

The logo for AWS re:Invent features the words "AWS" and "re:Invent" stacked vertically. "AWS" is in a smaller, sans-serif font above "re:Invent", which is in a larger, bold, sans-serif font.

AWS
re:Invent

ANT350 - R

What's New with Amazon Redshift ft. Dow Jones

Colleen Camuccio
VP, Program Management
Dow Jones

Vidhya Srinivasan
General Manager, Amazon Redshift
AWS



DOW JONES

Data challenges we face with customer data



Multiple versions of the truth

“

I see so many different versions of the same metric



Limited visibility into performance

“

I have no idea how my sales team performs vs. peer groups



Wasted time spent hunting for data

“

It takes four days to analyze my customers' usage patterns



Missing insights impairs decision-making

“

I have difficulty determining the value of the customers we are acquiring



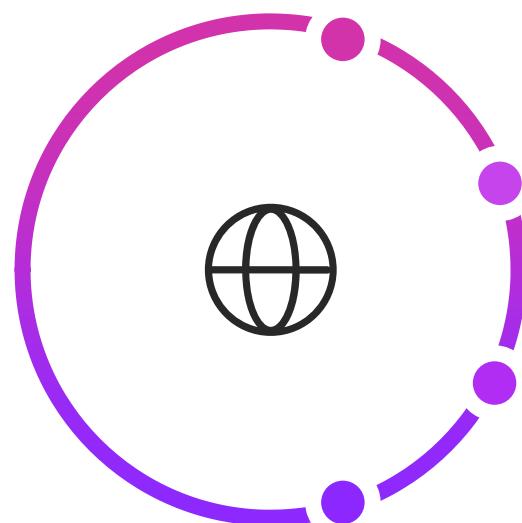
Inability to segment customers

“

I still can't tell how many business travelers subscribe

We set out to build a world-class data platform

Key objectives



Enable
new and existing
revenue streams

Improve
customer experience through
targeted communications

Reduce
operational cost

Promote
more informed
decision-making

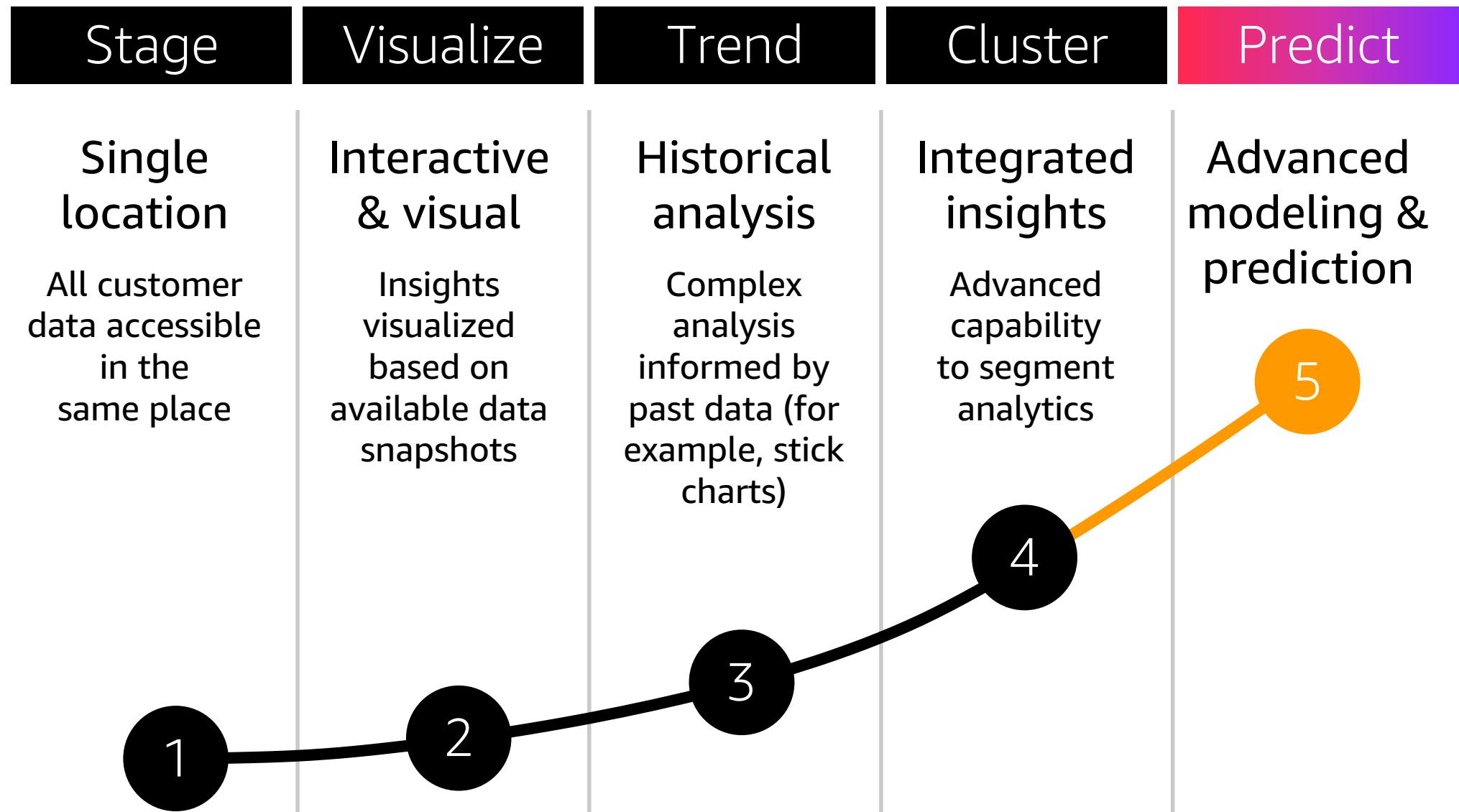
We ramped up an AWS Landing Zone



Created
an environmental
council to make key
architectural decisions

Partnered
closely with AWS
through the process

Platform evolution



End goal

Predictive analytics & machine learning

The pace of evolution for each customer lifecycle step is based on

Knowledge & Insight

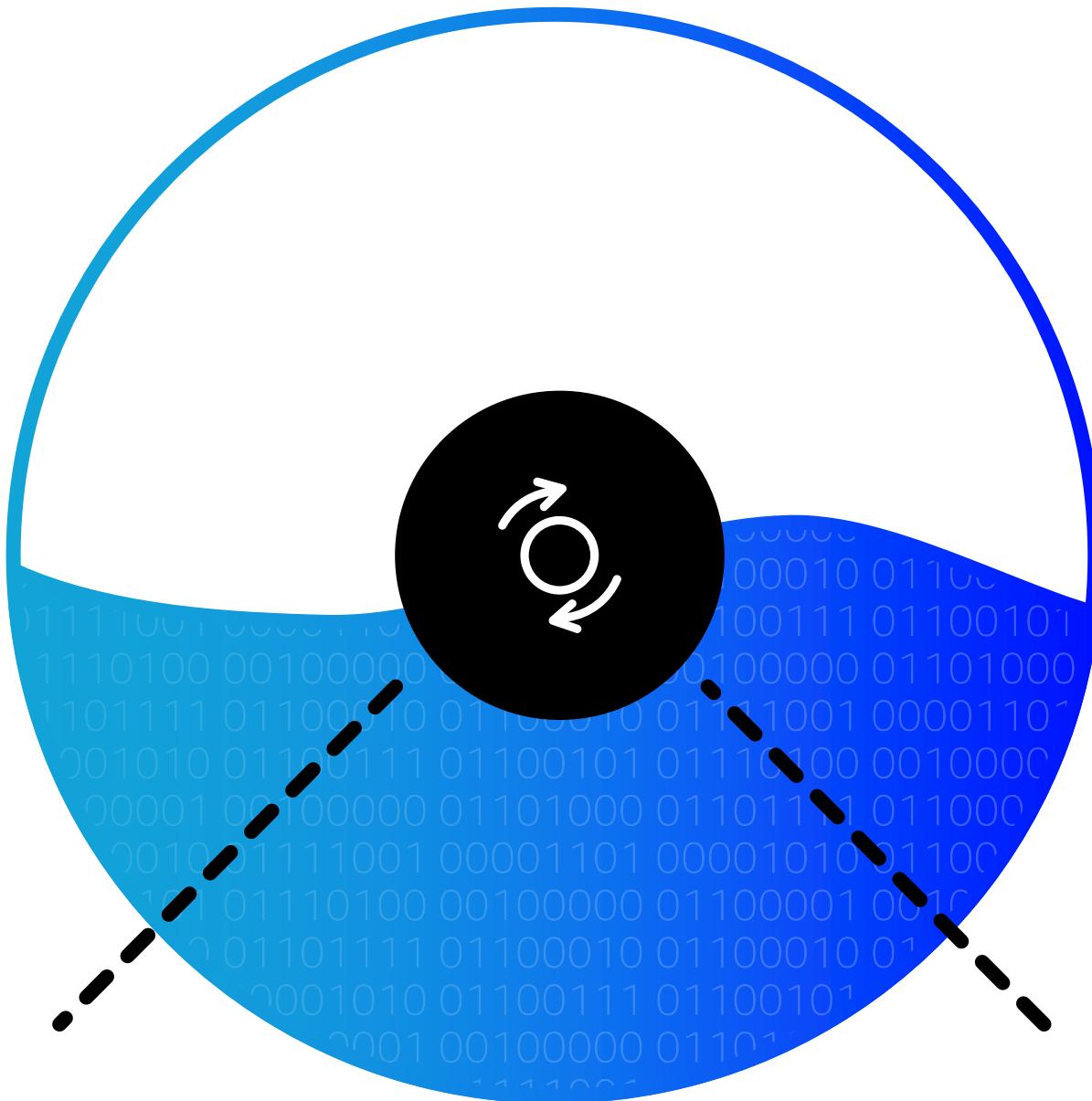
Risk & Compliance

News & Information



A **data lake** is a storage repository that holds a vast amount of raw data in its native format

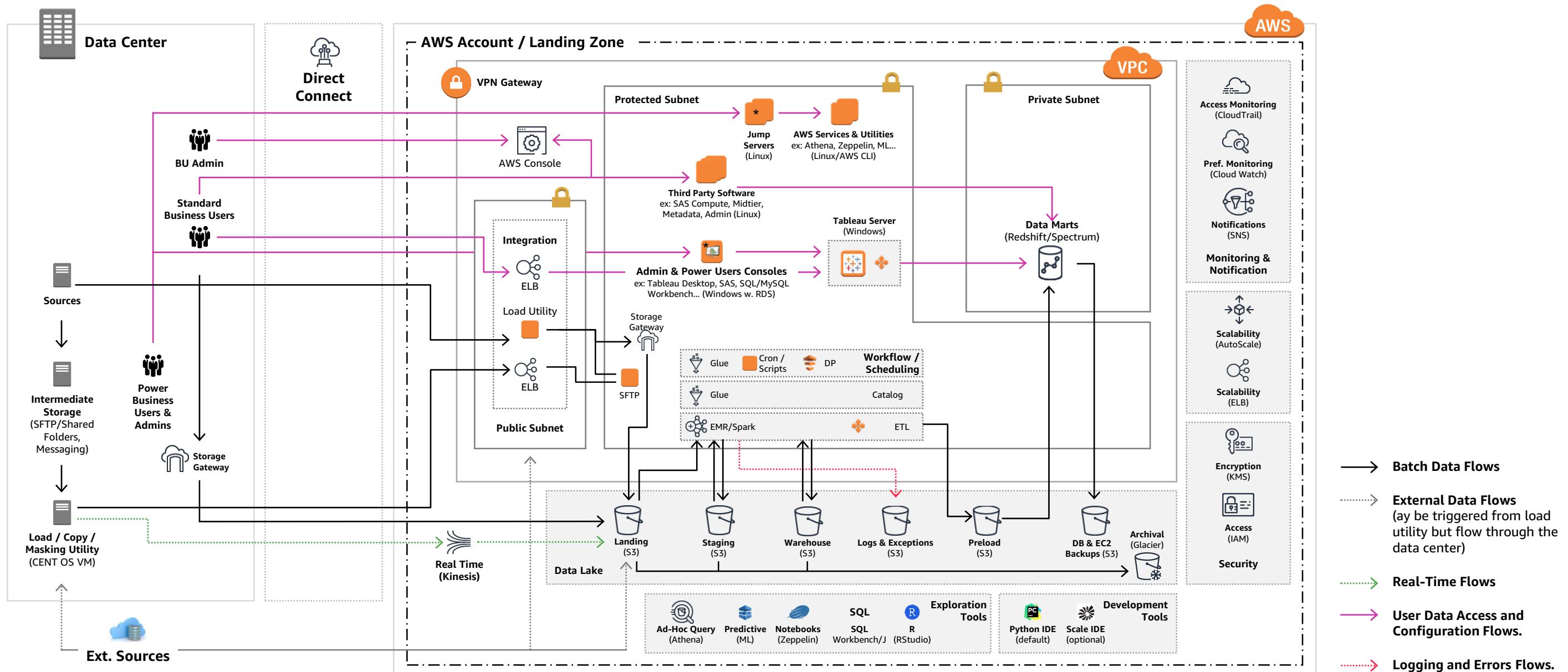
We ensure the data is loaded correctly each day. We aim to collect every relevant customer event that occurs for any of our members.



