AWS
re:Invent

ANT335-R2

# How to scale data analytics with Amazon Redshift

ANT-335-R2, Thursday, Dec 5, 2019, 12:15 PM - 1:15 PM — Venetian, Level 3, Lido 3005

**Vinay Shukla**
Principal Product Manager, AWS

**Maor Kleider**
Principal Product Manager, AWS

**Jonathan Burket**
Senior Software Engineer, Duolingo

AWS re:Invent

aws

# Large trends in data

**Migrations to the cloud**

**Exponential growth of event data**
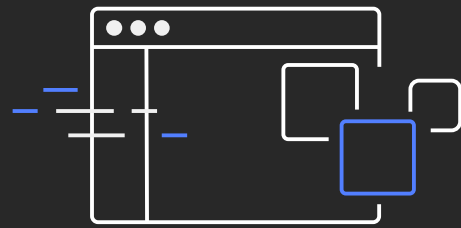
Data
010010010
01010001
100010100

**End-to-end insights from analyzing all your data**

# Challenges of data analytics at scale

**Data volume, variety, velocity**

**Performance, concurrency**

**Multiple analytics needs**

**Security, governance**

**Increasingly costly, inflexible**

# Amazon Redshift architecture

**Massively parallel, shared-nothing architecture**

## Leader node

SQL endpoint, stores metadata

Coordinates parallel SQL processing

Free for any cluster with two or more nodes

## Compute nodes
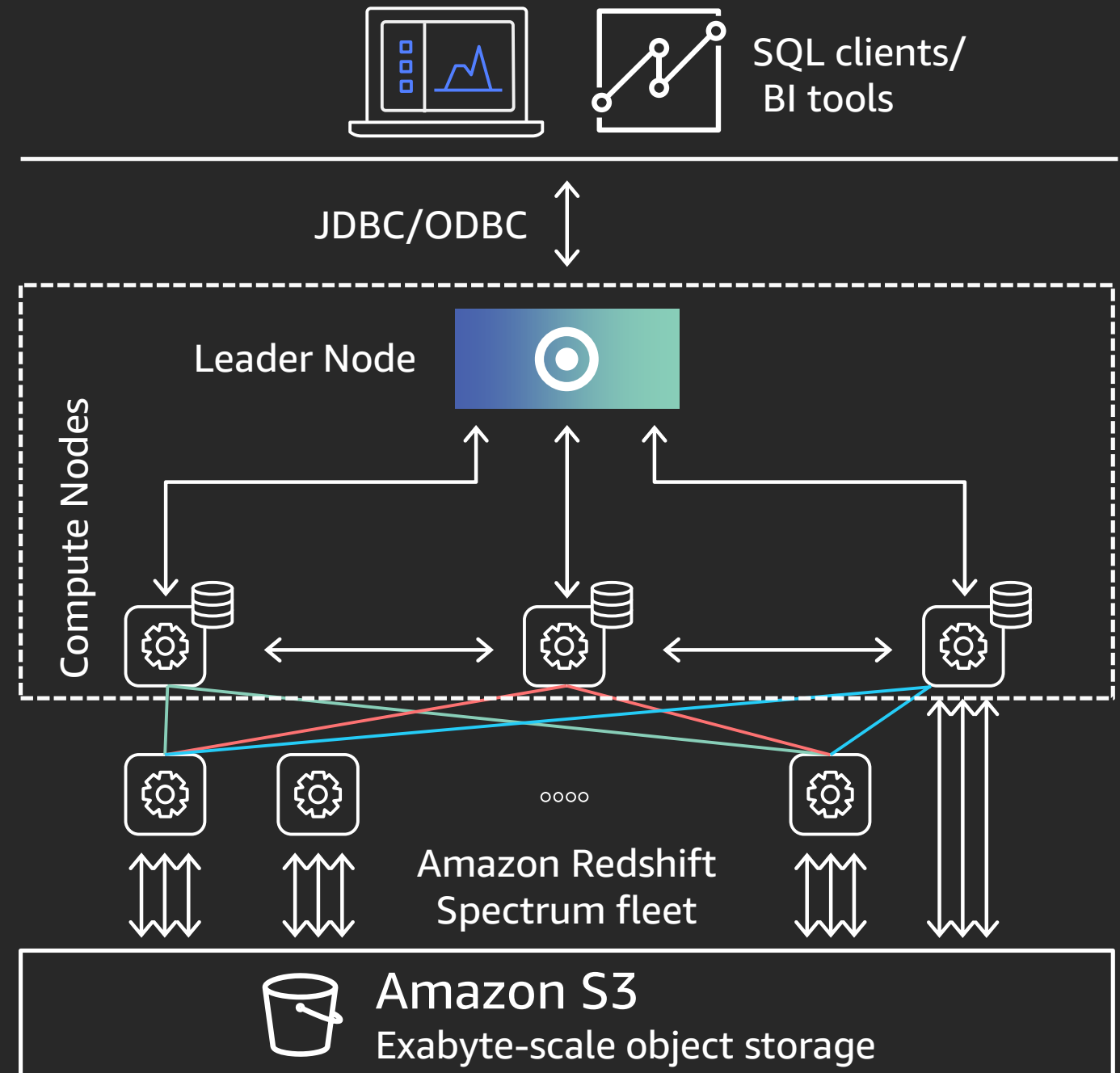
Local, columnar storage

Executes queries in parallel

Load, backup, restore

## Amazon Redshift Spectrum nodes

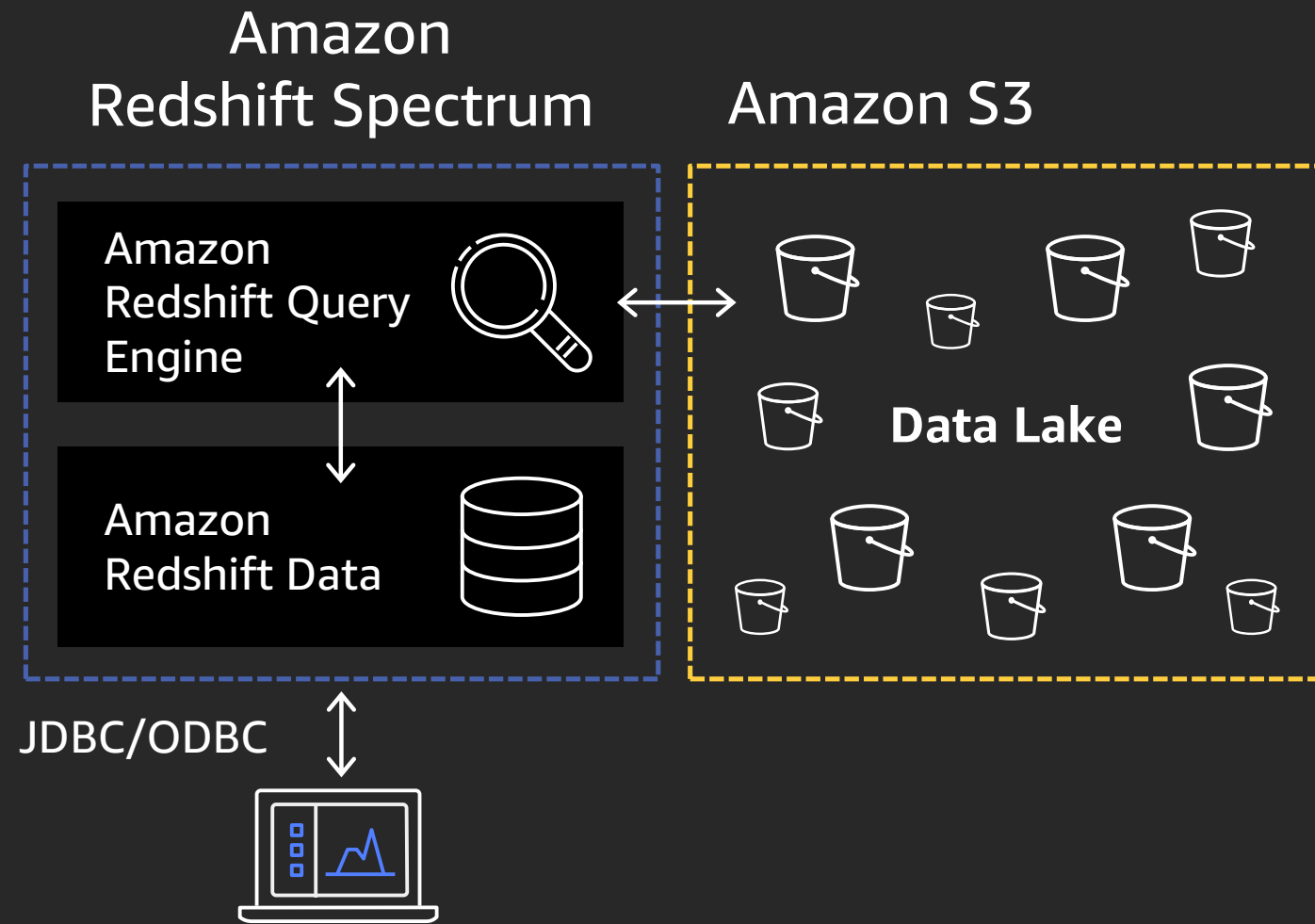Serverless, not managed by customer, bring power proportional to cluster slices

