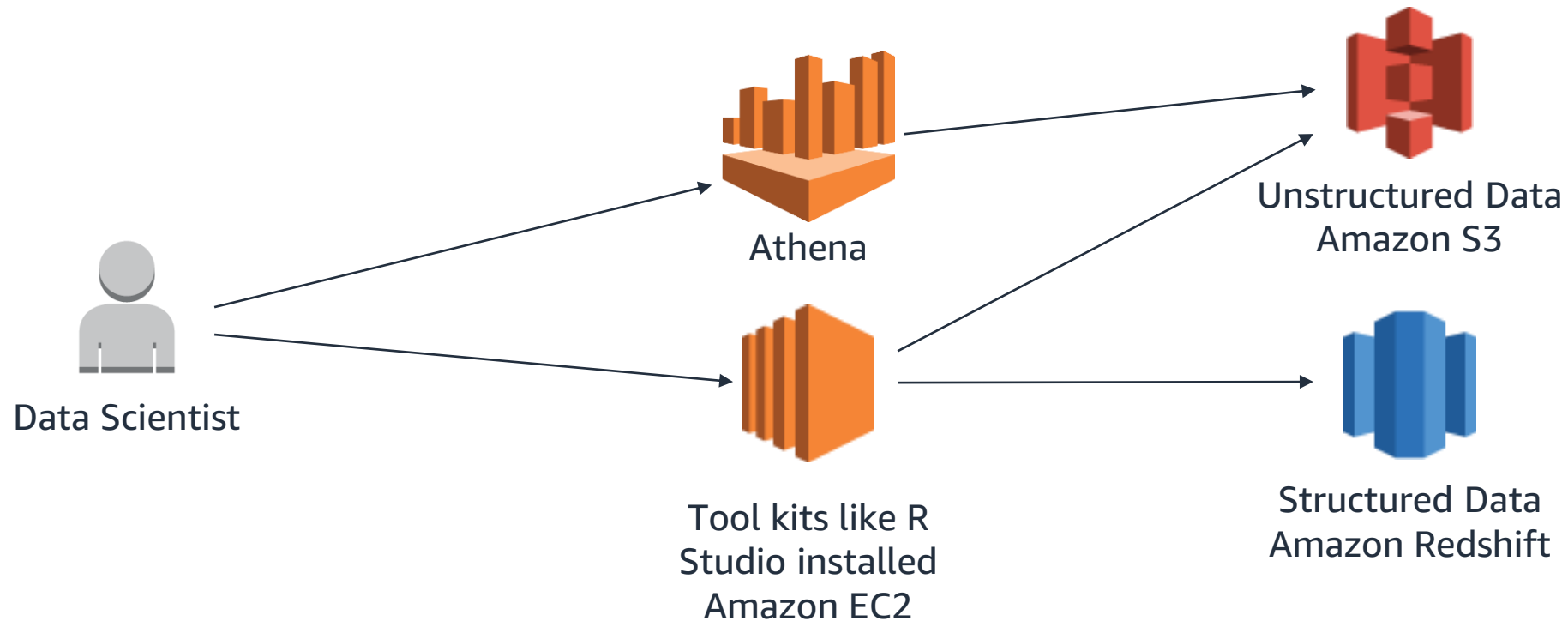
Sign-in





Data Scientist

Data Scientist with existing toolsets



- Work with unstructured datasets
- Use existing toolsets to connect to Redshift

What is R?