To make sense of this customer data, we will construct easy-to-understand **data marts**

that will normalize this data to a state where it can be consumed by self-serve tools and ad hoc analysis.

Reference architecture



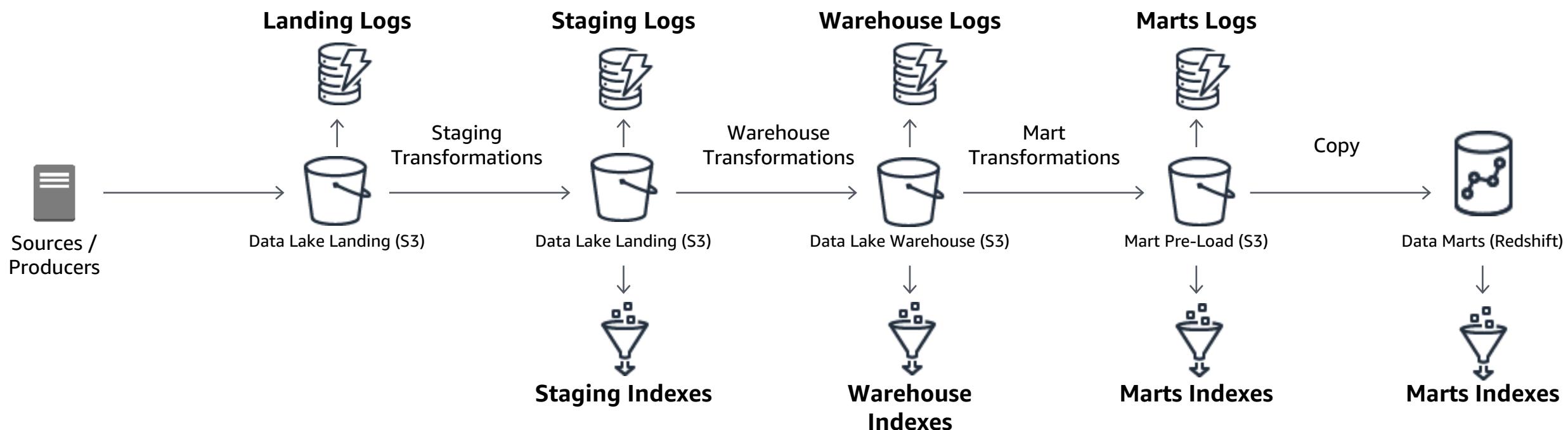
Creating our data platform

Use traditional data warehouse practices to process and store our data

Leverage the staging area where we prepare the source data, standardize it, catalog it and make it available as the data lake

Prepare intermediate tables in the warehouse layer and store them in aggregated data marts

Use the Copy command to push the data to Amazon Redshift



Amazon Athena Vs. Amazon Redshift



Access data in
Amazon Redshift

Best performance



All data is
stored in **Amazon S3**
but created as
external tables
in Amazon Redshift

Redshift Spectrum to query

More cost-effective measure



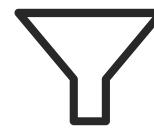
Use the **Copy**
command to copy
data from **Amazon S3**
to **Amazon Redshift**
as a physical table

Case-by-case basis

When performance
takes priority over cost

Analytics tools

connected directly
to Redshift create
custom dashboards
for our users

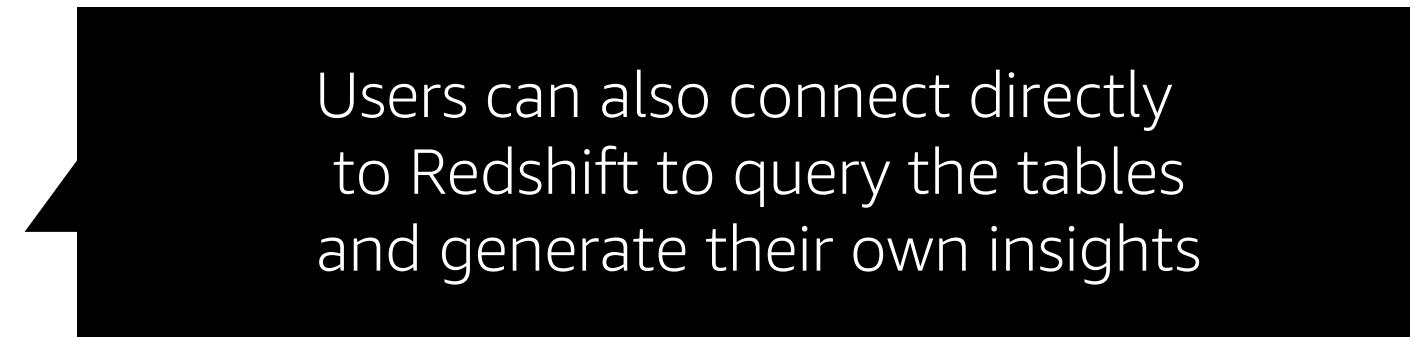


Harness the power of both the
Data Lake and the Data Marts

Create a BI layer where we join
multiple data sets to create custom
insights for our business partners

Analytics tools

connected directly
to Redshift create
custom dashboards
for our users

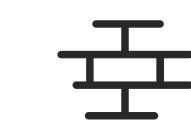


Users can download the tool of their choice, as long as it contains a JDBC/ODBC connection to Redshift



Users are armed with a data dictionary and a set of guidelines that allow them to harness the power of Redshift

Best Practices: Querying the data



**Select
a specific set
of columns**
instead
of Select *

**Conduct a
"Group By"**
instead of a
"Select Distinct"

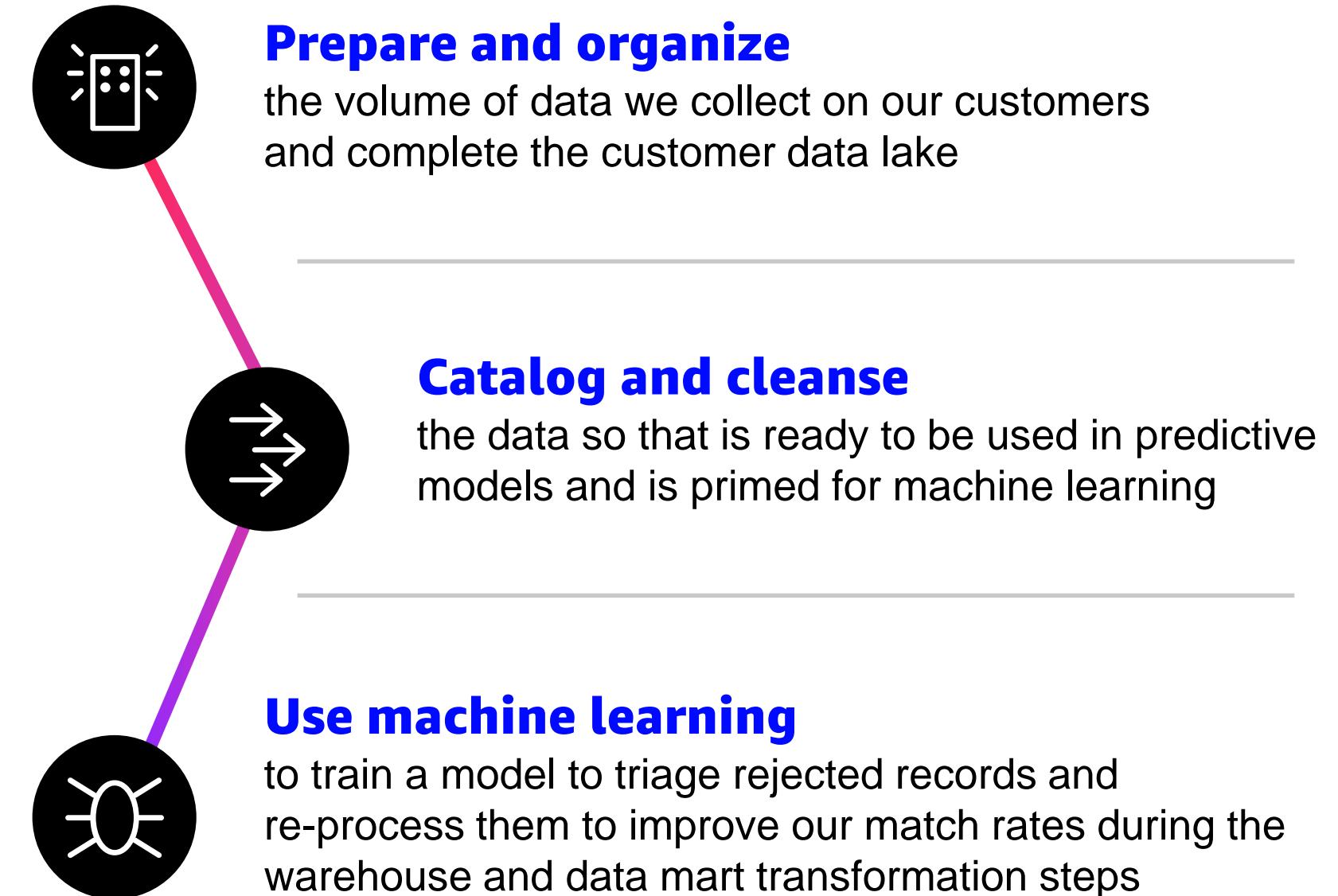
**Use
"With
Clauses"**
instead of
"Sub Queries"
as it will create
a temp table
which can
be re-used
or referenced
multiple times