Execute queries directly against data lake

SQL clients/
BI tools

JDBC/ODBC

Compute Nodes

Leader Node

Amazon Redshift
Spectrum fleet

Amazon S3
Exabyte-scale object storage

# Challenges with growing data: volume, variety, velocity

aws

# Data at any scale: Query all your data
## Unified view: Local storage and Amazon S3 Data Lake

**Amazon Redshift Spectrum**

Amazon Redshift Query Engine

Amazon Redshift Data

**Amazon S3**

Data Lake

JDBC/ODBC

Directly query exabytes in S3

No data loading, eliminate ingestion time

Unified view of data across Amazon Redshift and S3

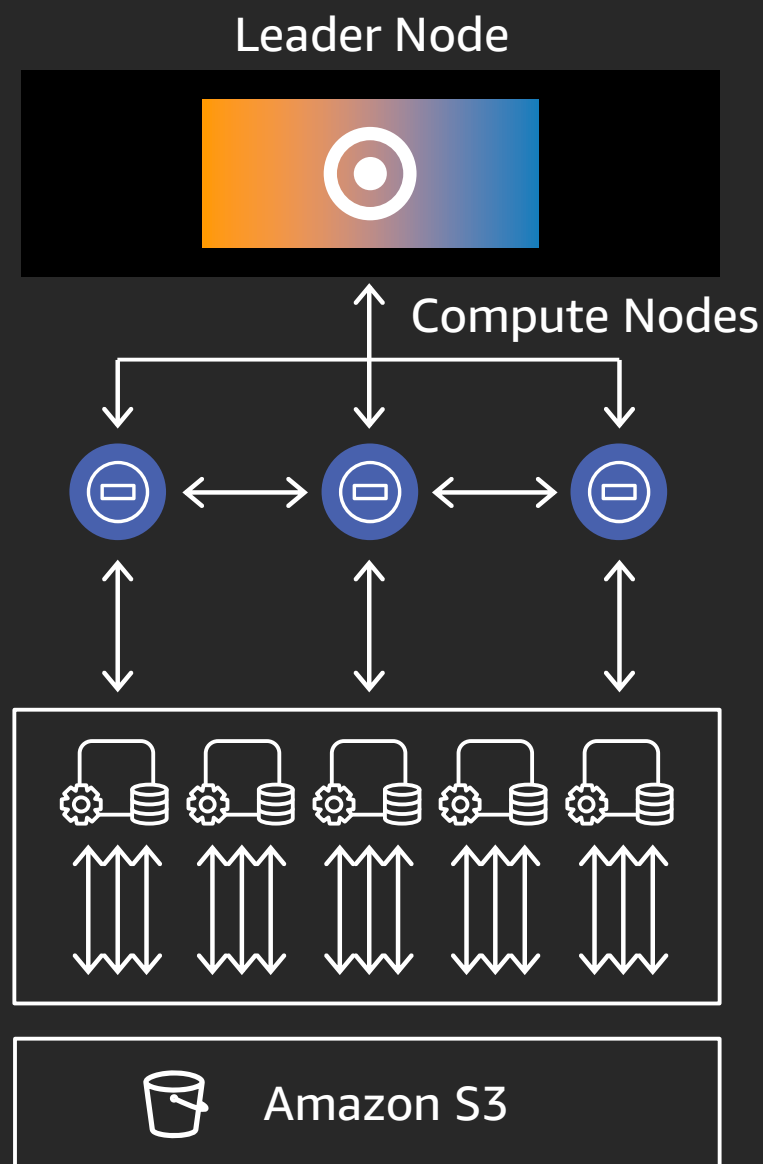Scale compute and storage separately

No server to maintain for S3 query

Support for Parquet, ORC, Avro, CSV, JSON, Grok, and other open file formats

Pay only for the amount of data scanned

# Challenges with rapidly growing data

## Amazon Redshift Architecture

**Leader Node**



**Compute Nodes**



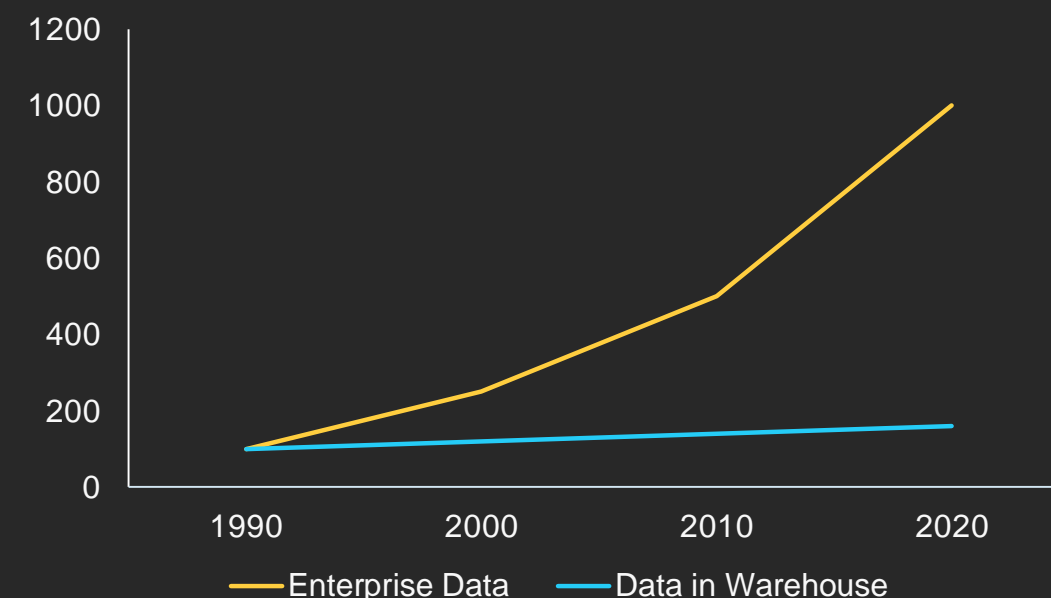Amazon S3

## Either

### Compute optimized

DC2.8xlarge

2.56 TB SSDs storage

DC2.large

.16 TB SSDs storage

## Or

### Storage optimized

DS2.8xlarge

16 TB HDDs storage

DS2.XL

2 TB HDDs storage

## Growing dark data



— Enterprise Data    — Data in Warehouse

## Solution until NOW

Add nodes
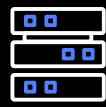
Delete old data

Unload data to data lake

# 3rd generation compute instance: RA3
## Scale compute and storage independently

**New!**

**Managed storage**
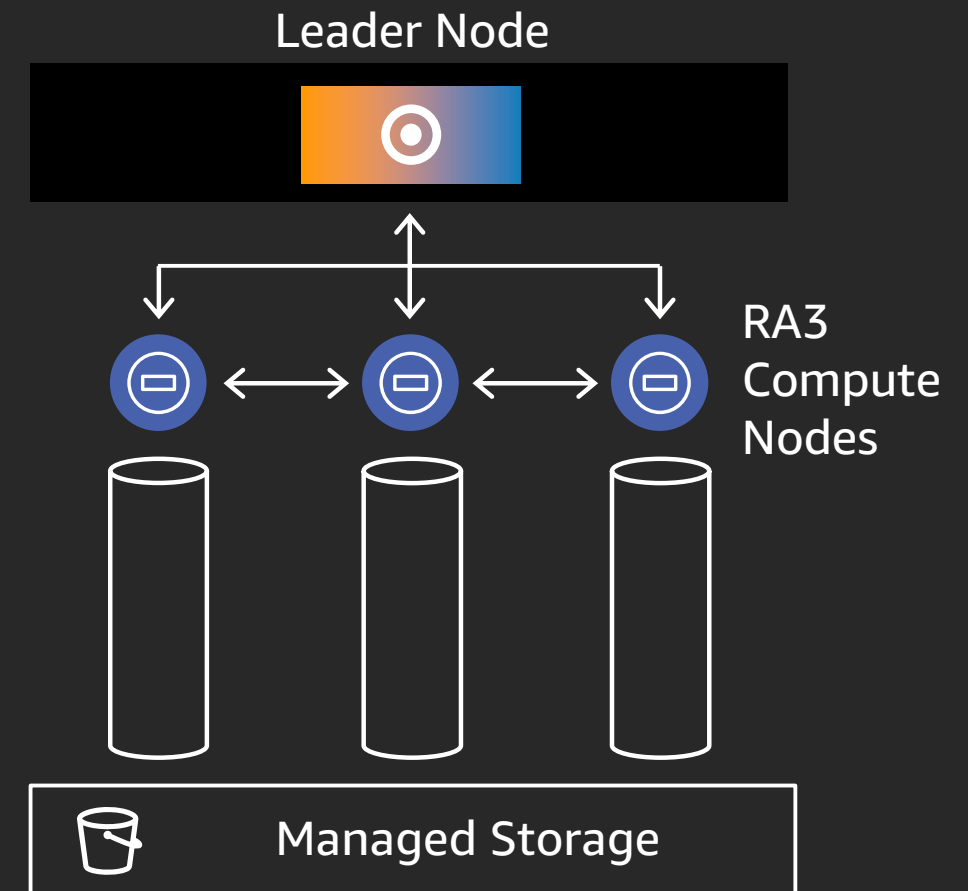
**Large high-speed cache**

**High-bandwidth networking**

**Size data warehouse** only based on steady-state compute needs

**Scale and pay independently** for compute and storage

**Automatic,** no changes to any workflows, no need to manage storage

Leader Node

RA3 Compute Nodes

Managed Storage

# RA3: Unmatched performance at unbeatable price

## RA3.16xl

Can scale to tens of PB of data (8 PB compressed)

On demand price — $13.04/node/hr

For storage pay $0.024/GB/month

**Coming soon
RA3.4xl**

2x performance and 2x storage capacity compared to DS2.8XL at the same on-demand price

3x price-performance compared to any other Cloud DW

Up-to 64 TB in managed storage per RA3.16xl node
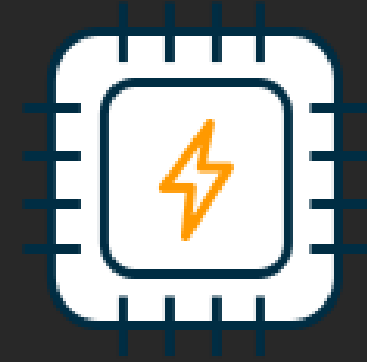
# RA3: Node specification

Node size: <u>ra3.16xlarge</u>

Node counts: 2-128

vCPUs: 48

Memory (GiB): 384

Managed storage quota: 64 TB (compressed)

Largest cluster: 8 PB (compressed)



AWS Nitro System

Breaks apart hypervisor, storage, networking, and management
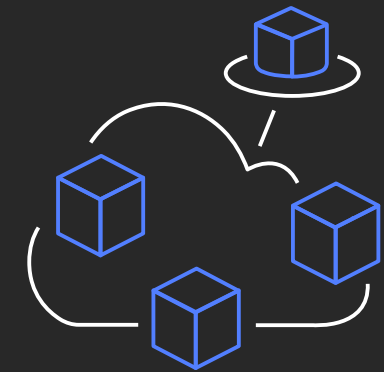
Offloads to dedicated hardware and software

# RA3: Migration from DS2

**Most DS2.8XL** clusters will get **up to 2x performance** and **2x storage** with RA3.16XL for the same on-demand price **(in 2:1 ratio)**
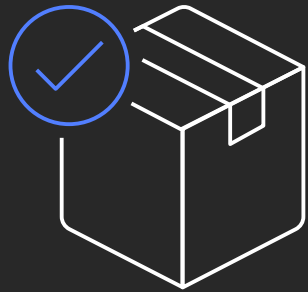
Can migrate in 2:1, 3:1, or even 4:1 node count ratio (DS2.8XL:RA3.16XL)

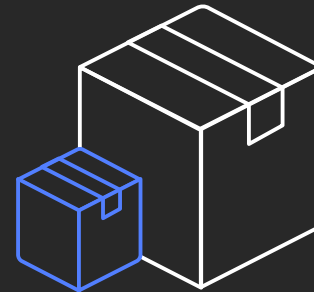Smaller DS2 clusters with under 10 TB, best suited for RA3.4XL
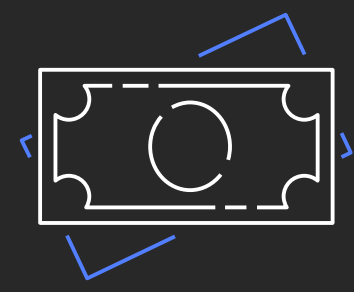
# RA3: Migration from DC2

**Larger DC2.8XL** clusters who **need more storage capacity** can potentially benefit from RA3

Can migrate in **3:1 node count ratio** (DC2.8XL to RA3:16XL) for price-equivalent

At price-equivalent RA3 provides similar performance to DC2 but provide **8x more storage capacity**

For smaller clusters with **5-10 TB of data, stay with DC2** for best price-performance

# RA3 evaluation results: 3 examples

➢ **Customer 1**

   ➢ Compared price-equivalent *14 nodes DS2.8XL to 7 nodes of RA3.16XL*; *most queries* were up-to *2.1x faster*.

➢ **Customer 2**

   ➢ Compared price-equivalent *15 nodes of DC2.8XL to 5 nodes of RA3.16XL*; *most queries were 1.25x faster*, some queries were *.8x slower*.

➢ **Customer 3**

   ➢ Compared price-equivalent *16 node DS2.8XL to 8 nodes of RA3.16XL*; most queries and ETL were *1.3x faster*.