Open source programming language and software environment designed for statistical computing, data analysis, and visualization



Open source IDE for R



Shiny Server - Visualization R package for creating interactive dashboards

The New York Times

Search All NYTimes.com

Business Computing

WORLD U.S. N.Y./REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION ARTS STYLE TRAVEL JOBS REAL ESTATE AUTOS

Data Analysts Captivated by R's Power

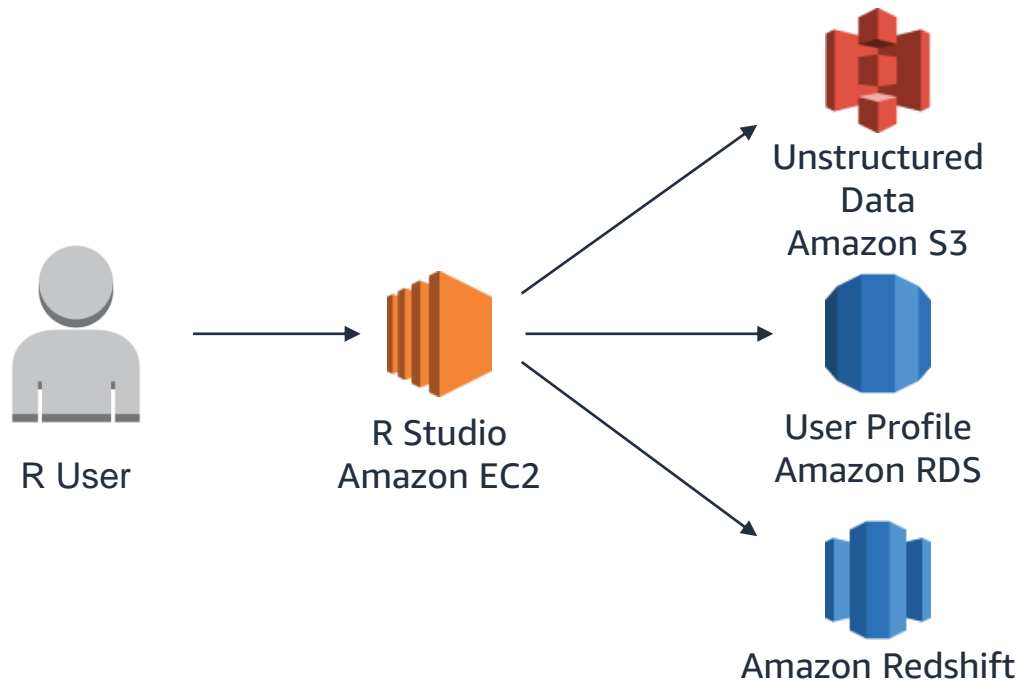


Left, Stuart Isett for The New York Times; right, Kieran Scott for The New York Times

R first appeared in 1996, when the statistics professors Robert Gentleman, left, and Ross Ihaka released the code as a free software package.

By ASHLEE VANCE
Published: January 6, 2009

Querying Redshift with R Packages



- **RJDBC** – supports SQL queries
- ***dplyr*** – Uses R code for data analysis
- **RPostgreSQL** – R compliant driver or Database Interface (DBI)

Connecting R with Amazon Redshift blog post:

<https://blogs.aws.amazon.com/bigdata/post/Tx1G8828SPGX3PK/Connecting-R-with-Amazon-Redshift>

Introducing Amazon Athena

Amazon Athena is an **interactive query service** that makes it easy to analyze data directly from Amazon S3 using Standard SQL

Query Data Directly from Amazon S3

- No loading of data
- Query data in its raw format
 - Text, CSV, JSON, weblogs, AWS service logs
 - Convert to an optimized form like ORC or Parquet for the best performance and lowest cost
- No ETL required
- Stream data from directly from Amazon S3
- Take advantage of Amazon S3 durability and availability

Use ANSI SQL

- Start writing ANSI SQL
- Support for complex joins, nested queries & window functions
- Support for complex data types (arrays, structs)
- Support for partitioning of data by any key
 - (date, time, custom keys)
 - e.g., Year, Month, Day, Hour or Customer Key, Date