**Amazon
Redshift is
best at
aggregated
queries**
use aggregate
functions
whenever
possible

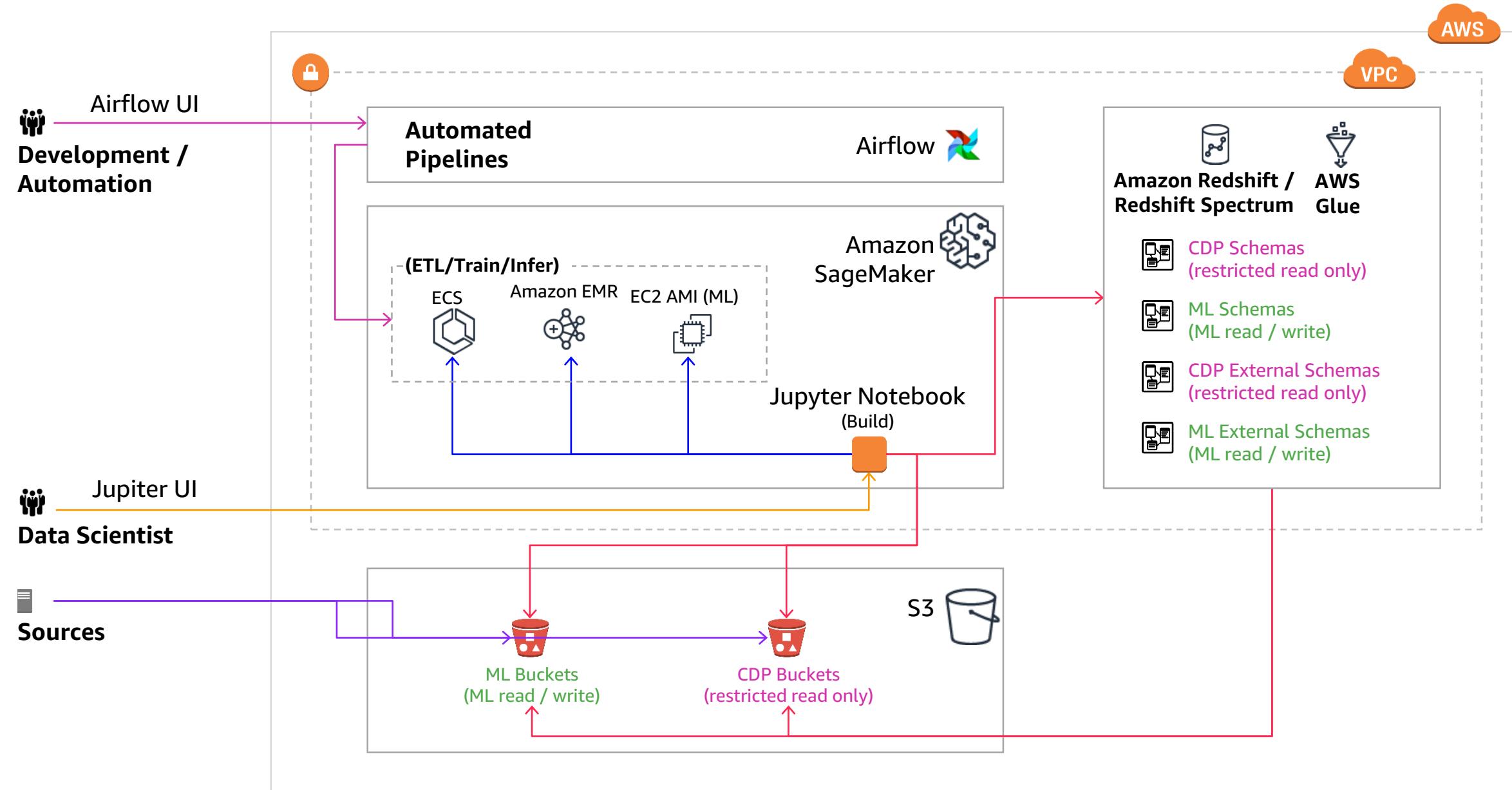
**When
investigating
the data,
use filters**
to make the
data set size
smaller or limit
the number
of records to
the first 1,000
records or so

**Always use
the data
partitions**
when querying
the data

What is next:
Our platform
is entering into
the final stages
of evolution



Reference machine learning architecture



Amazon SageMaker: “Analyst Sandbox” within Amazon S3



Query the
Data Lakes and
Data Marts
in **Amazon Redshift**



Create curated
data sets and
store them in **Amazon S3**
to be cleansed,
labeled and provided
as input into the data mart

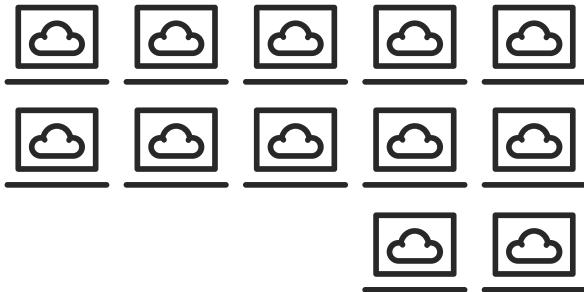


Access pre-setup
Jupyter notebook
instances to
both query and
prepare data

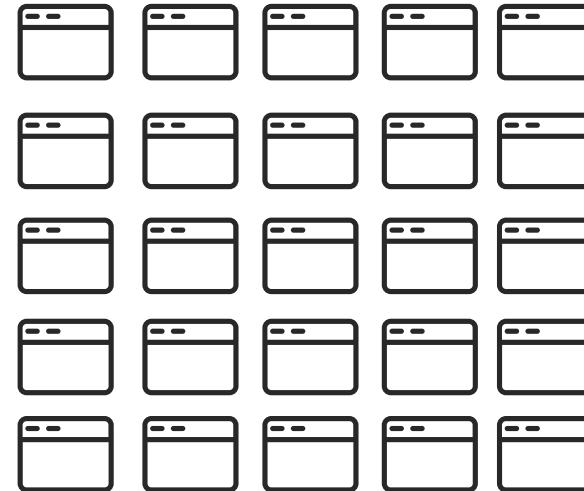
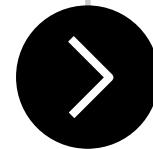


Amazon SageMaker will be responsible for managing the AWS resources
behind the scenes while the analysts focus on the key modeling activities

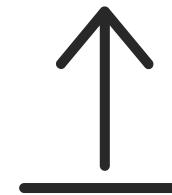
Key Results



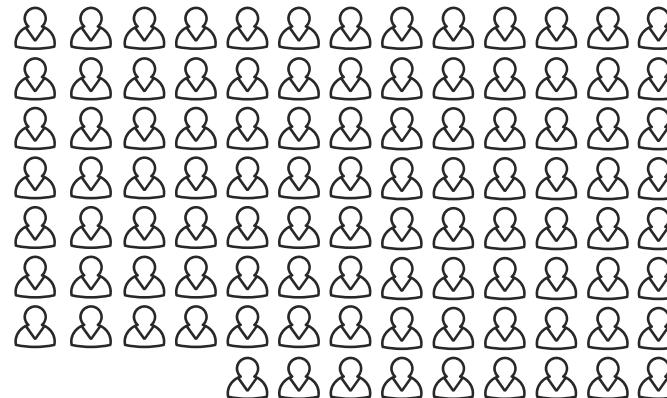
We have spun up
12
different dashboards
to support our stakeholders



These dashboards
are powered by
30
different fact & dimension tables in Redshift



Across the whole of the warehouse we hold upwards of
118 TB of Data



The platform is currently accessed by
100+ users



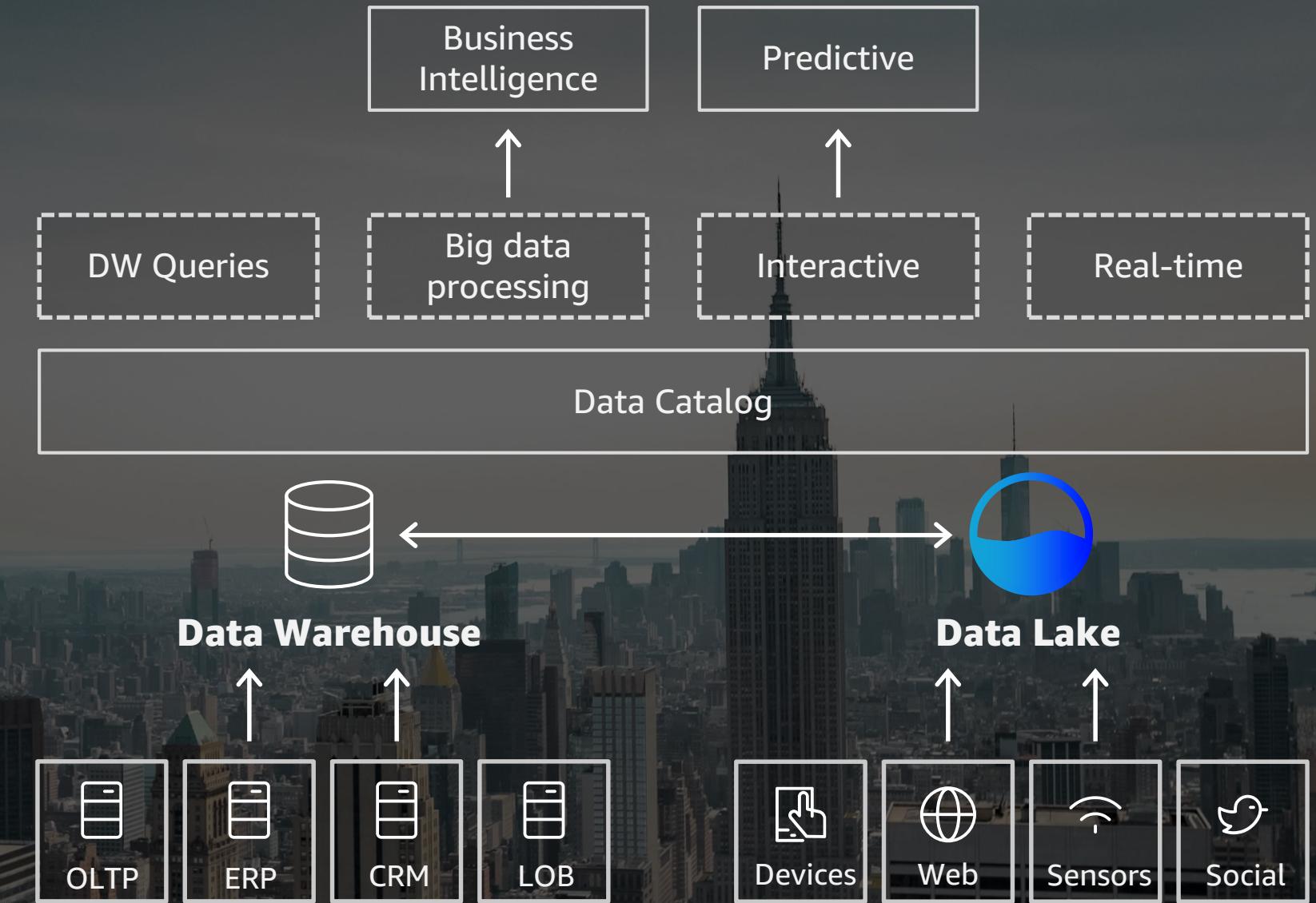
Thank you.

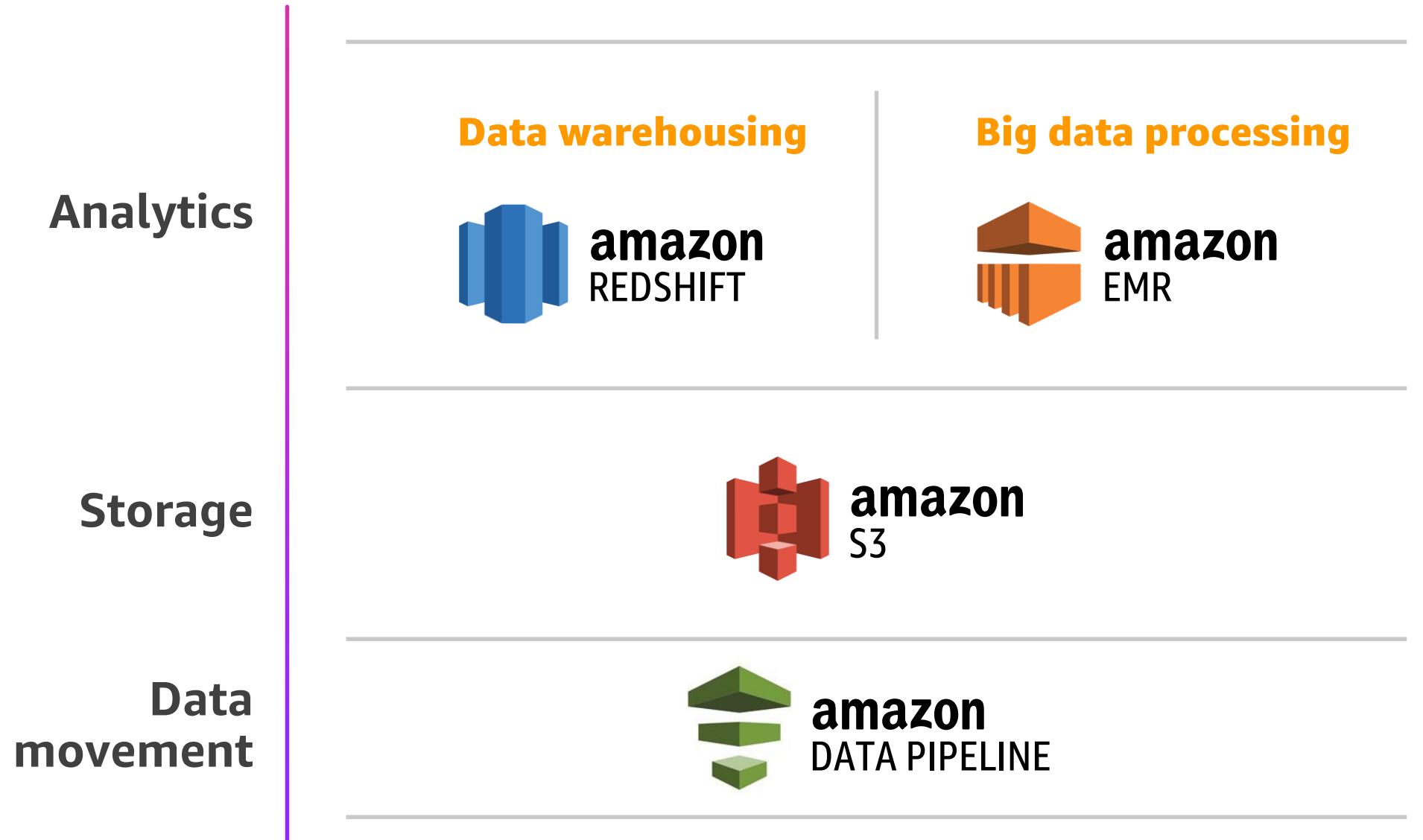
Bucket strategy matters

Glue and EMR/Spark are great tools for data transformation

Redshift provides a central location to view high-performance reports and run queries with data in S3