# RA3: Migration considerations

## Migrate using restore from snapshot

- Get a new RA3 cluster in minutes
- Validate the new RA3 cluster and delete the old cluster
- Use modify cluster to rename the RA3 cluster to old cluster's name
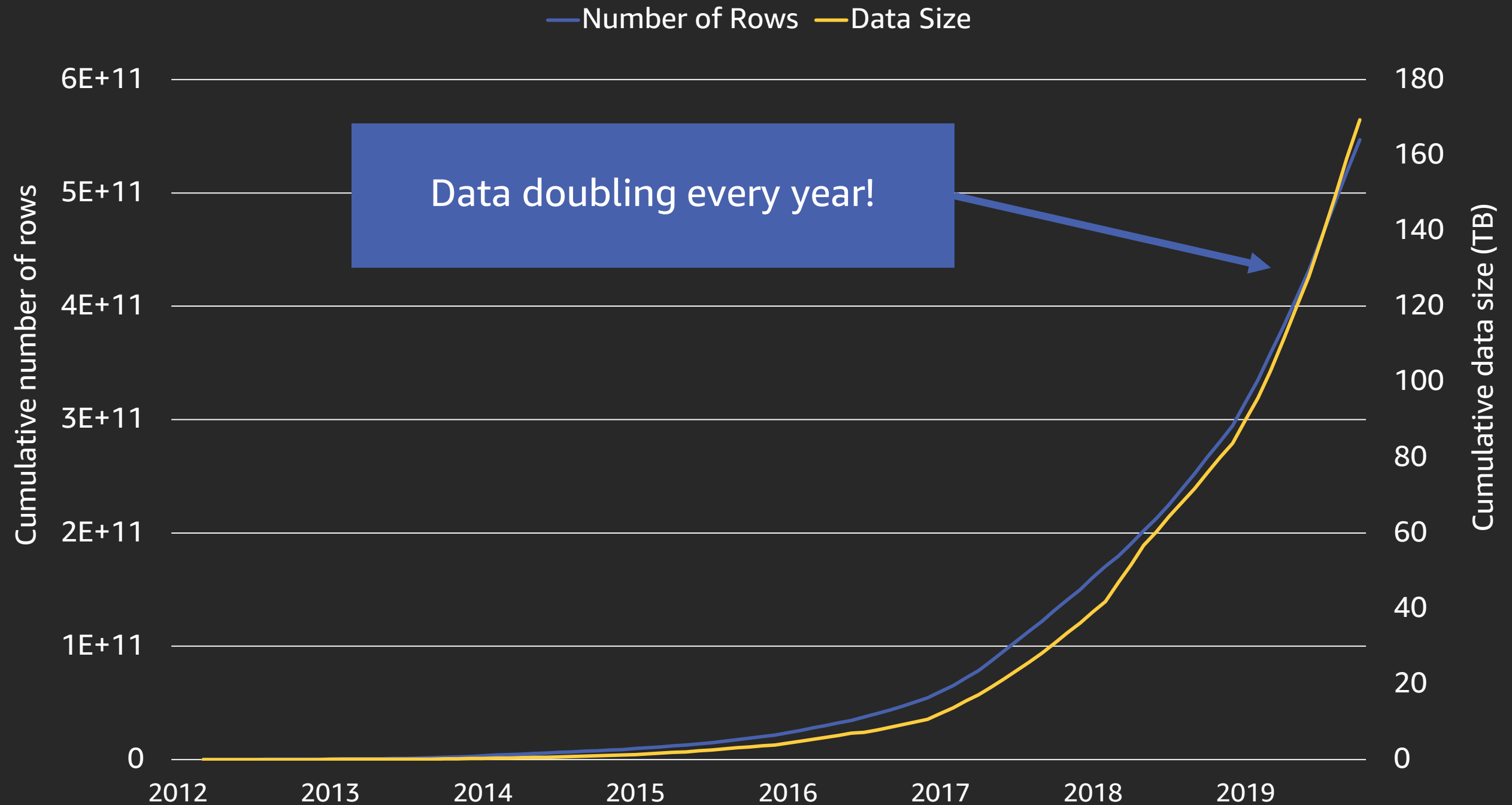- Reduces the flexibility of Elastic Resize

## Another option is classic resize

- Classic resize copies data from old to new cluster and renames the cluster upon completion (Classic resize is slower than restore)
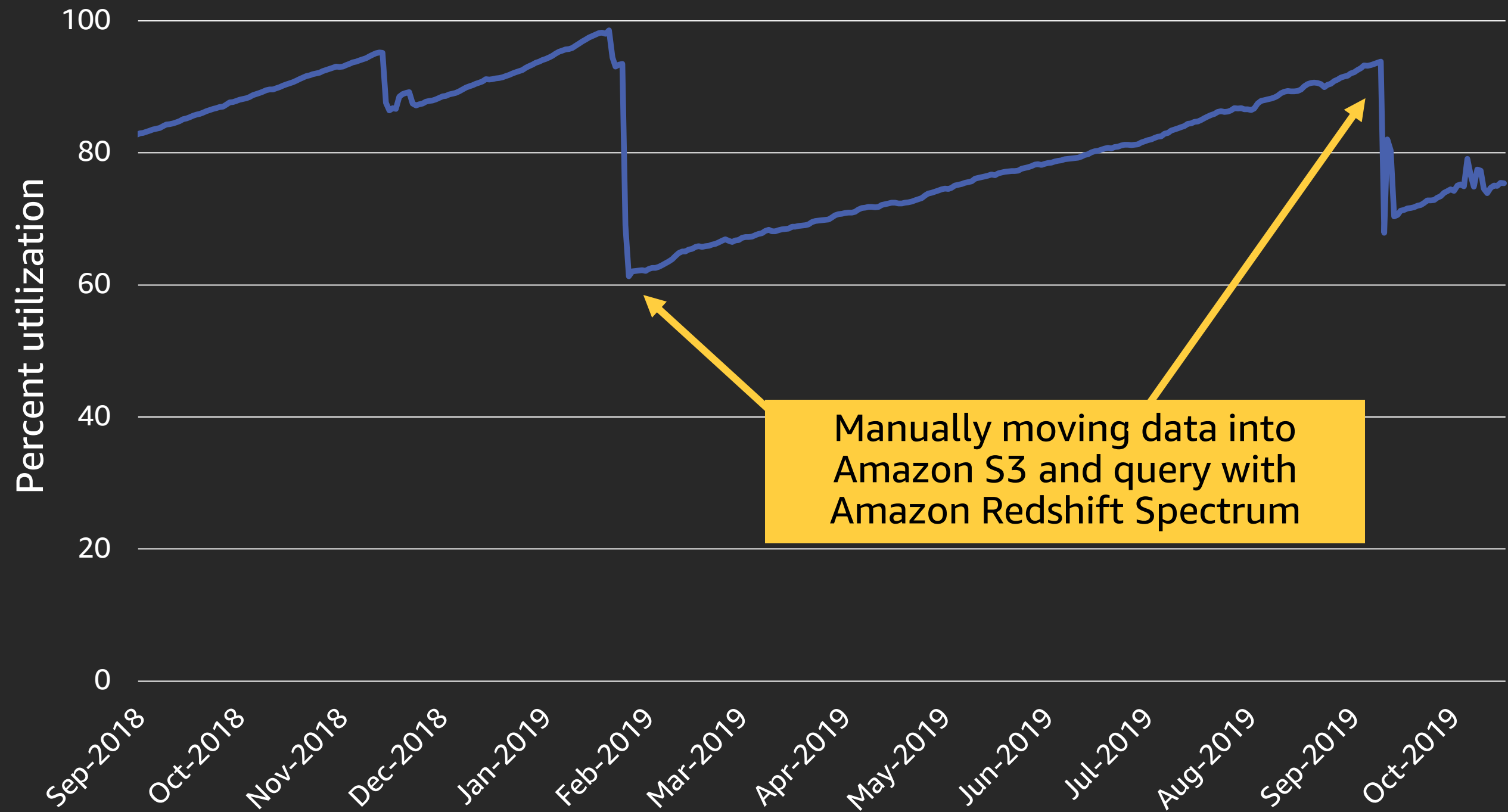- Retains full flexibility of Elastic Resize

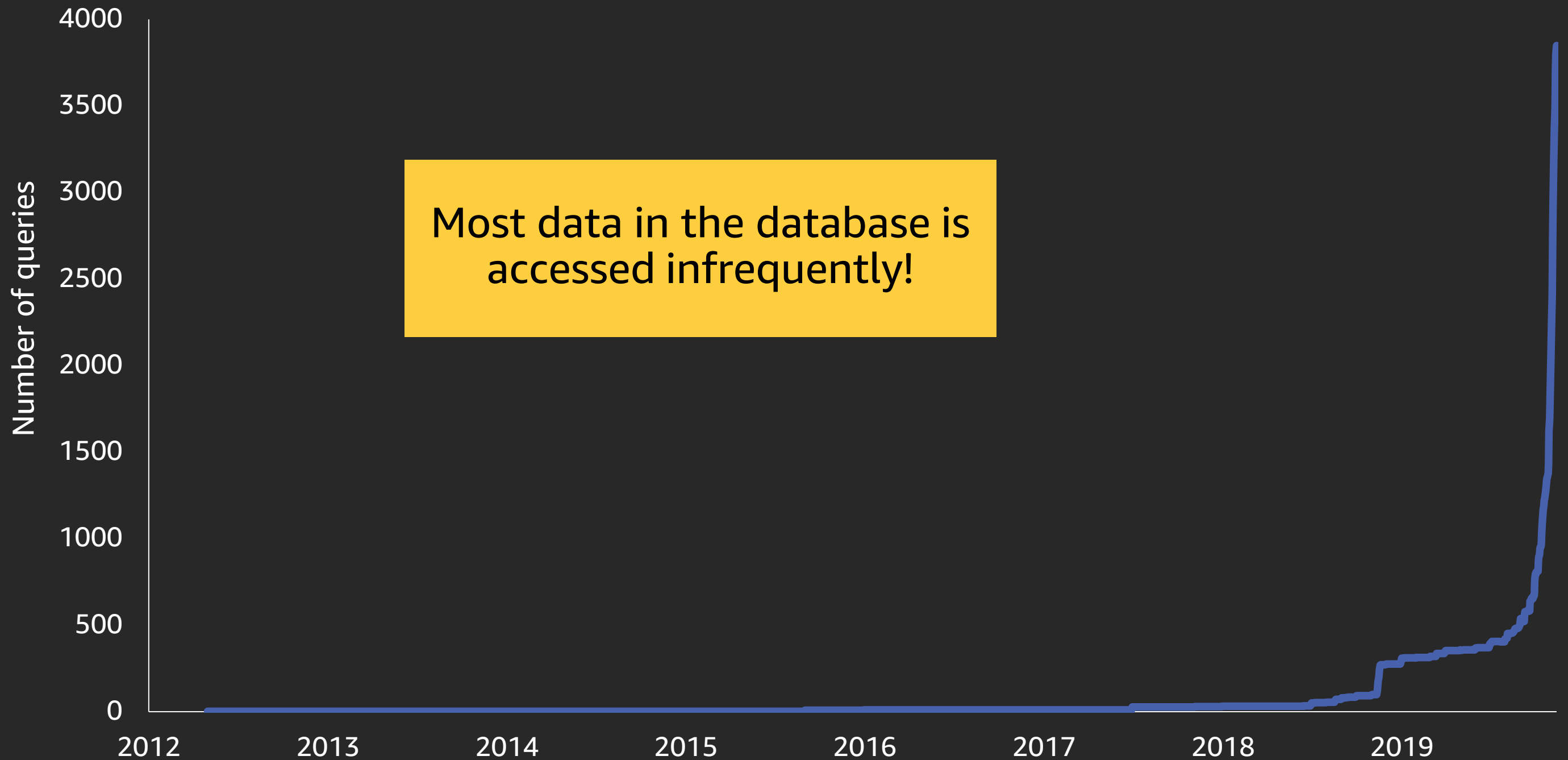Analytics data stored in Amazon Redshift

# Redshift disk utilization



Manually moving data into
Amazon S3 and query with
Amazon Redshift Spectrum

# Experience with RA3
## Queries above yellow line were slower; below were faster

- **2x Faster** COPY performance

- 78% of **ad-hoc queries** performed faster (**median improvement: 2.1x**)