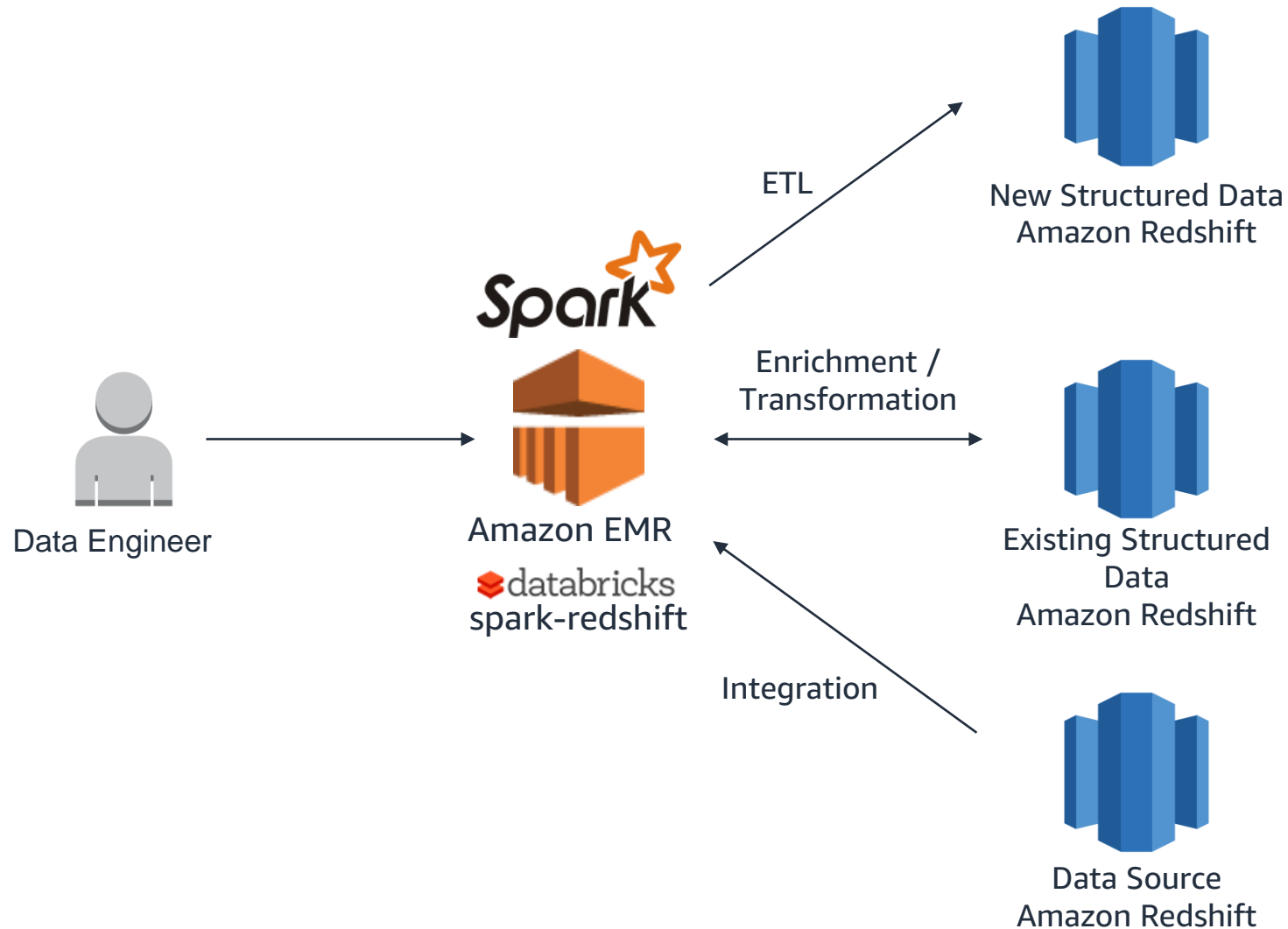
```
1 WITH q21_tmp1_cached AS
2 (SELECT l_orderkey,
3    count(DISTINCT l_supkey) AS count_supkey,
4    max(l_supkey) AS max_supkey
5 FROM lineitem_parq
6 WHERE l_orderkey IS NOT NULL
7 GROUP BY l_orderkey),
8 q21_tmp2_cached AS
9 (SELECT l_orderkey,
10    count(DISTINCT l_supkey) AS count_supkey,
11    max(l_supkey) AS max_supkey
12 FROM lineitem_parq
13 WHERE l_receiptdate > l_commitdate
14 AND l_orderkey IS NOT NULL
15 GROUP BY l_orderkey)
16 SELECT s_name,
17    count(1) AS numwait
18 FROM
19 (SELECT s_name
20 FROM
21 (SELECT s_name,
22    t2.l_orderkey,
23    l_supkey,
24    count_supkey,
25    max_supkey
26 FROM q21_tmp2_cached t2
27 RIGHT OUTER JOIN
28 (SELECT s_name,
29    l_orderkey,
30    l_supkey
31 FROM
32 (SELECT s_name,
33    t1.l_orderkey,
34    l_supkey,
35    count_supkey,
36    max_supkey
37 FROM q21_tmp1_cached t1
38 JOIN
39 (SELECT s_name,
40    l_orderkey,
41    l_supkey
42 FROM orders_parq o
43 JOIN
44 (SELECT s_name,
45    l_orderkey,
46    l_supkey
47 FROM nation_parq n
48 JOIN supplier s ON s.s_nationkey = n.n_nationkey
49 AND n.n_name = 'SAUDI ARABIA'
50 JOIN lineitem_parq l ON s.s_supkey = l.l_supkey
51 WHERE l.l_receiptdate > l.l_commitdate
52 AND l.l_orderkey IS NOT NULL) t1 ON o.o_orderkey = t1.l_orderkey
53 AND o.o_orderstatus = 'F') t2 ON t2.l_orderkey = t1.l_orderkey) a
54 WHERE (count_supkey > 1)
55 OR ((count_supkey=1)
56 AND (l_supkey <= max_supkey))) t3 ON t3.l_orderkey = t2.l_orderkey) b
57 WHERE (count_supkey IS NULL)
58 OR ((count_supkey=1)
59 AND (l_supkey = max_supkey))) c
60 GROUP BY s_name
61 ORDER BY numwait DESC,
62    s_name LIMIT 100;
```