Scale on a per job or per query basis





aws
2012
Analytics
Portfolio

AWS databases and analytics

Broad and deep portfolio, built for builders

Business Intelligence & Machine Learning						AWS Marketplace 250+ solutions
 QuickSight	 SageMaker	 Comprehend	 Rekognition	 Lex	 Transcribe	 DeepLens
Databases	Analytics			Blockchain		
 QLDB Ledger Database NEW	 Neptune Graph	 Redshift Data warehousing	 Athena Interactive analytics	 Managed Blockchain NEW	730+	Database solutions
 ElastiCache Redis, Memcached	 DynamoDB Key value, Document	 EMR Hadoop + Spark	 Kinesis Analytics Real-time	 Blockchain Templates	600+	Analytics solutions
 Aurora MySQL, PostgreSQL	 Timestream NEW Time Series	 Elasticsearch service Operational Analytics			25+	Blockchain solutions
 RDS MySQL, PostgreSQL, MariaDB, Oracle, SQL Server	 RDS on VMWare NEW					
S3/Glacier	Lake Formation Data Lakes NEW	Data Lake	Glue ETL & Data Catalog		20+	Data lake solutions
Data Movement					30+	solutions
Database Migration Service Snowball Snowmobile Kinesis Data Firehose Kinesis Data Streams Data Pipeline Direct Connect						



There is **more data** than people think

Data

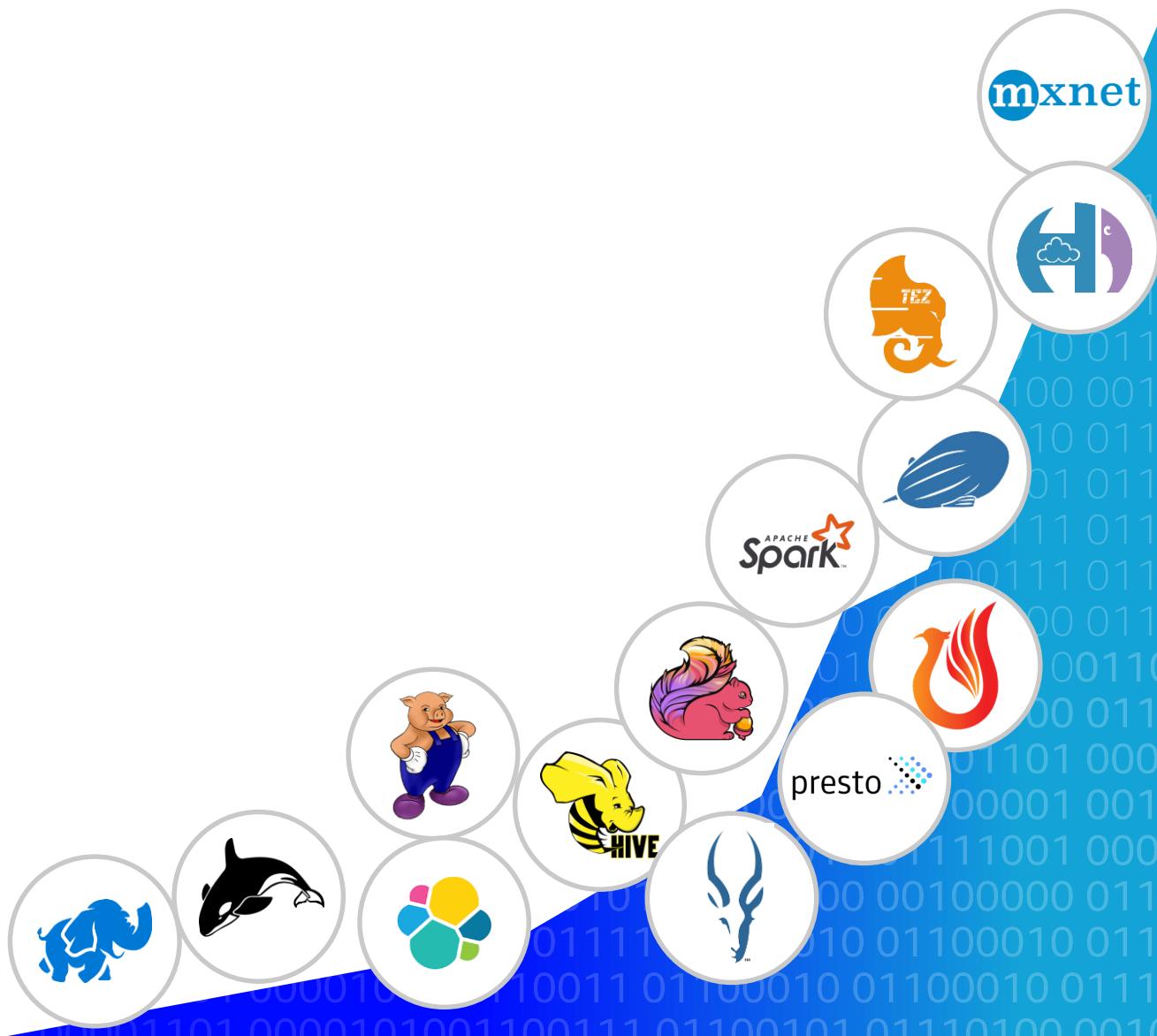
grows
>10x
every 5 years

Data platforms need to

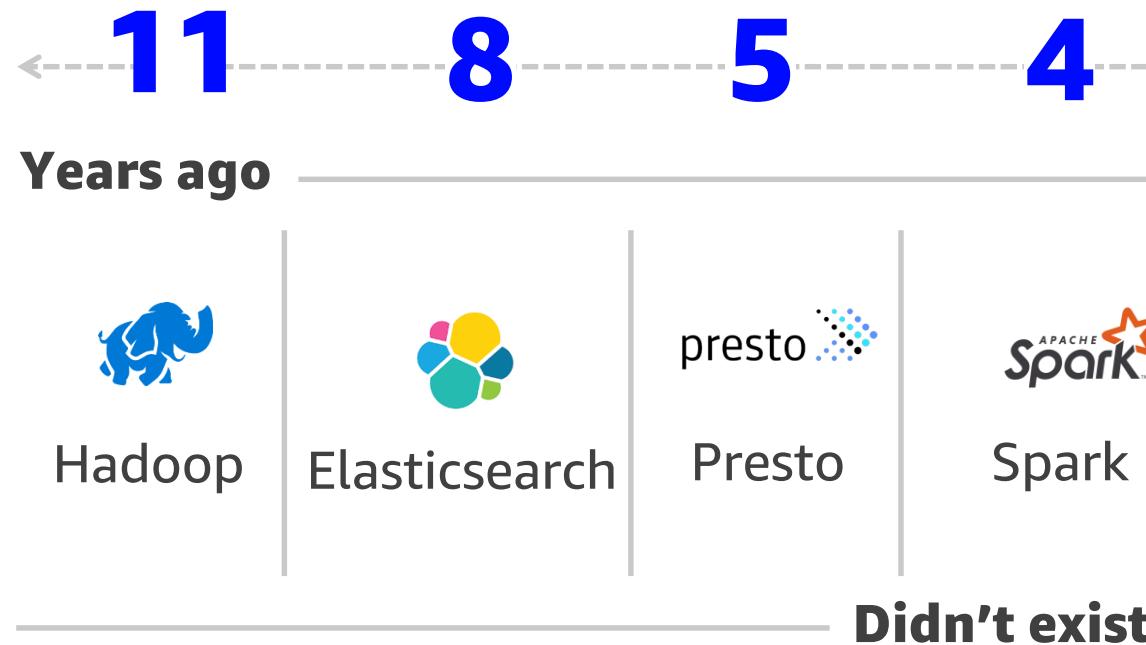
live for
15
years

scale
1,000x

* IDC, Data Age 2025: The Evolution of Data to Life-Critical Don't Focus on Big Data, Focus on the Data That's Big, April 2017.

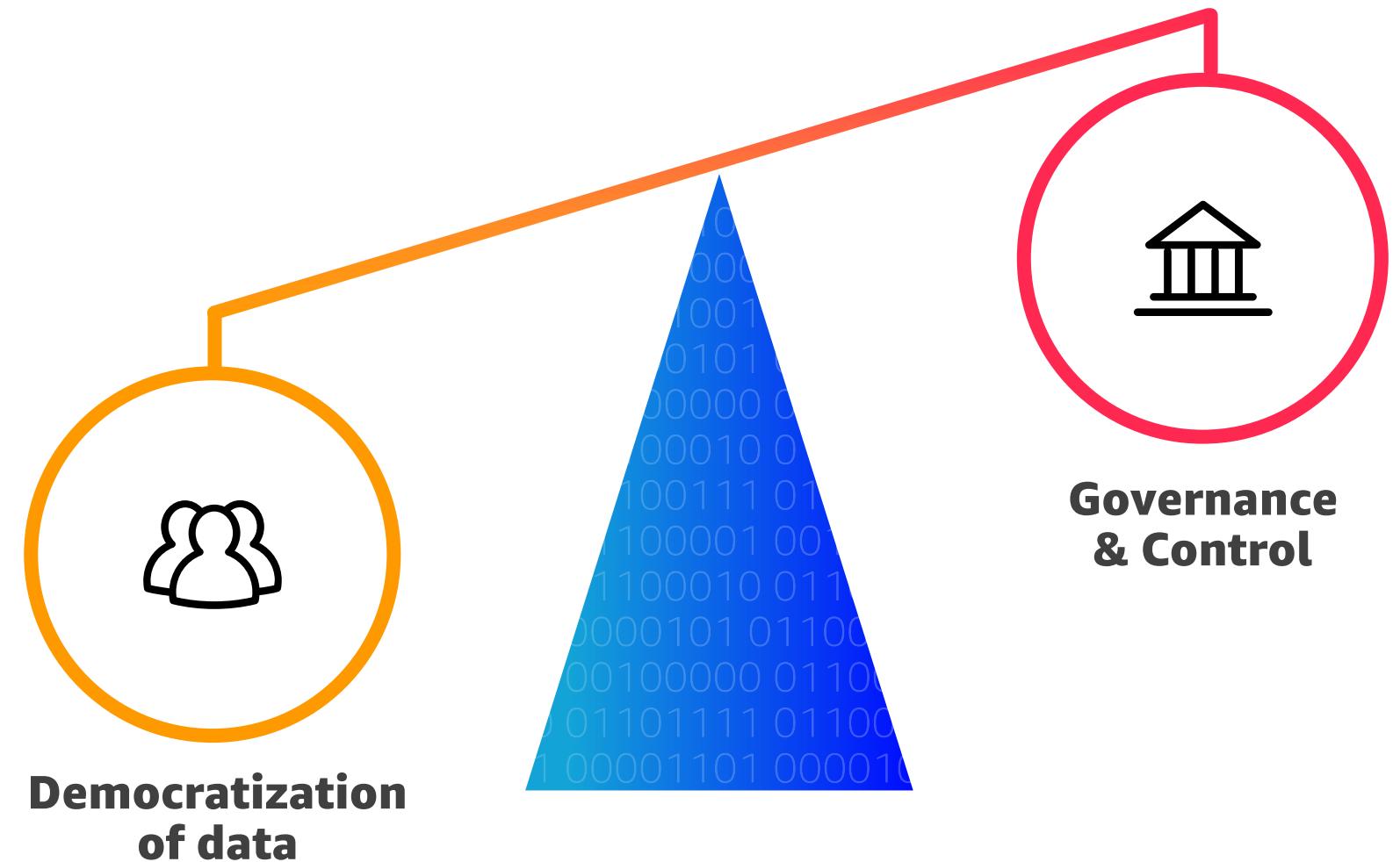


There are **more ways** to analyze data than ever before



There are **more**
people working
with data than
ever before

How do I provide democratized
access to data to enable
informed decisions while at the
same time enforce data
governance and prevent
mismanagement of the data?