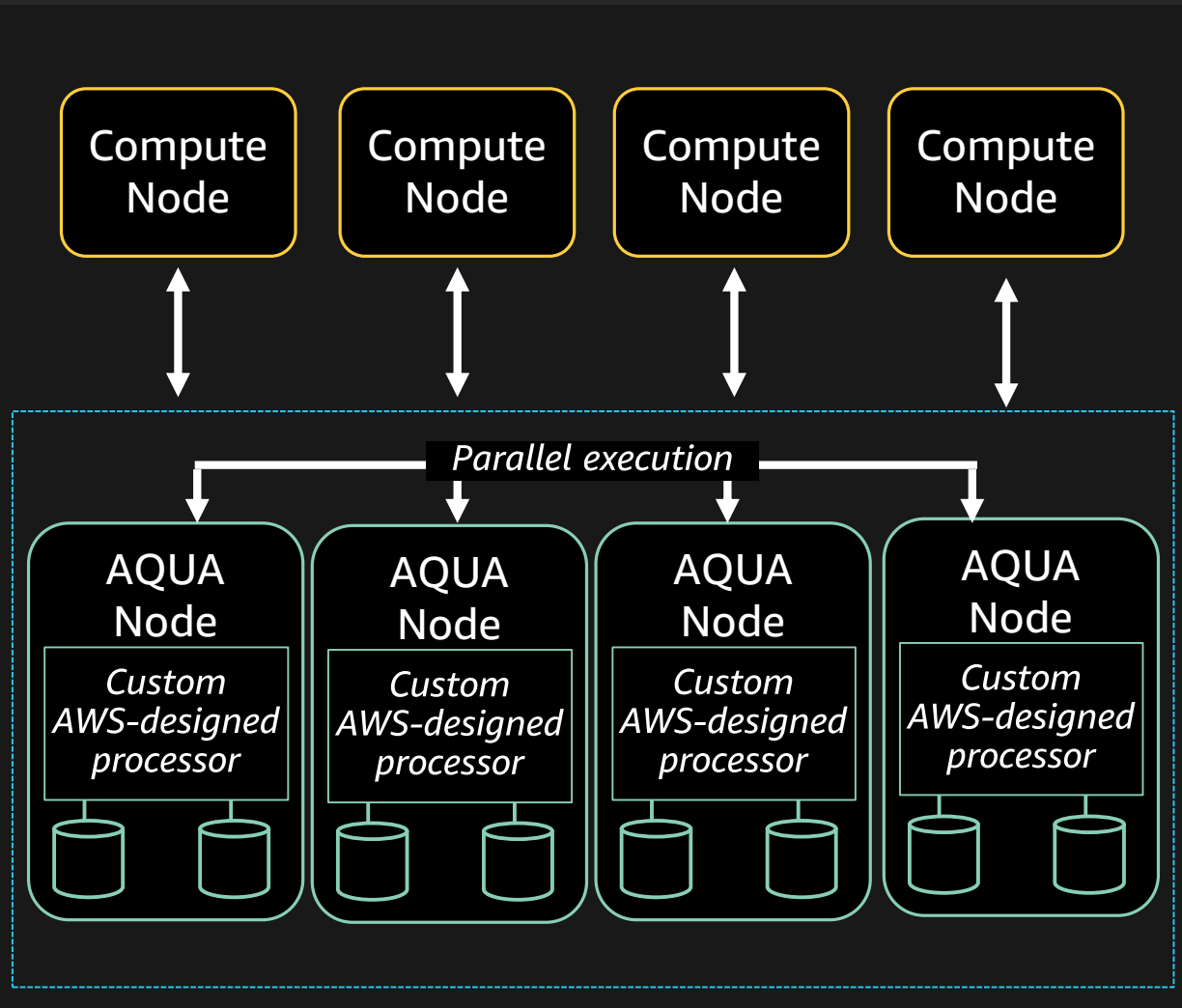- **2.3x average runtime improvement** for our query benchmark



Query runtime (seconds)

RA3 with managed storage vs. Current Amazon Redshift cluster

# Performance and concurrency at ever increasing scale

# AQUA for Amazon Redshift - Advanced Query Accelerator

A new distributed and hardware-accelerated processing layer that will make Amazon Redshift **10x faster** than any other cloud data warehouse without increasing cost

**Preview!**

| Compute Node | Compute Node | Compute Node | Compute Node |

*Parallel execution*

| AQUA Node | AQUA Node | AQUA Node | AQUA Node |
| *Custom AWS-designed processor* | *Custom AWS-designed processor* | *Custom AWS-designed processor* | *Custom AWS-designed processor* |

Minimize data movement over the network by pushing down operations to AQUA Nodes

AQUA Nodes with custom AWS-designed analytics processors to make operations (compression, encryption, filtering, and aggregations) faster than traditional CPUs

Available only with RA3, no code changes required. Available in preview.

# Two forms of compute elasticity

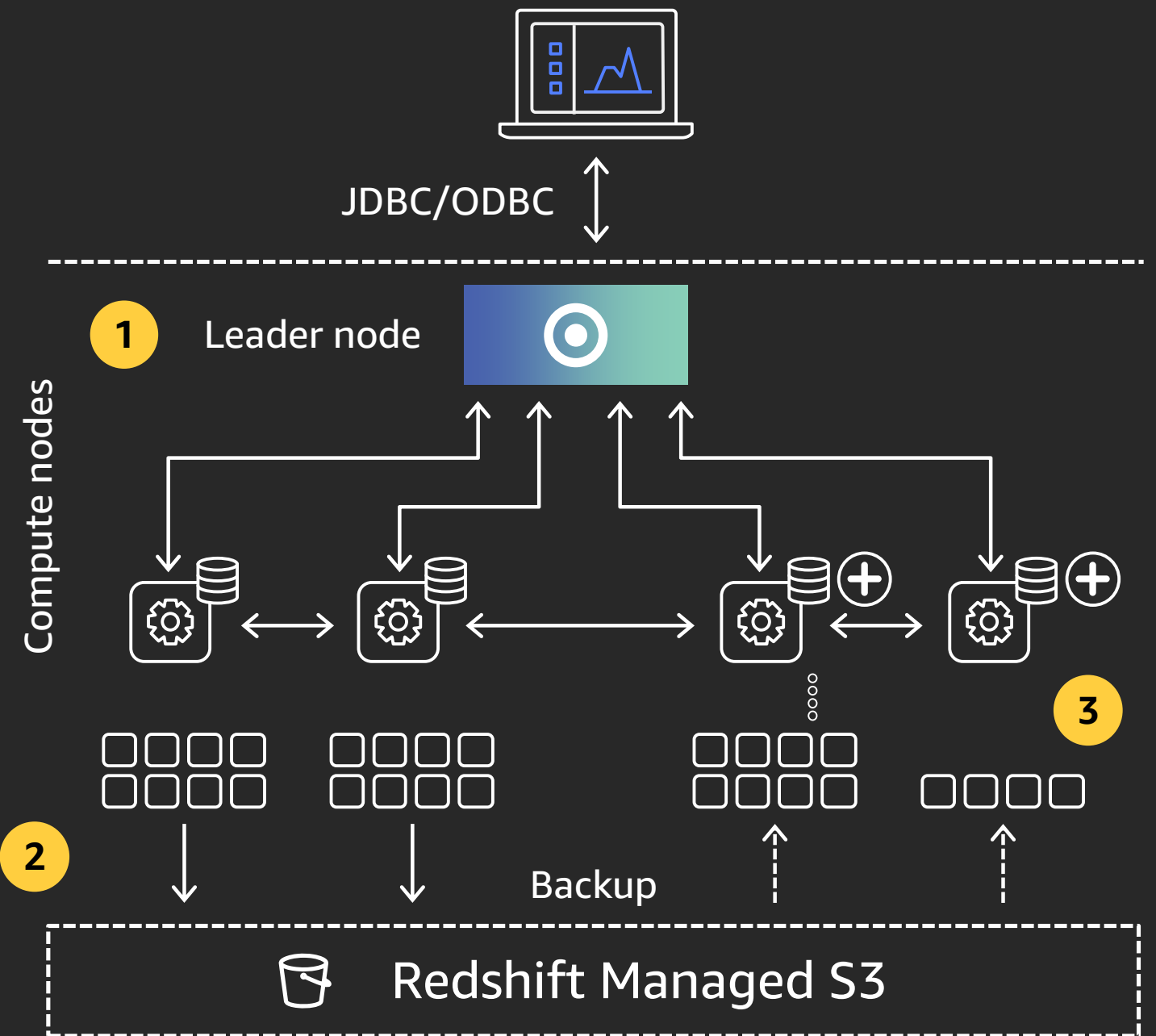|  | Vertical scaling | Horizontal scaling |
|---|---|---|
| Question | How can I speed up my running jobs? | How do I support spikes in users without provisioning for peak demand? |
| Answer | Add more nodes with Elastic Resize | Enable concurrency scaling |

# Elastic resize: Change cluster performance
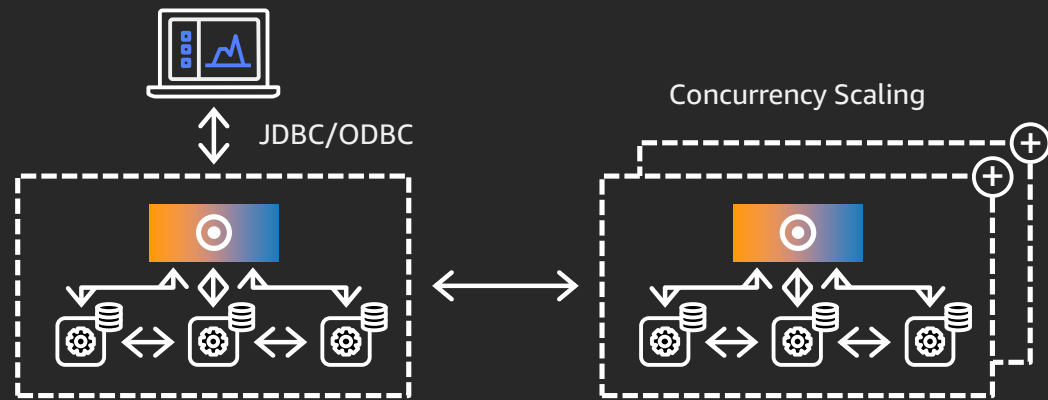
Add or remove compute
nodes to an existing cluster

Completes within few minutes

Minimal disruption to sessions
and queries running

JDBC/ODBC

**1** Leader node

Compute nodes

**2**

**3**

Backup

Redshift Managed S3

# Concurrency scaling: Eliminate wait time for bursts of users

JDBC/ODBC

Concurrency Scaling

Scale-out to multiple Amazon Redshift clusters from a single endpoint in seconds
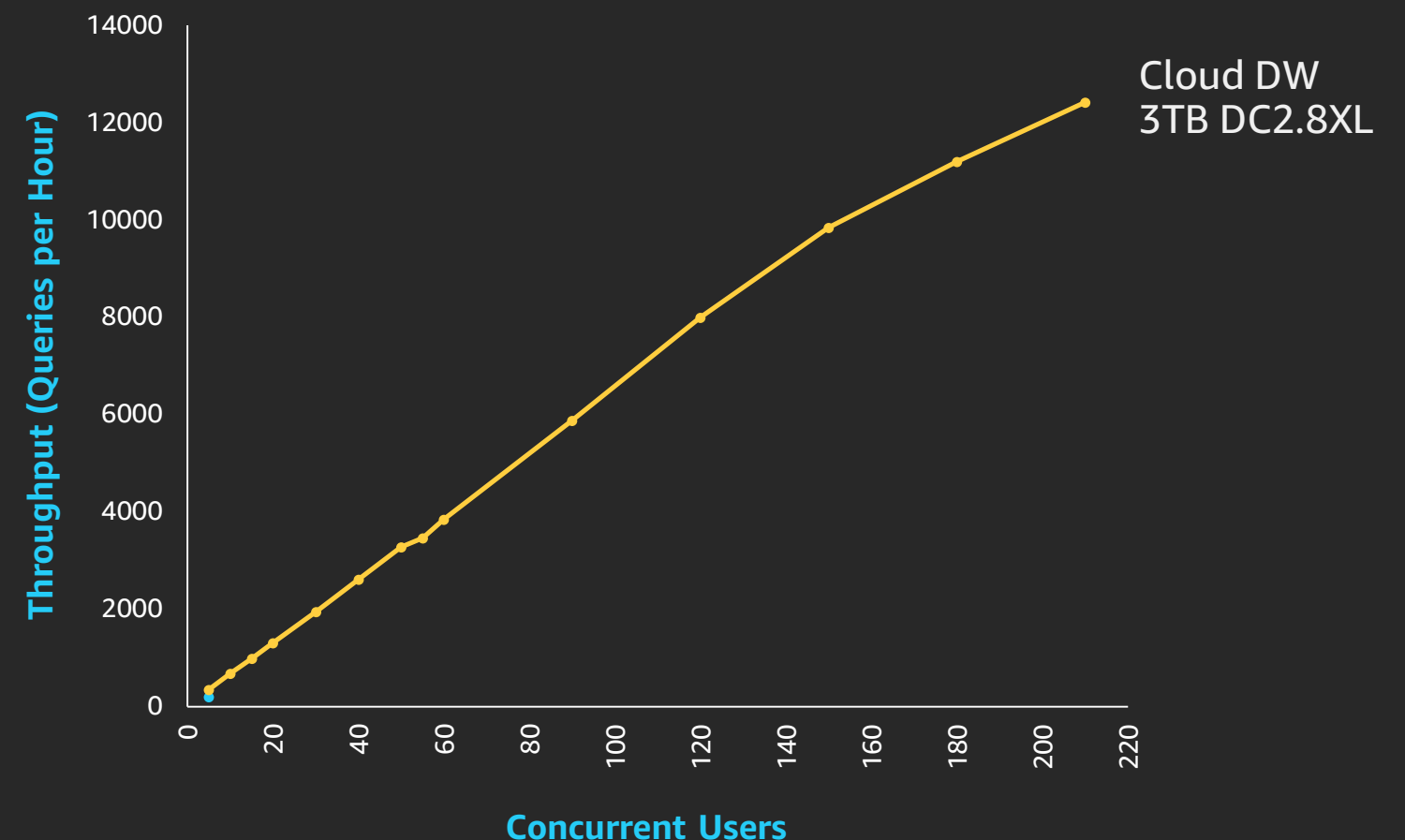
Support virtually unlimited concurrent users and queries while maintaining SLAs