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

Data Engineer

Data Engineer familiar with Hadoop and Spark





Spark SQL +
DataFrames

Streaming

MLlib
Machine Learning

GraphX
Graph Computation

Spark Core API

R

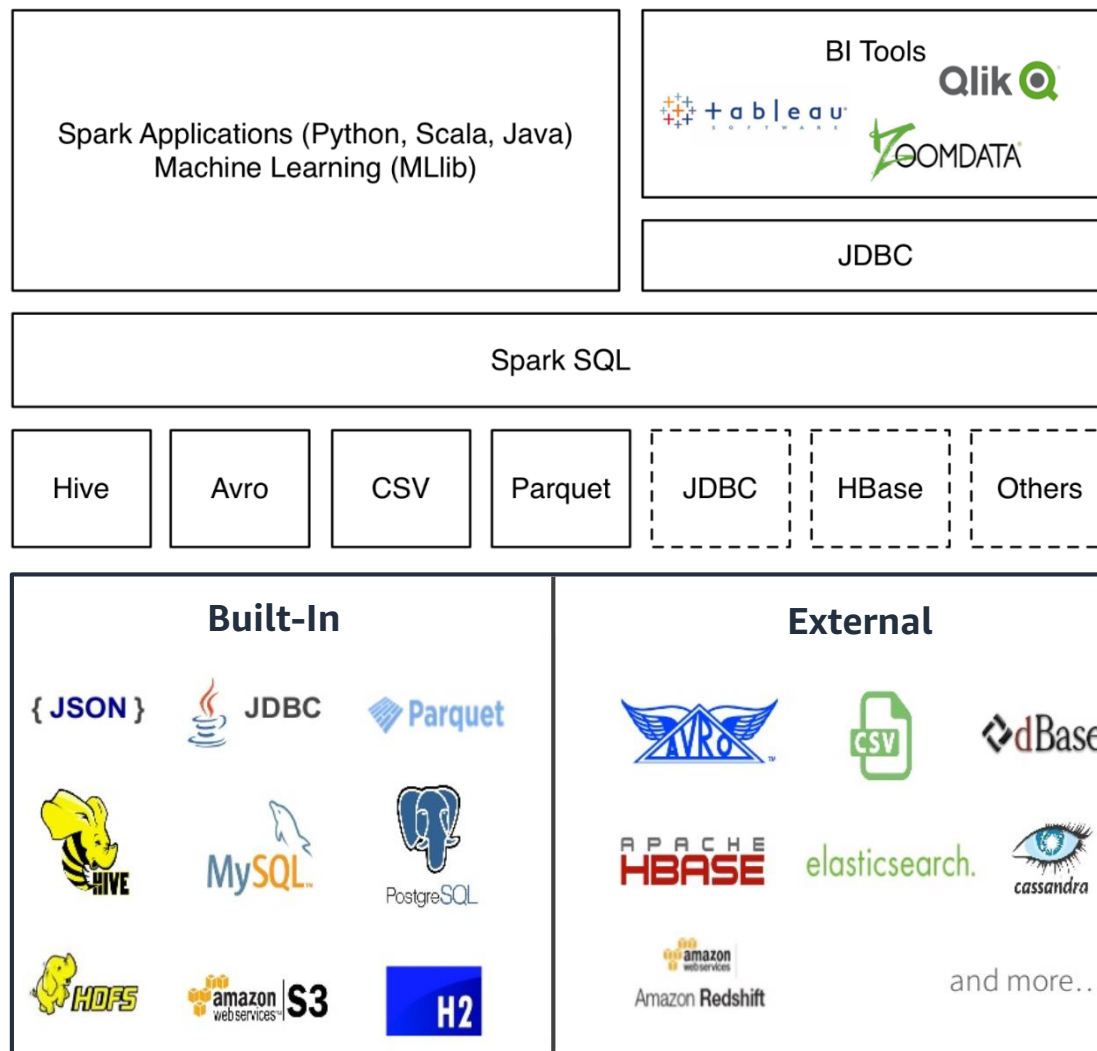
SQL

Python

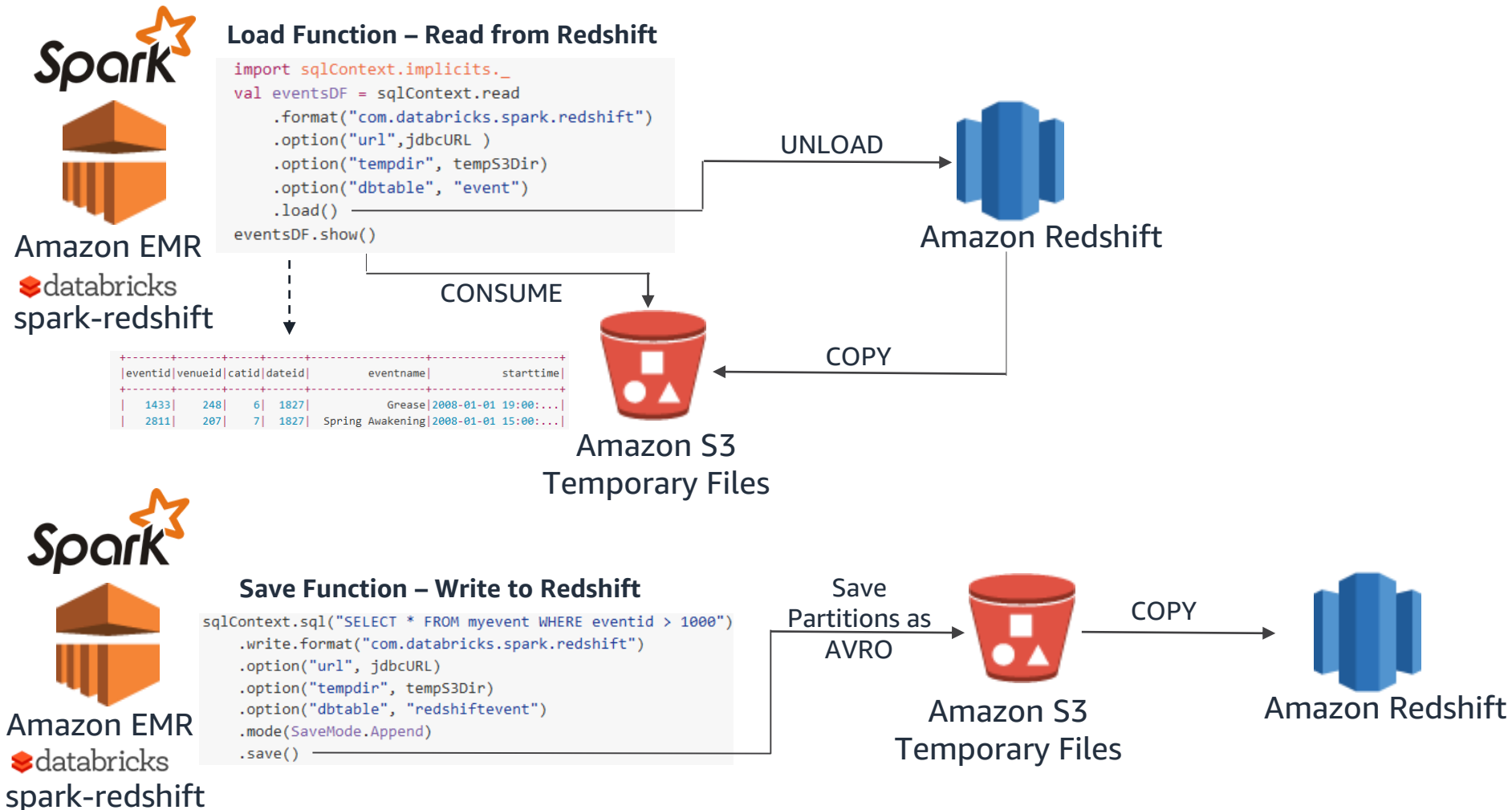
Scala

Java

Spark SQL Data Sources API



Manipulating Redshift with spark-redshift



Databricks Tutorial: <https://github.com/databricks/spark-redshift/tree/master/tutorial>

Learn from AWS experts. Advance your skills and knowledge. Build your future in the AWS Cloud.



Digital Training

Free, self-paced online courses built by AWS experts



Classroom Training

Classes taught by accredited AWS instructors



AWS Certification

Exams to validate expertise with an industry-recognized credential

Ready to begin building your cloud skills?
Get started at: <https://www.aws.training/>

With deep expertise on AWS, APN Partners can help your organization at any stage of your Cloud Adoption Journey.



AWS Managed Service Providers

APN Consulting Partners who are skilled at cloud infrastructure and application migration, and offer proactive management of their customer's environment.



AWS Competency Partners

APN Partners who have demonstrated technical proficiency and proven customer success in specialized solution areas.



AWS Marketplace

A digital catalog with thousands of software listings from independent software vendors that make it easy to find, test, buy, and deploy software that runs on AWS.



AWS Service Delivery Partners

APN Partners with a track record of delivering specific AWS services to customers.

Ready to get started with an APN Partner?
Find a partner: <https://aws.amazon.com/partners/find/>
Learn more at the AWS Partner Network Booth

Thank You for Attending AWS Innovate

We hope you found it interesting! A kind reminder to **complete the survey.**

Let us know what you thought of today's event and how we can improve the event experience for you in the future.



aws-apac-marketing@amazon.com



twitter.com/AWSCloud



facebook.com/AmazonWebServices



youtube.com/user/AmazonWebServices



slideshare.net/AmazonWebServices



twitch.tv/aws