Managing the evolving data landscape



Flexible
Open APIs
and open
data formats



Choice
Use the best
analytic tool for
the job, without
data movement



Scale
Platforms
that scale up
to 1,000x



Secure
Full auditability,
access controls,
and data
governance

Amazon Redshift

The 4 things that matter most



Speed



Simplicity



Scale



Security

More customers use

Amazon Redshift

for their data
warehouse
workloads than
anyone else





iflix



LOYALTY LAB
every company deserves loyal customers

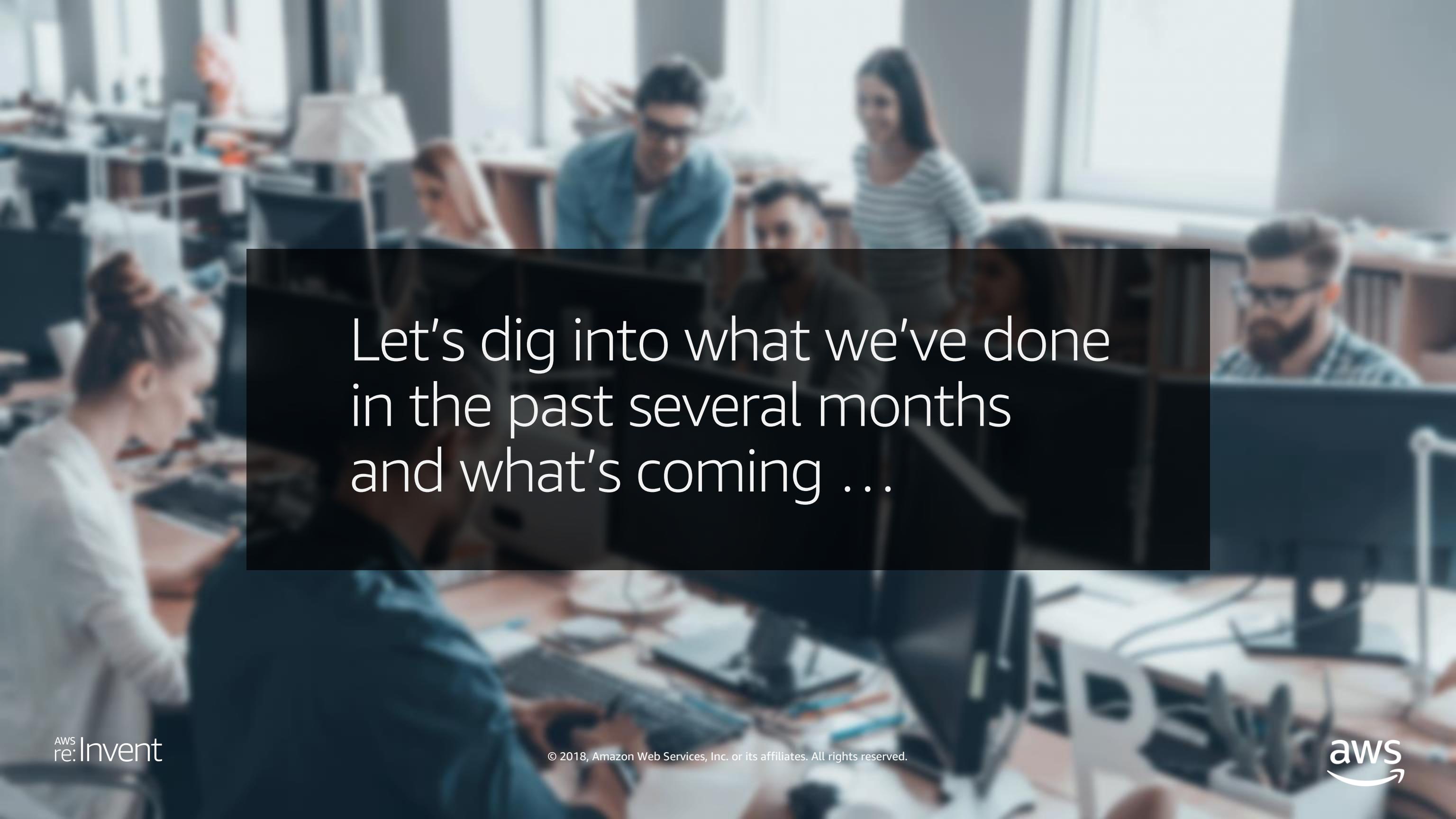


AWS
re:Invent

© 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved.

aws





Let's dig into what we've done
in the past several months
and what's coming ...

Redshift is **growing fast** and **innovating faster**

Since we last spoke ...

Automatically enabled short query acceleration

Support for lateral column alias reference

New Quick Starts

New CloudWatch metrics

Customized Recommendations with Advisor

Current and trailing tracks for release update

Federated authentication with single sign-on

Improved performance for commits

COPY from Parquet and ORC file formats

Support for Parquet and ORC in Kinesis Data Firehose

Improved workload management console experience

Query Editor

Support for late-binding views

SQL Scalar user-defined functions

Integration with AWS Glue

Support for Nested Data with Spectrum

Spectrum support for DATE data type

Improved performance for UNION ALL queries

Free upgrade from DC1 to DC2 RIs

Query monitoring rules (QMR)

Support for Zstandard high compression encoding

Query processing improvements

220+

◆ features and
enhancements
released*

Support for Python UDF logging module

Enhanced VPC routing

Additional Spectrum regions

Support for Scalar JSON and Ion data types

Late materialization for faster query processing

Support for DATE data type with Spectrum

Short Query Acceleration

Utilization reports

Machine learning integration to accelerate dashboards and interactive analysis

Improved resource management for memory-intensive queries

Faster string manipulation

Automatically hopping queries without restarts

Support for uppercase column names

Result Caching for Repeat Queries

Support for LISTAGG DISTINCT

Support for ORC and Grok file formats

Integration with QuickSight

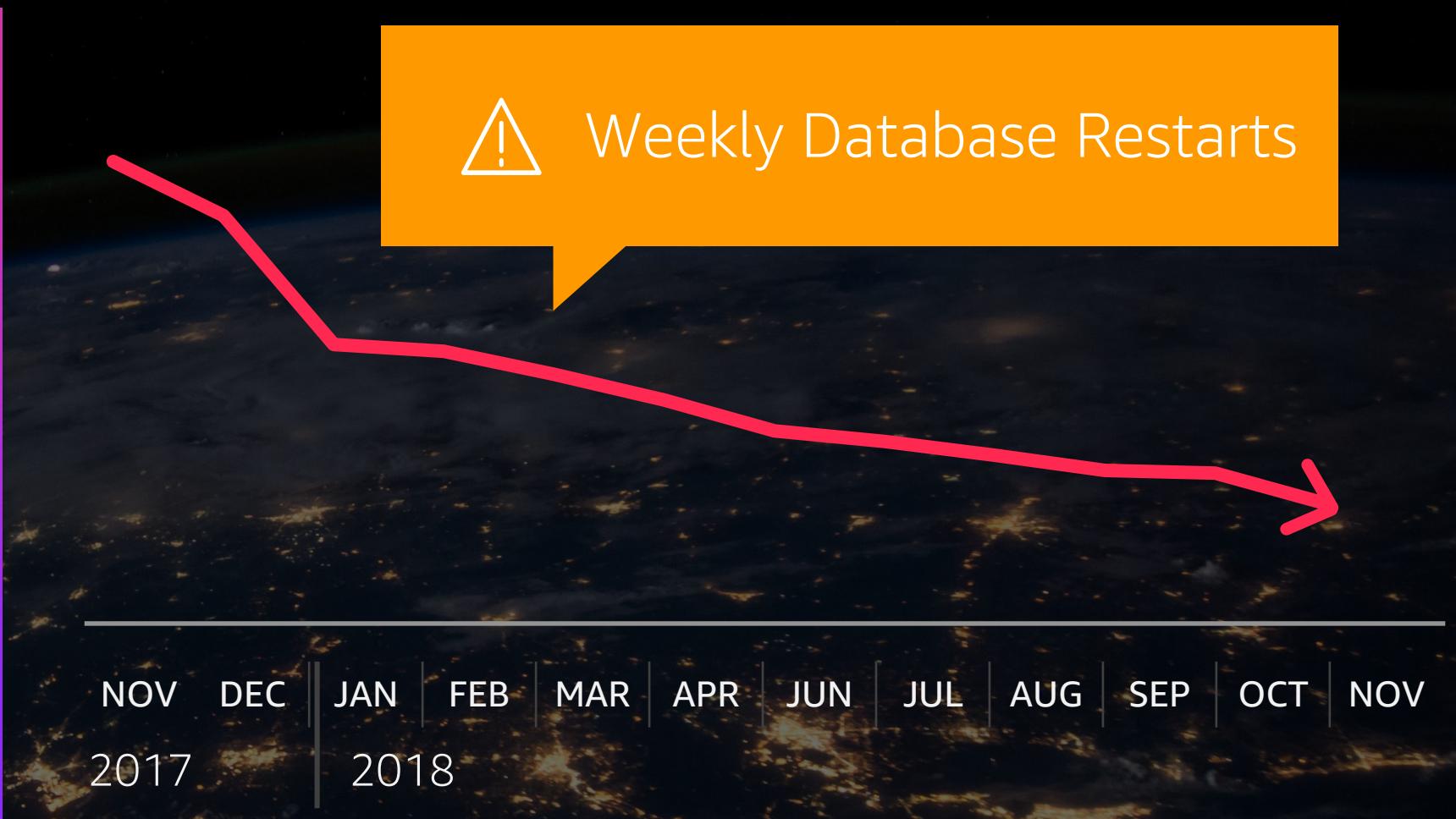
DMS support with Redshift

3.5x Improved Throughput

Improved performance for repeat queries

*Since re:Invent 2017

Improvements
in availability
since we
last talked



Compiled code cache

Support for lateral column alias reference

Query planning

Cluster resize operations

Short query acceleration

Improvements for the COPY operation when ingesting data from Parquet and ORC formats

Commit processing enhancements

COPY operation when ingesting data from Parquet and ORC formats

Query processing improvements

Resource management for memory-intensive queries

Faster string manipulation

Query rewrites that pushdown selective joins into a subquery

Result caching

Queries operating over CHAR and VARCHAR columns

2x the number of tables in a cluster

Late materialization

Single-row inserts

Complex EXCEPT subqueries

DC2 nodes

Hash join memory utilization optimizations and cache line prefetching

Expressions on the partition columns of external tables

Queries that refer to stable functions with constant expressions

Performance improvement for queries that refer to stable functions over constant expressions