Per-second billing for additional clusters used

Free 1hr usage per day
(free for 97% of clusters)

**35x improvement
in throughput in 2019**

**Scalability improvements**

Cloud DW
3TB DC2.8XL

Throughput (Queries per Hour)

14000

12000

10000

8000

6000

4000

2000

0

0    20    40    60    80    100    120    140    160    180    200    220

**Concurrent Users**

# Multiple analytics needs

# Enable all your analytical workloads: Choose best tool for the job



QuickSight
EMR
Athena
Kinesis
Elasticsearch Service
Redshift
AI Services

S3

Snowball
Snowmobile
Kinesis Video Streams
Kinesis Data Firehose
Kinesis Data Streams

## Exabyte scale

## Store and analyze relational and non-relational data

## Purpose-built analytics tools

## Cost effective
Store at 2.3 cents per GB/month in Amazon S3
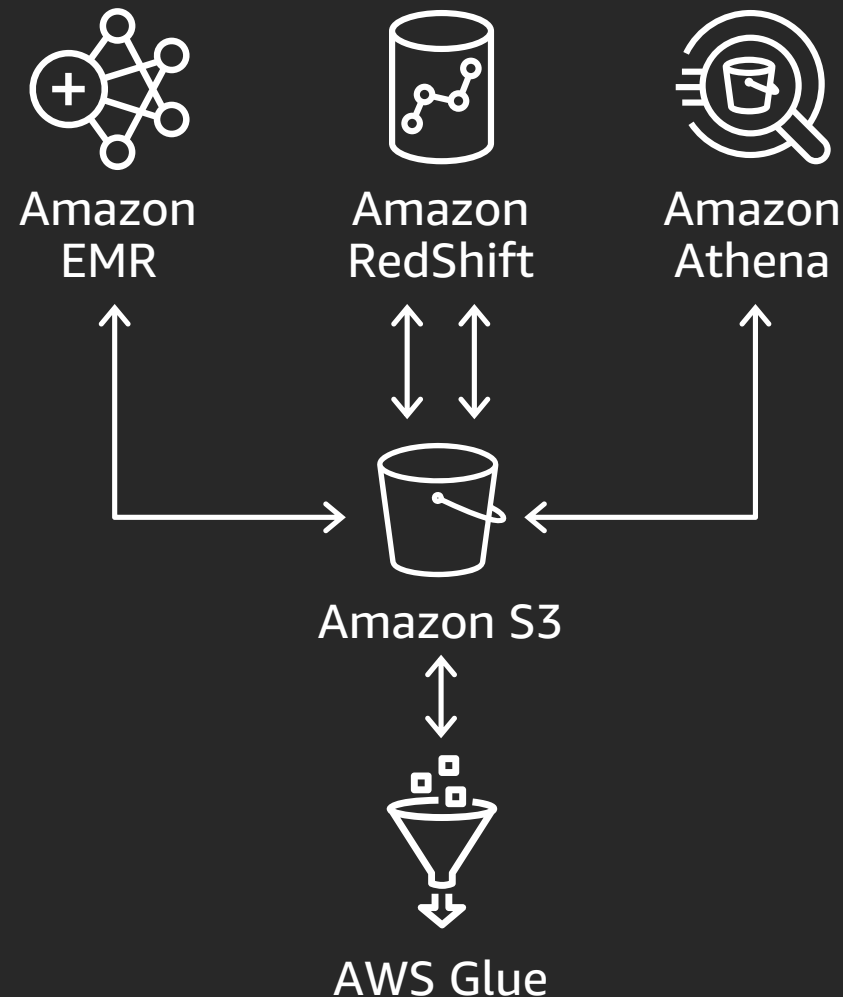Query with Amazon Athena at ½ cent per GB scanned
DW with Amazon Redshift for $1,000/TB/year

## Give access to everyone
Amazon QuickSight: $.0.30/session up to max of $5 per month. No usage, no fee; little usage, little fee; max $5 per month

# Export Amazon Redshift data as Parquet to S3

Amazon Redshift now supports exporting data to S3 in Parquet format. This makes **sharing data across the data lake easier and faster, without conversion.**

Amazon EMR

Amazon RedShift

Amazon Athena

Amazon S3

AWS Glue

**Parquet is an open data format** supported by EMR, Athena, and Redshift

Amazon Redshift Unload command now supports Parquet format. This allows data in Redshift to be exported as Parquet to be processed by EMR or Athena without any data conversion.

# Security

aws

# Amazon Redshift: Security is built in at no extra cost
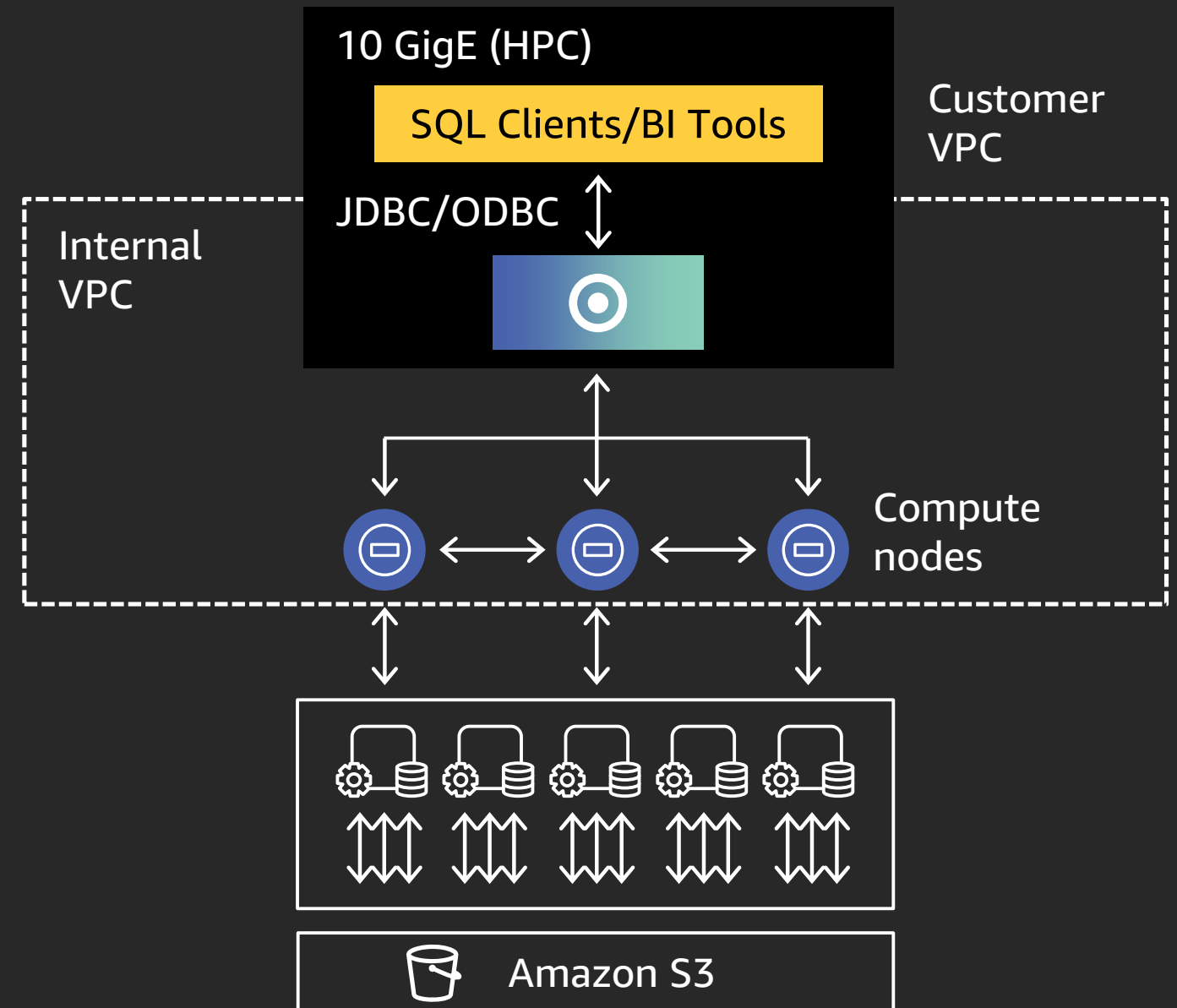
**AWS IAM integration**

**End-to-end encryption**

**Integration with AWS Key Management Service**
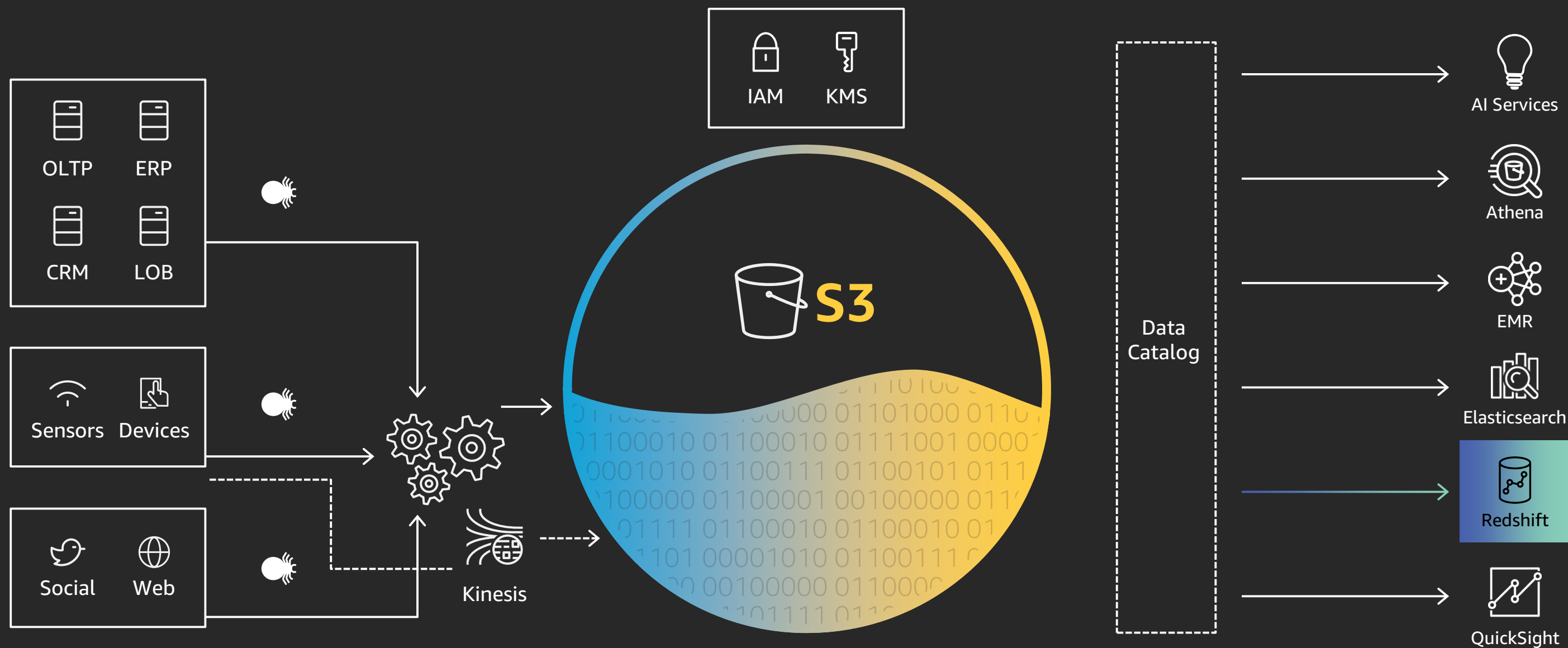
**Select compliance certifications***