Joins involving large numbers of NULL values in a join key column

*Since re:Invent 2017

Increases in performance in **real-world workloads**

17x

for repetitive
queries

10x

for bulk-deletes

3x

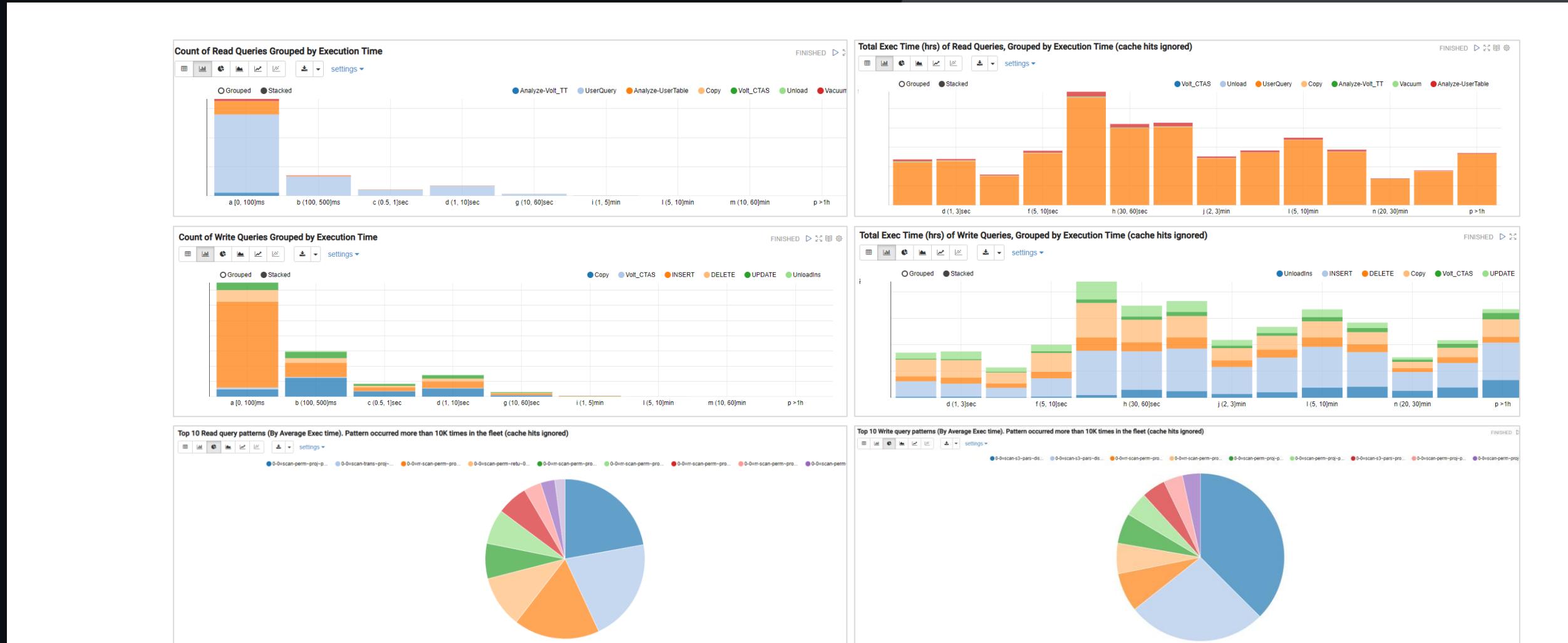
for single-row
inserts

2x

for commits

How do we
improve real-world
performance?

How we leverage **fleet telemetry**



66

Redshift's query performance and scalability has been increasing, even though our data has grown. In the last 10 months, we have seen commit performance **increase by 500%** without any increase in cost.

99

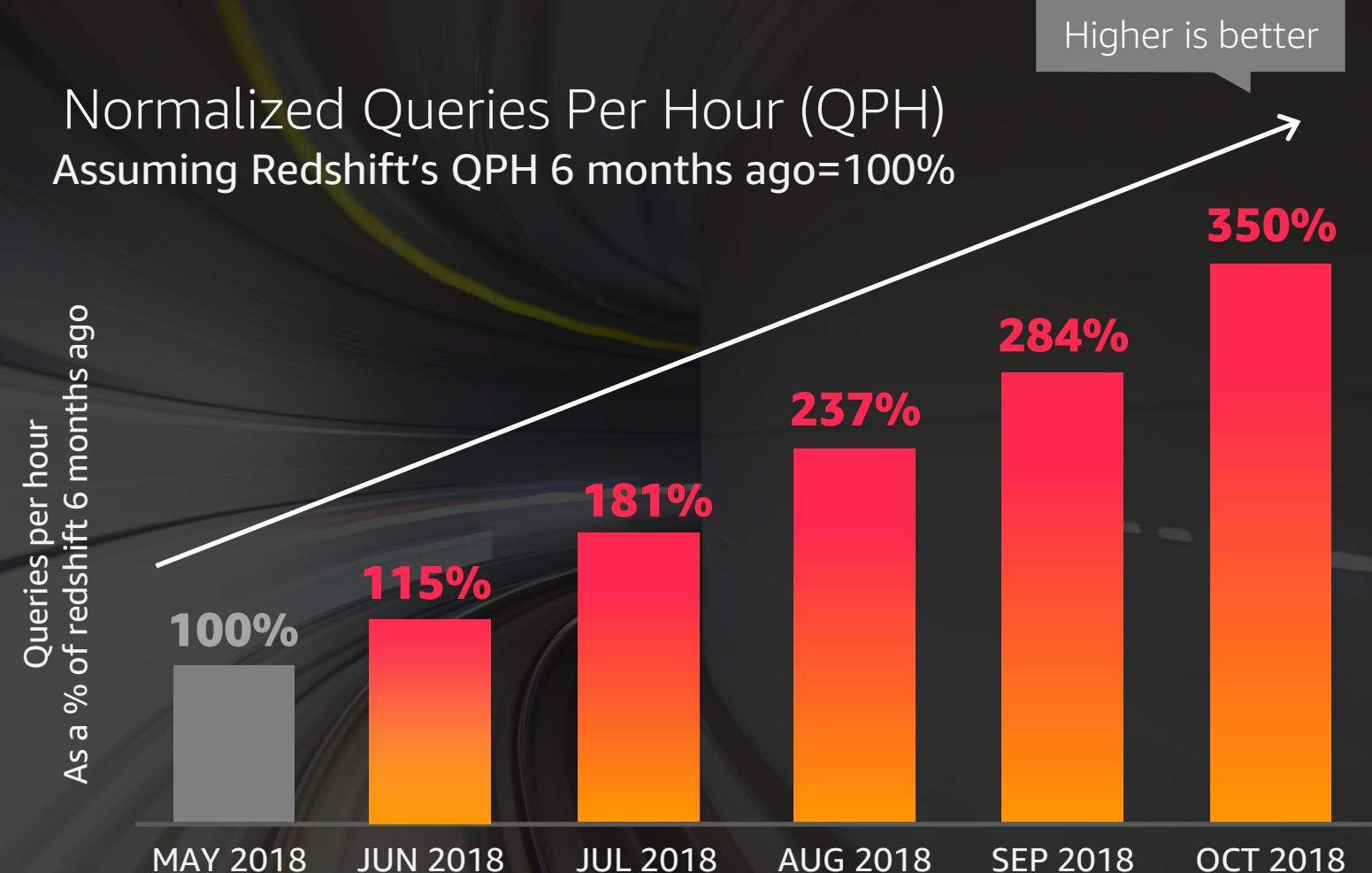
- **Minero Aoki**

Senior Data Engineer, Cookpad Inc.



cookpad

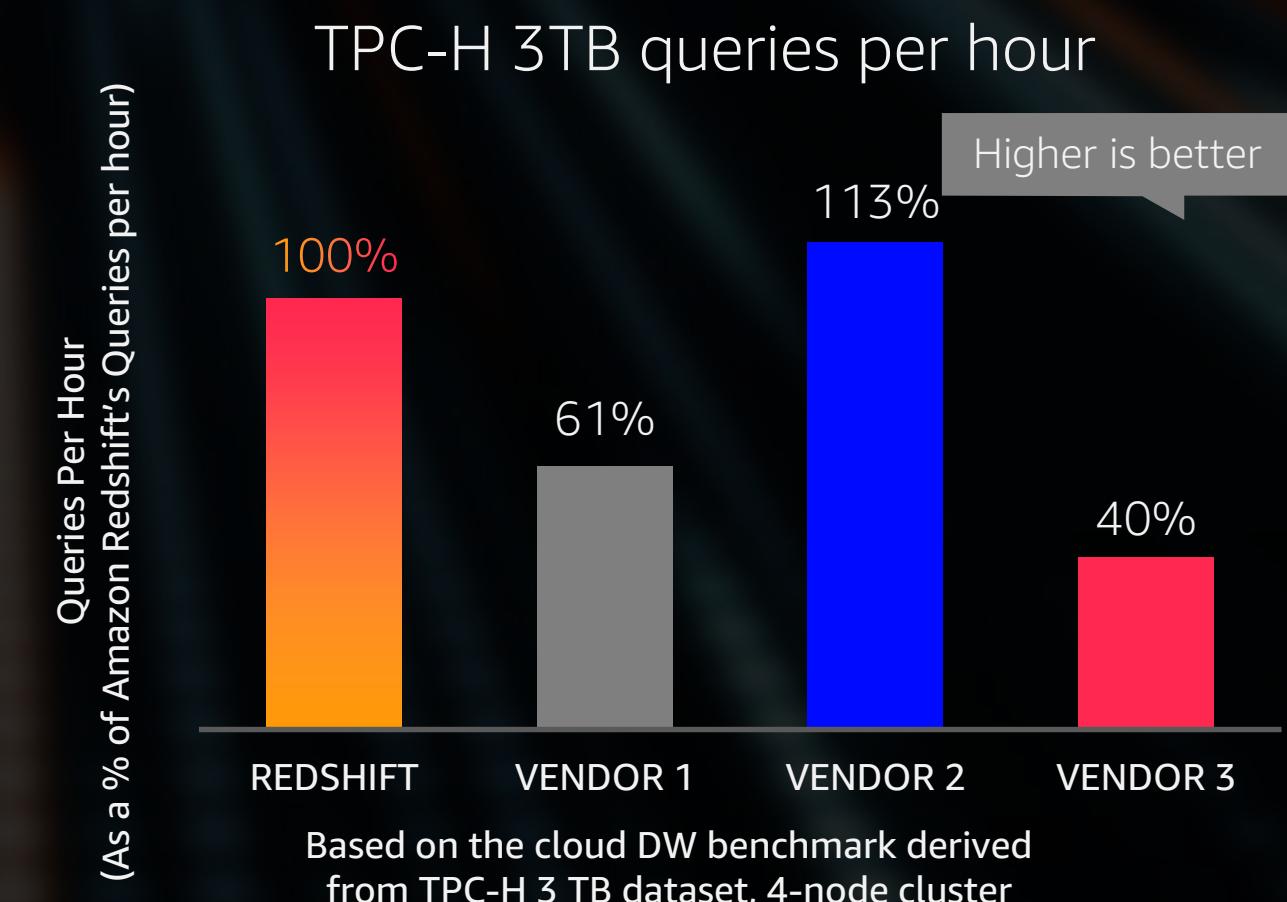
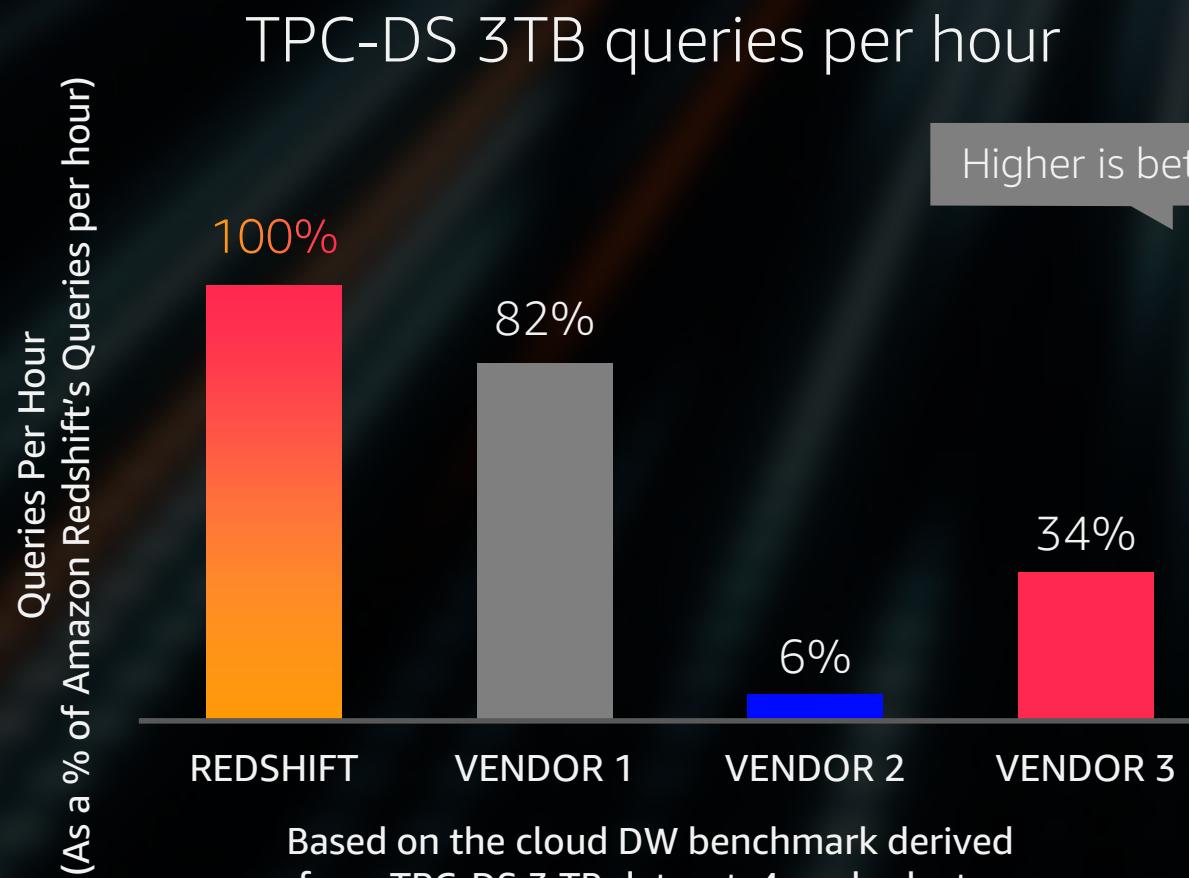
Amazon Redshift is now over 3x faster on standard benchmarks than 6 months ago



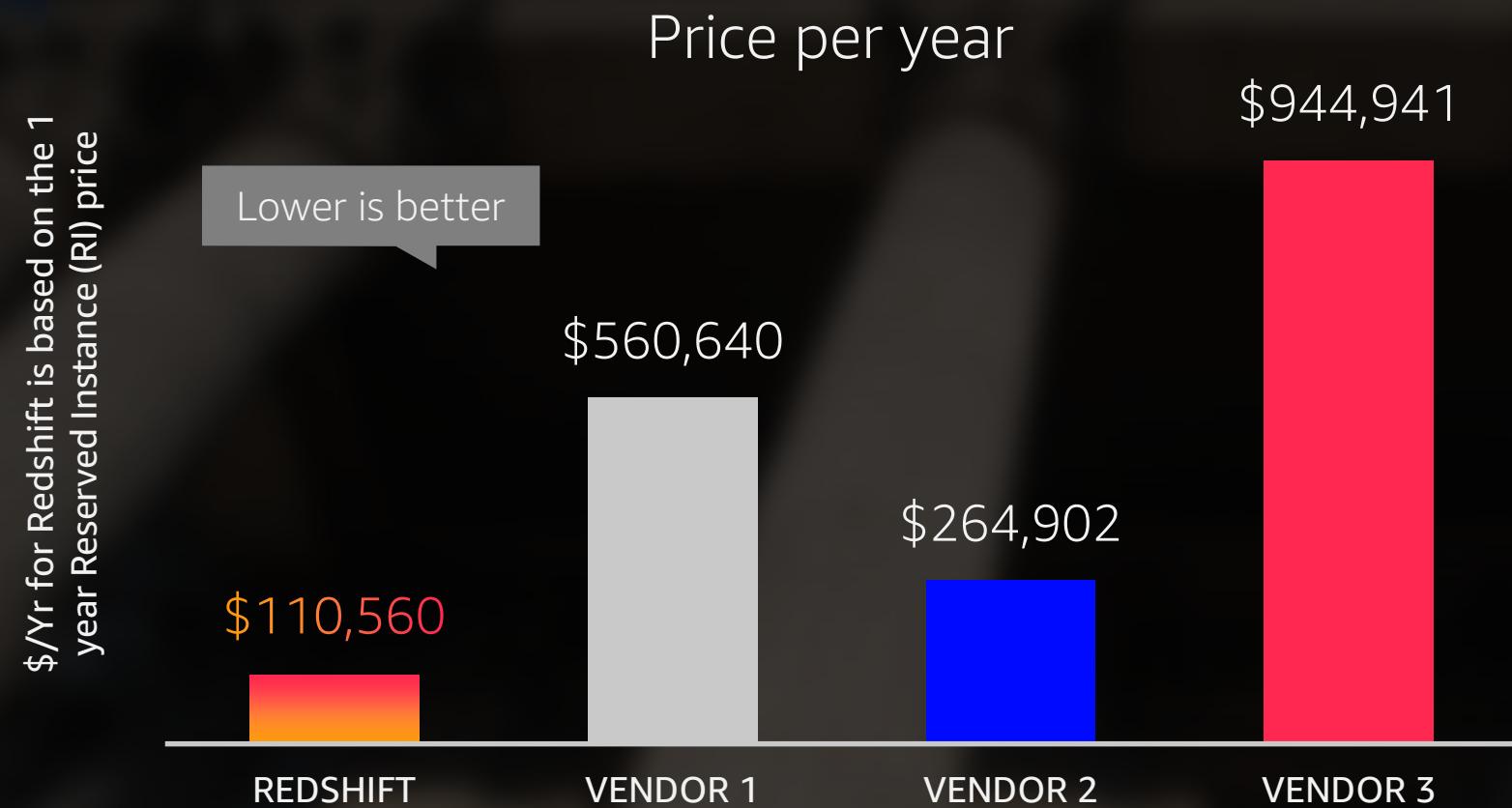
Amazon Redshift

up to

16x faster



Amazon Redshift is the most cost-effective cloud data warehouse



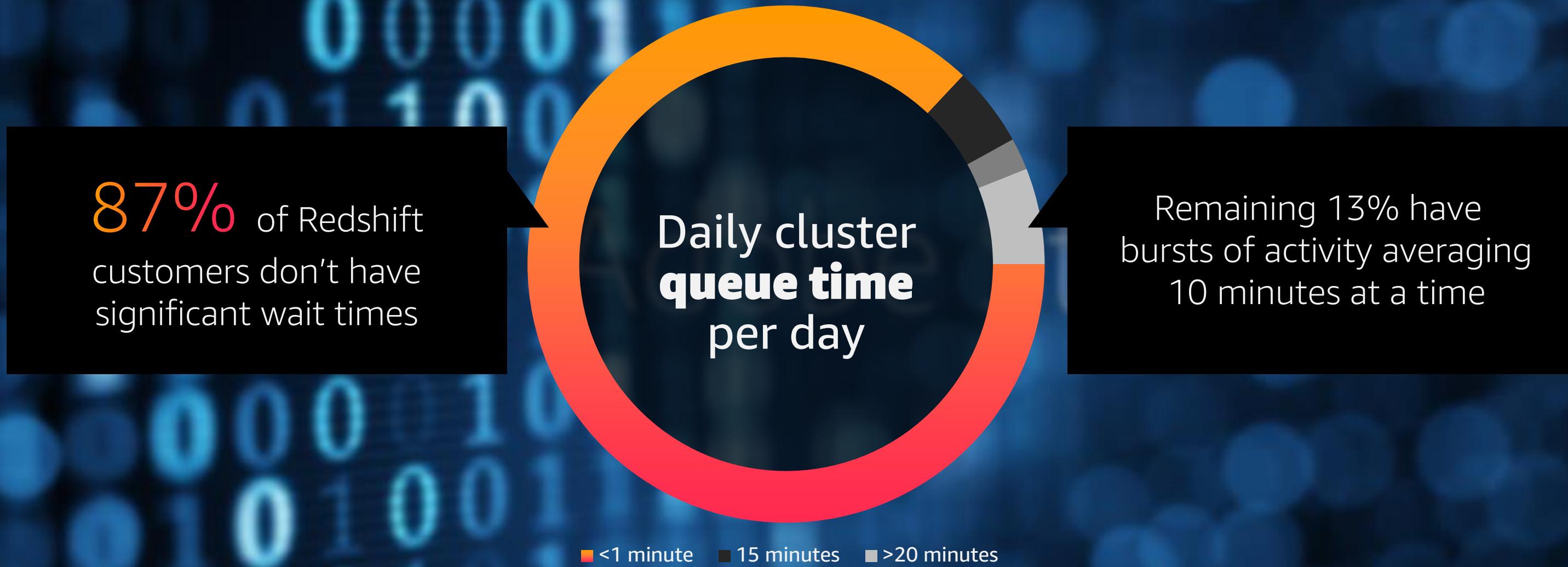
The only data warehouse with reserved instances saving



up to **75%**

The best price-to-performance

Fleet telemetry on query wait times



Concurrency Scaling for bursts of user activity (Preview)

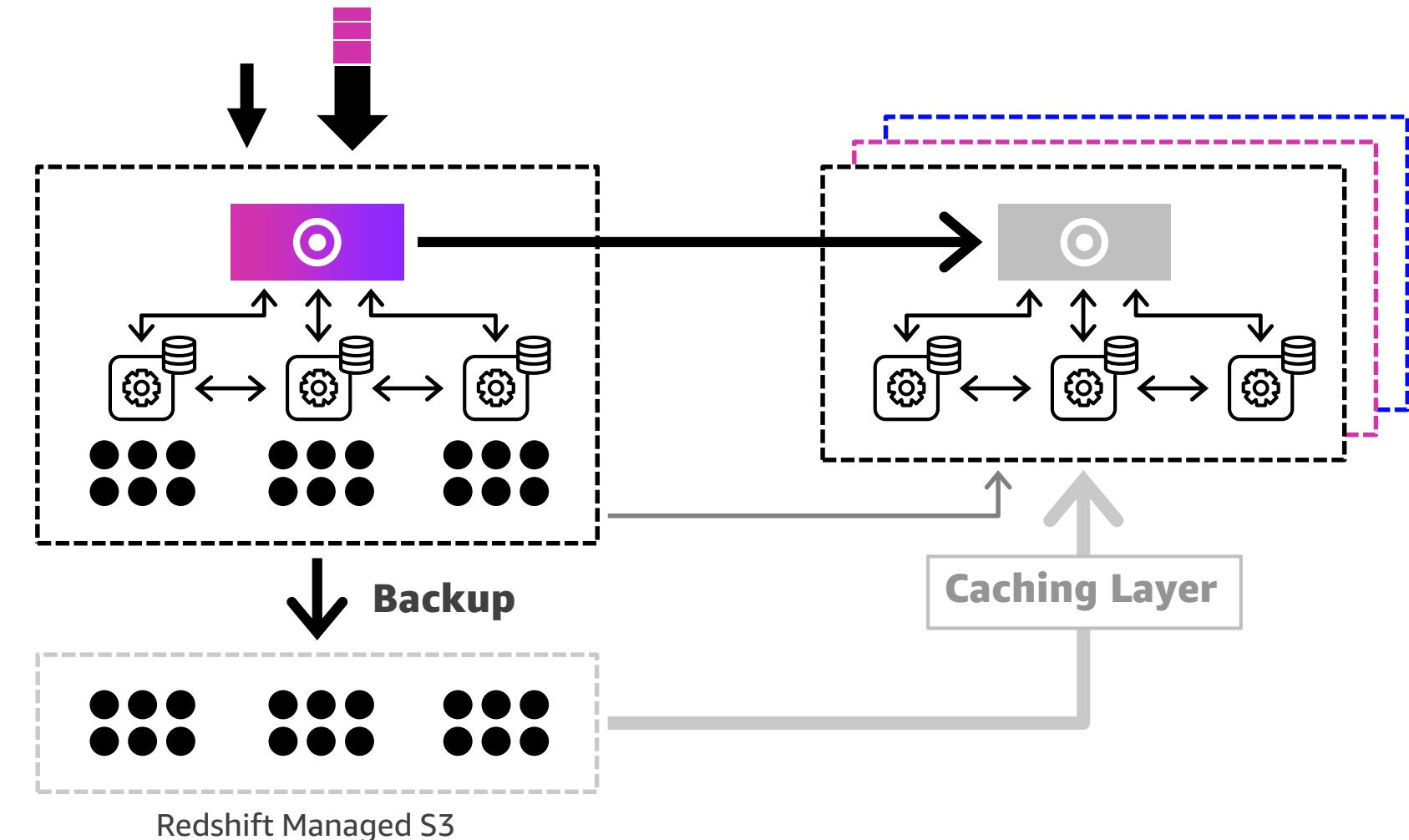
New!