AICPA SOC (Formerly SAS 70 Reports)

DEPARTMENT OF DEFENSE · UNITED STATES OF AMERICA

ISO

PCI DSS COMPLIANT

HIPAA COMPLIANT

FR FedRAMP

*Full list of compliance certifications is available here: https://aws.amazon.com/compliance/

## Network isolation

10 GigE (HPC)

SQL Clients/BI Tools

JDBC/ODBC

Customer VPC

Internal VPC

Compute nodes

Amazon S3

# Unified column level access control for the data lake

**New!**

IAM  KMS

OLTP  ERP

CRM  LOB

Sensors  Devices

Social  Web

Kinesis

**S3**

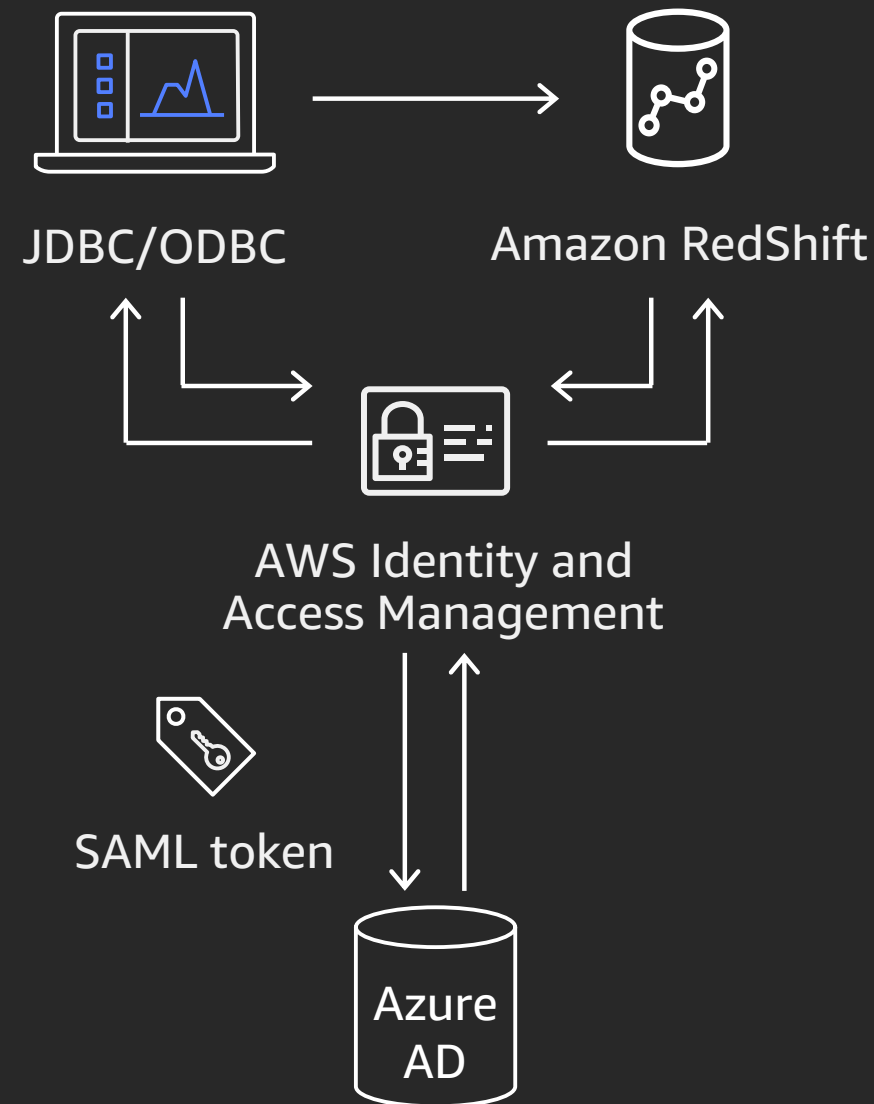Data Catalog

AI Services

Athena

EMR

Elasticsearch

Redshift

QuickSight

# Single-Sign On with Azure Active Directory

Amazon Redshift now
integrates with
Azure Active Directory
to provide
Single-Sign On

JDBC/ODBC

Amazon RedShift

AWS Identity and
Access Management

SAML token

Azure
AD

## Single-Sign On with Azure Active Directory

SAML compliant Single-sign On.
Redshift ODBC/JDBC drivers support
industry standard SAML workflows
and integrate with both on-premise
and Cloud SSO providers. Azure AD,
Active Directory Federation Services,
Okta, Ping Federate.

### Benefits

Simple: Re-use corporate
identity with Redshift

Compliance: use Azure AD base
password policies, password
rotation, onboarding etc

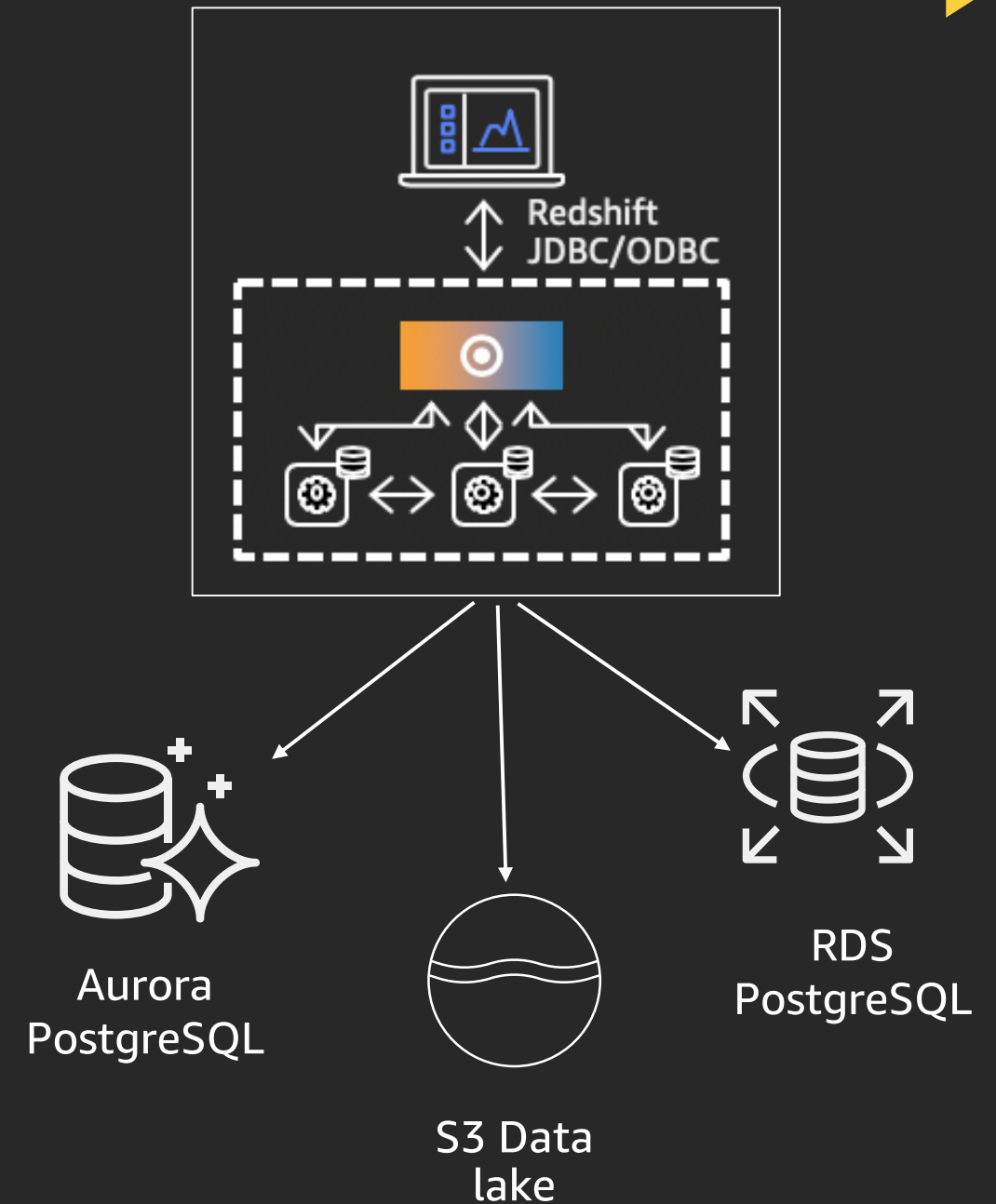# Reduce TCO: easier Amazon Redshift

# Amazon Redshift Federated Query

Query and join data from one or more RDS and Aurora PostgreSQL databases

Analytics on operational data without data movement and ETL delays

Integrate operational data with data warehouse and S3 data lake

Flexible and easy way to ingest data avoiding complex ETL pipelines

Intelligent distribution of computation to remote sources to optimize performance

Redshift JDBC/ODBC

Aurora PostgreSQL

S3 Data lake

RDS PostgreSQL

# Materialized views
## Compute once, query many times

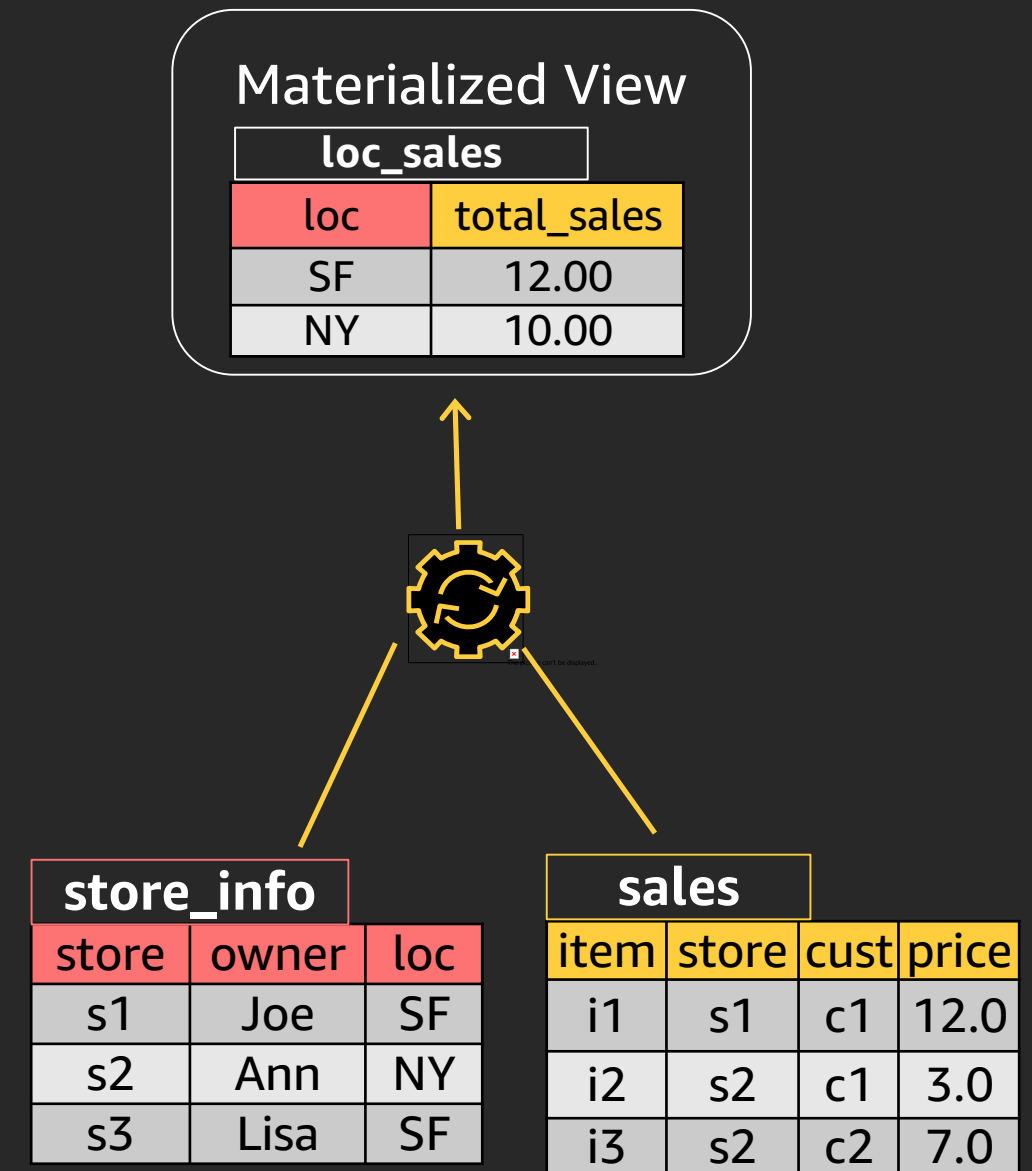Speed up queries by orders of magnitude
- Joins, filters, aggregations, and projections

Simplify and accelerate ETL/BI pipelines
- Incremental refresh
- User triggered maintenance

Easier and faster migration to Amazon Redshift

**Materialized View**

| loc_sales | |
|---|---|
| **loc** | **total_sales** |
| SF | 12.00 |
| NY | 10.00 |

**store_info**

| store | owner | loc |
|---|---|---|
| s1 | Joe | SF |
| s2 | Ann | NY |
| s3 | Lisa | SF |

**sales**

| item | store | cust | price |
|---|---|---|---|
| i1 | s1 | c1 | 12.0 |
| i2 | s2 | c1 | 3.0 |
| i3 | s2 | c2 | 7.0 |

# Amazon Redshift automates tuning and maintenance

**Simplified** user experience

---

**Optimizes** for **peak performance** as workloads and data scale

---

**Automatic data layout changes** and smart **recommendations based on continuous analysis** of workloads

Automatic
Analyze

Automatic Table
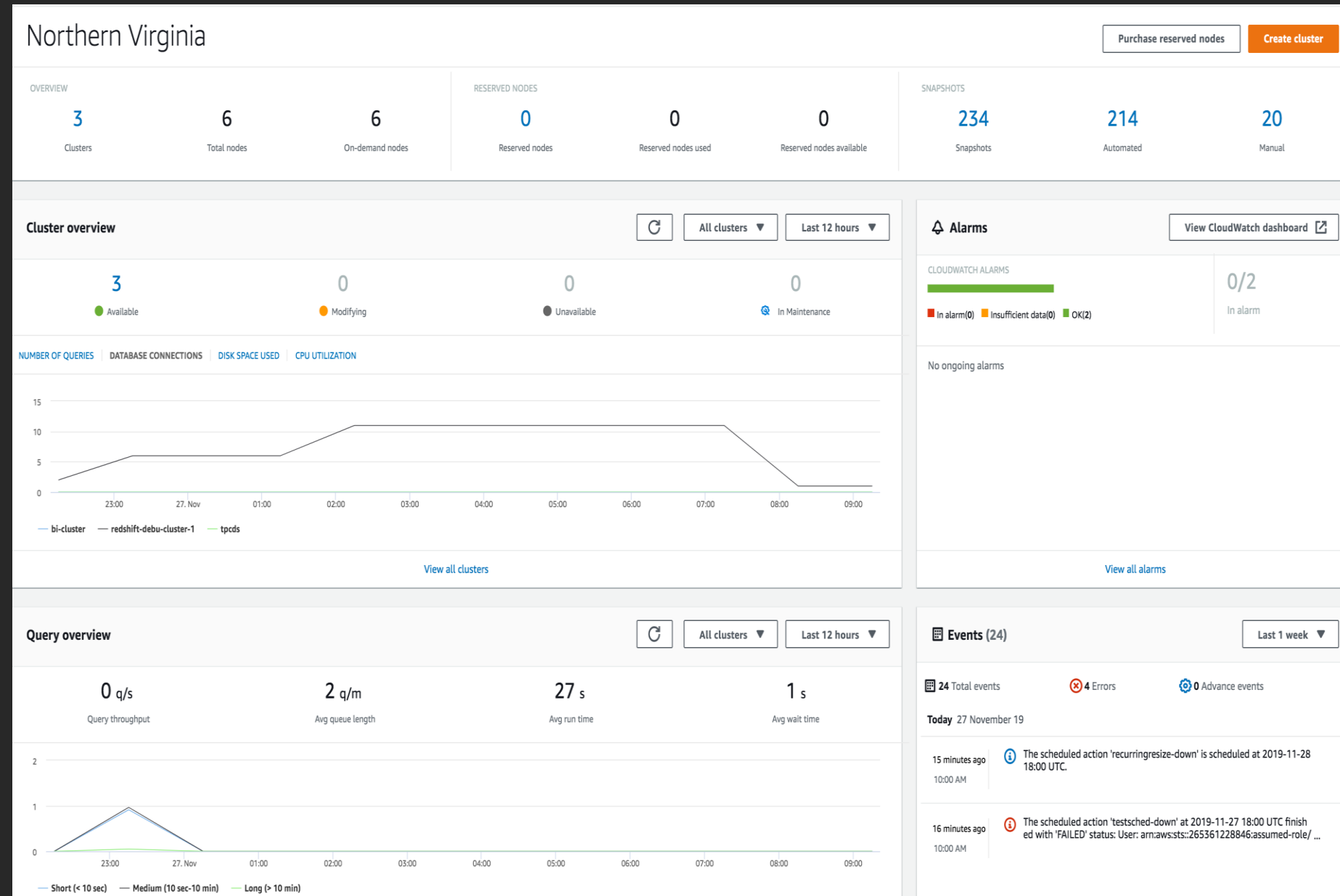Distribution Style

Distribution/Sort key
advisors

Automatic
Vacuum Delete

Automatic
Table Sort

# New Amazon Redshift console
## Modernizes interface and enhances user experience



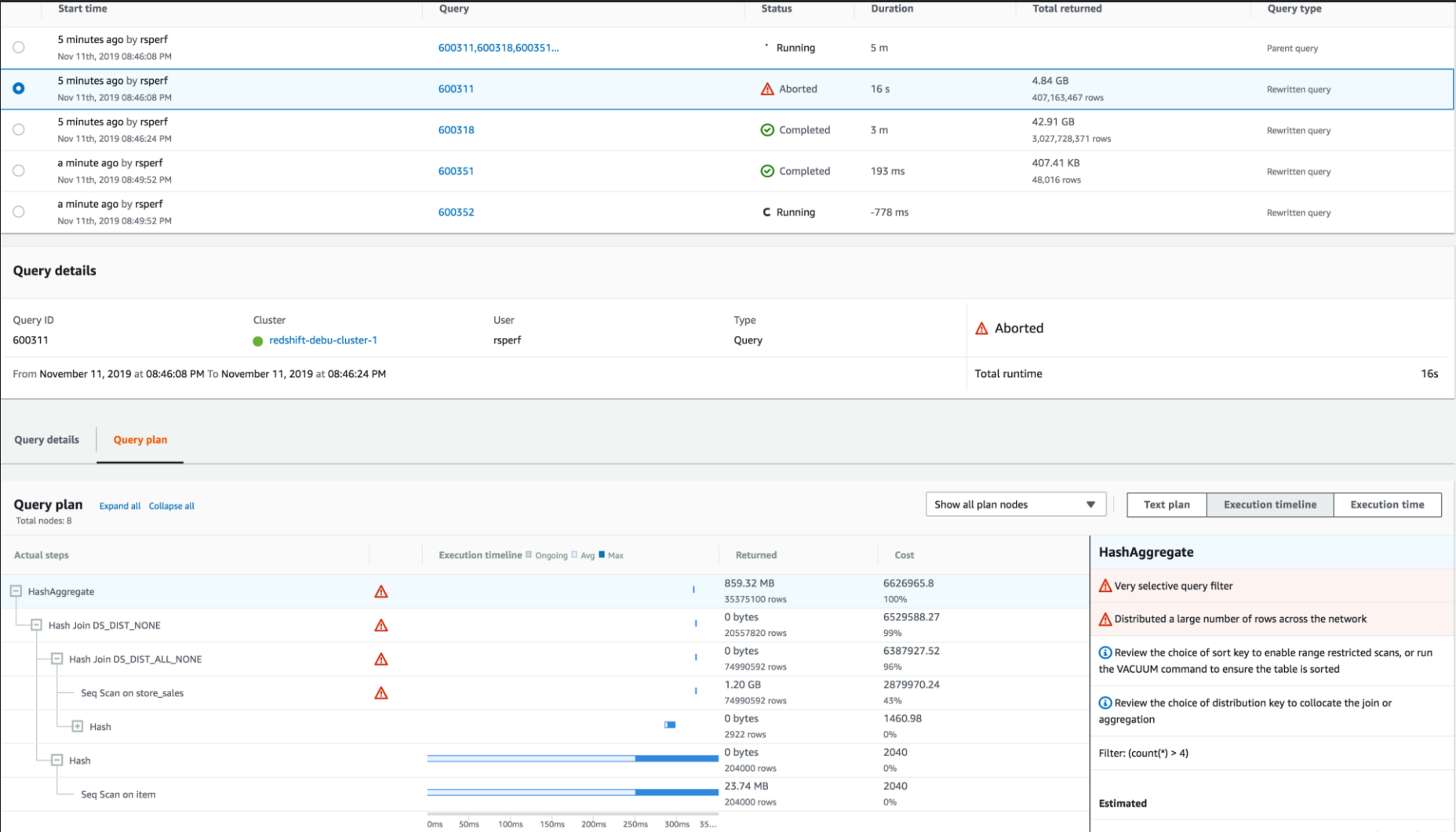Gain visibility to health of all clusters in your account

Simplify creation and management of clusters

Reduced time to diagnose user query performance issues

Share Query Editor with non-admin users
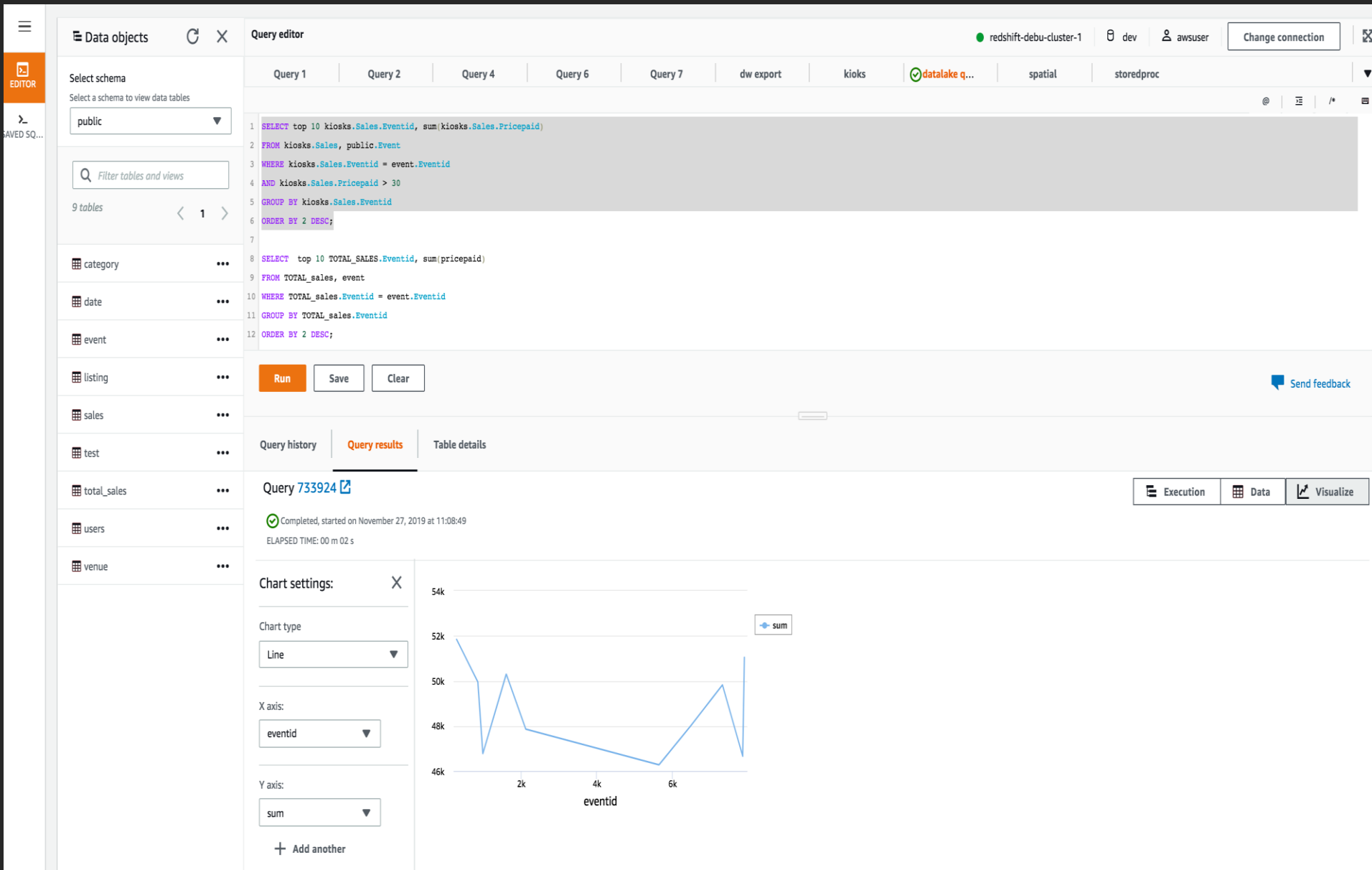
# Monitoring of User Queries
Diagnose query performance faster



Monitor your queries

View all rewritten query in context

Visual analysis of query plan

Correlate with cluster performance

View in-place recommendation

# Query Editor
Easier to run and analyze your queries



Share the query editor to non-admin users as a separate URL

Command assist, auto-complete and keyboard short-cuts

Visually analyze your query results

In-place analysis of query plan

# Stored procedures support to simplify migrations

Use Schema Conversion Tool to automatically convert your stored procedures

Migrating to Amazon Redshift is even easier!

Amazon Redshift supports Stored Procedures in PL/pgSQL format

**Stored procedures used for ETL, data validation, and custom business logic close to data.**

```
CREATE OR REPLACE PROCEDURE test_sp1(f1 int, f2 varchar)
AS $$

BEGIN

        RAISE INFO 'f1 = %, f2 = %', f1, f2;

END;

$$ LANGUAGE plpgsql;


call test_sp1(5, 'abc');


INFO: f1 = 5, f2 = abc
```

# Spatial processing

Spatial Analytics at scale — ingest,
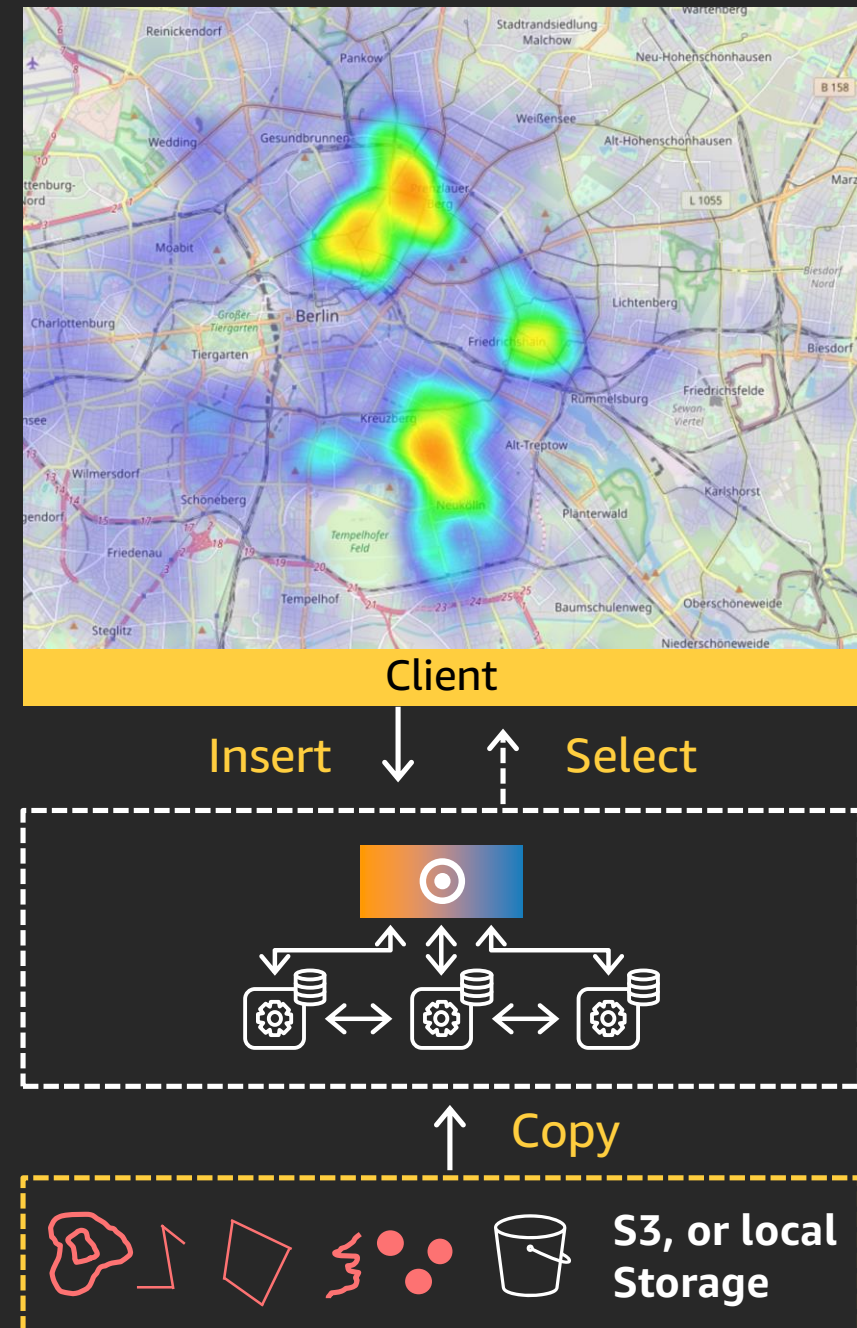store and analyze spatial data

Seamlessly integrate spatial and business data

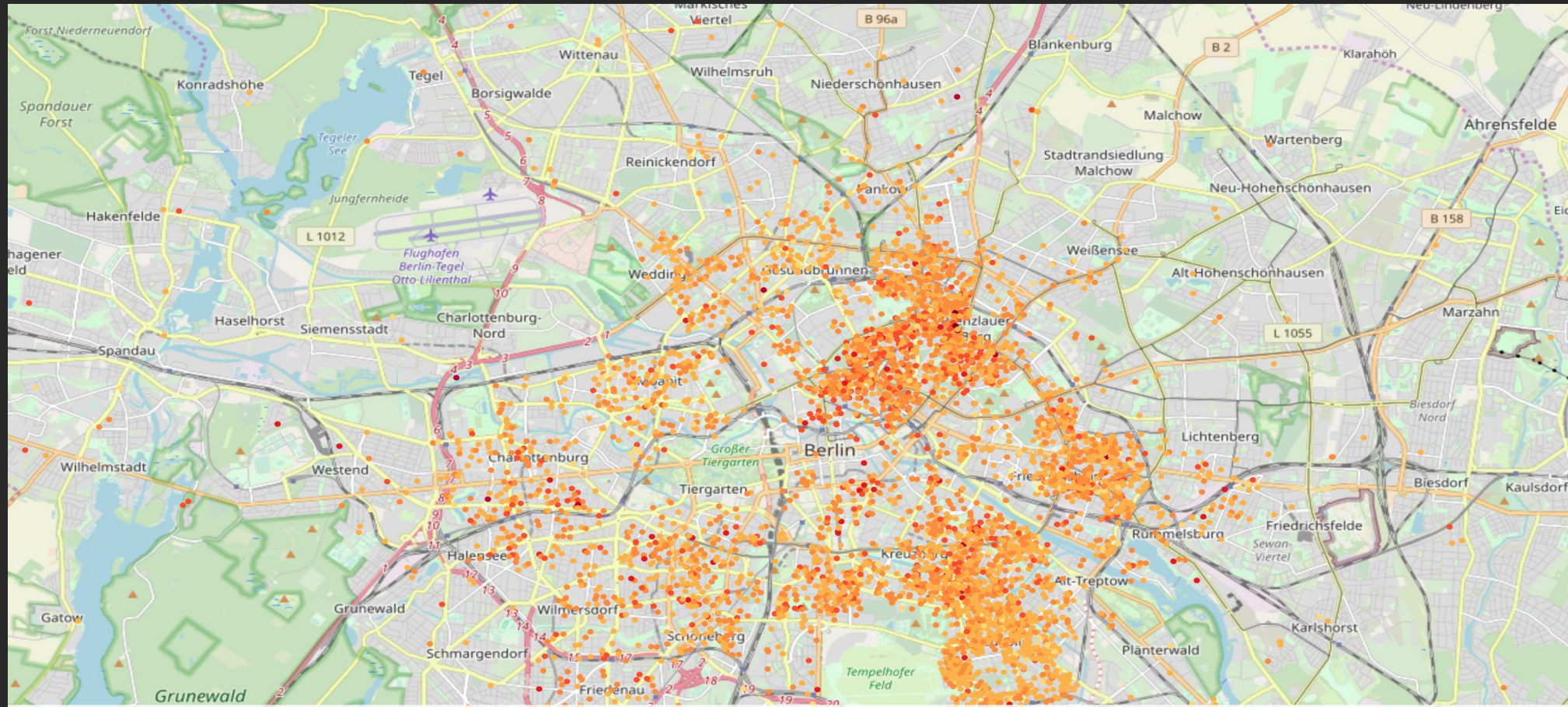Get new dimension of insights and value

New data type GEOMETRY

40+ SQL spatial functions
Accessors, Constructors, Predicates

# Spatial processing — sample query



**Data Types**
GEOMETRY
Point, Linestring, Polygon,
MultiPoint, MultiLinestring,
MultiPolygon,
GeometryCollection

**Spatial Accessors**
ST_NumGeometries,
ST_GeometryType,
ST_Dimension, …

**Spatial Predicates**
ST_Covers, ST_Equals,
ST_Within, ST_DWithin,…

**Spatial Functions**
ST_Distance,
ST_Azimuth, …

**Spatial Formats**
WKT/WKB, EWKT/EWKB, GeoJSON
Ingestion: CSV

```
SELECT name, ST_X(shape) as lng, ST_Y(shape) as lat, price
FROM accommodations
WHERE ST_Within(shape, ST_GeomFromText( 'POLYGON((13.11839294433596
52.428594259606 3, 13.11839294433596 52.60117089057946, 52.428594259606 3))', 4326))
LIMIT 5000
```
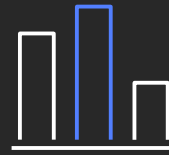
# Amazon Redshift benefits
Tens of thousands of customers use Amazon Redshift and process exabytes of data per day

### Data lake & AWS integrated
Lake Formation catalog and security,
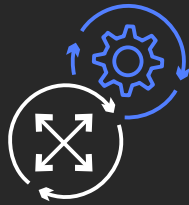Exabyte scale query (spectrum & federated),
AWS integrated (DMS, CloudWatch)

### Best performance
3x faster than other
cloud data warehouses

### Lowest cost
75% less expensive than all other cloud
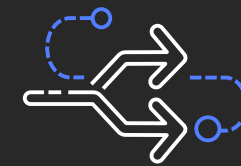data warehouses and predictable costs

### Most scalable
Virtually unlimited
concurrency, scale compute and storage
independently

### Most secure & compliant
AWS-grade security, (e.g. VPC, encryption
with KMS, Cloud Trail), Certifications such
as SOC, PCI, DSS, ISO, FedRAMP, HIPAA

### Fully managed
Easy to provision and manage, automated
backups, AWS support, 99.9% SLAs

# Thank you!

aws

Please complete the session survey in the mobile app.