Automatically creates more clusters on-demand

Consistently fast performance even with thousands of concurrent queries

No advance hydration required

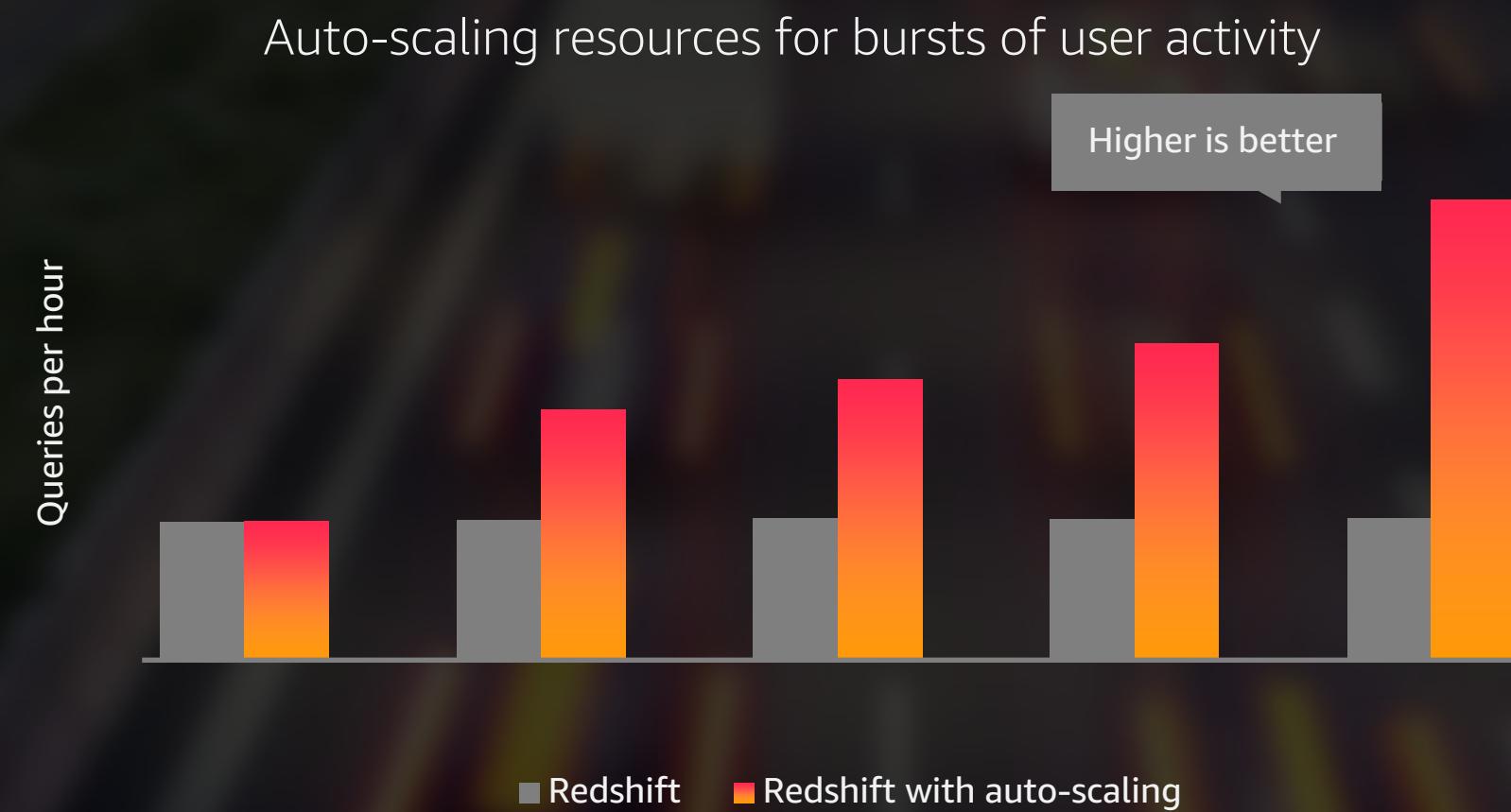
Quickly scale to serve changing query workload



Results with Concurrency Scaling

Concurrency Scaling is free for more than 97% of Redshift customers.

For every 24 hours your main cluster is in use, we'll provide a one-hour credit for concurrent cluster usage.



Redshift Elastic Resize (GA)

New!

Scale up and down in minutes

Adds additional nodes

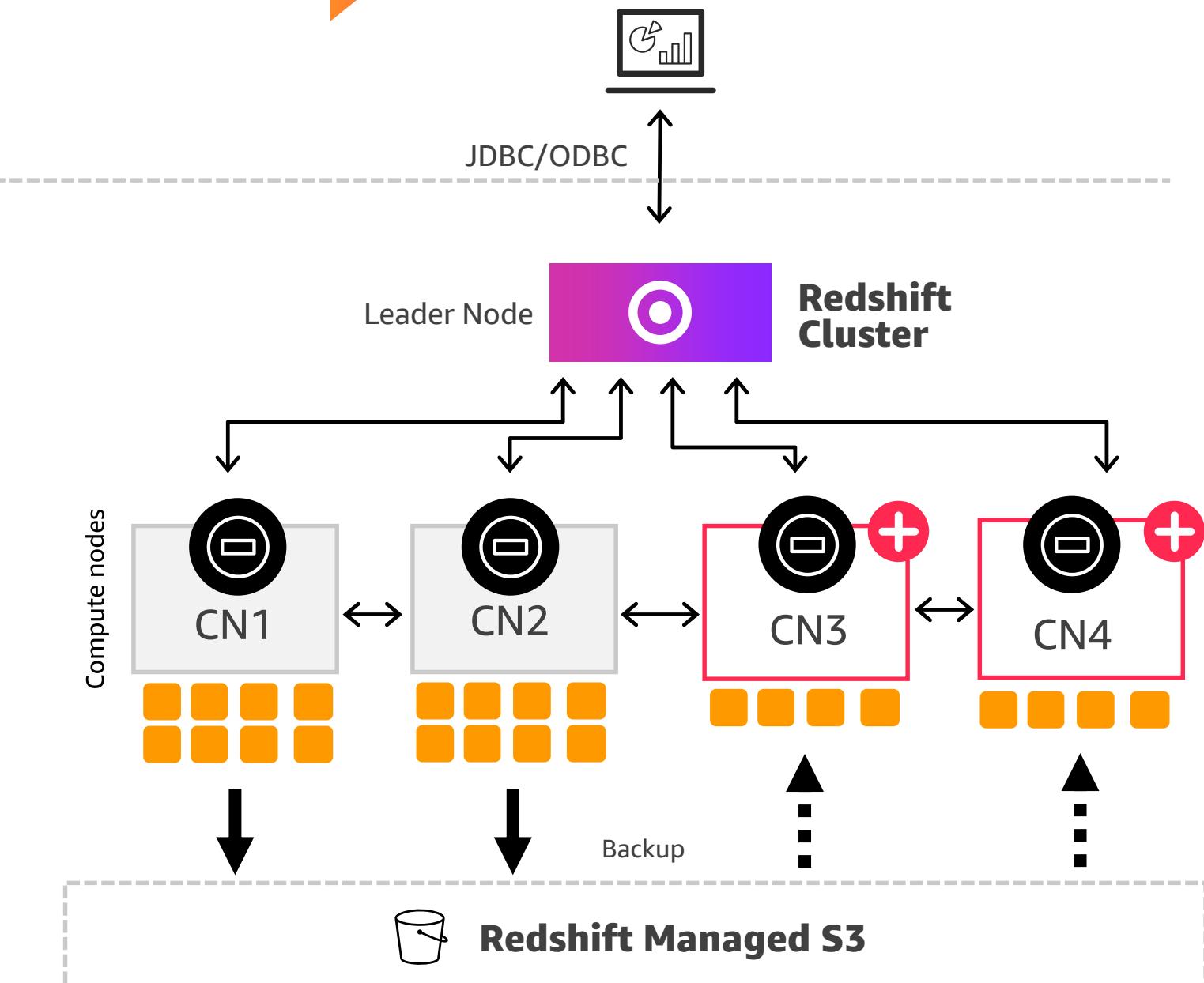
to Redshift cluster

Minimal
transition time

Distributes data

across new configuration in minutes

Scale compute and storage on-demand



Efficiency of backup performance

CloudWatch metrics for
Workload Execution
Breakdown

Automatic vacuum delete

Cluster resize

Enhancements to
VACUUM DELETE

Manage components
of a multi-part query
in the AWS console

Redshift Advisor for best practice recommendations

Current and trailing tracks
for release updates

Lateral column alias reference

Stream real-time data in
Parquet or ORC formats
using Kinesis Data Firehose

Free upgrade from DC1 RIs to DC2

CloudWatch metrics for Query
Throughput, Query Duration

CloudWatch metrics
for Query Duration
by WLM Queues

Cluster resize operations

Query Monitoring Rules (QMR)
now support 3x more rules

CloudWatch
Query Runtime Breakdown metric

Query Editor

Improvements to simplicity

**Short query
acceleration is
self-optimizing**

**DISTSTYLE AUTO
distribution style**

***Since re:Invent 2017**

Redshift Query Editor

Launched in October!

The screenshot shows the AWS Redshift Query Editor interface. On the left, there's a navigation sidebar with options like Redshift dashboard, Clusters, Query editor (which is selected), Saved queries, Snapshots, Security, Parameter groups, Workload management, Reserved nodes, Advisor (Beta), Events, Connect client, and What's new. The main area has tabs for Internal schema, My Spectro... (which is selected), New Query 1, and New Query 2. Below these tabs is a code editor containing a SQL query:

```
1 /* Join external and internal table */
2
3 SELECT
4     myexternalschema.sales.eventid,
5     sum(myexternalschema.sales.pricepaid)
6 FROM
7     myexternalschema.sales,
8     myinternalschema.event
9 WHERE
10    myexternalschema.sales.eventid = myinternalschema.event.eventid
11    AND myexternalschema.sales.pricepaid > 50
12 GROUP BY
13     myexternalschema.sales.eventid
14 ORDER BY
15     1 DESC;
```

Below the code editor are buttons for Run query, Save as, Save, and Clear. The results section shows a table titled "Query results" with a note "Query completed in 4.679 seconds". It includes buttons for Download CSV, View execution, and a link "Showing row(s) 1 - 100". The table has two columns: eventid and sum. The data is as follows:

	eventid	sum
1	8798	10191.0
2	8797	12128.0
3	8796	21483.0
4	8795	7930.0
5	8794	11481.0
6	8793	5935.0
7	8792	8903.0
8	8791	639.0
9	8790	7855.0
10	8789	7020.0
11	8788	9184.0
12	8787	6033.0
13	8786	6271.0
14	8785	9258.0
15	8784	17713.0
16	8783	14354.0
17	8782	26245.0



**Query data
directly from
the AWS console**

Results are instantly
visible within the console

No need to install
and setup an external
JDBC/ODBC client

Redshift Advisor

Launched in July!



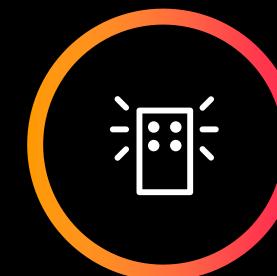
Provides automated recommendations to help optimize database performance and decrease operating costs



>96% of clusters have tailored feedback



Actionable WLM
COPY, storage, and system maintenance advice



Intelligent recommendations for tuning based on continuous workload analysis

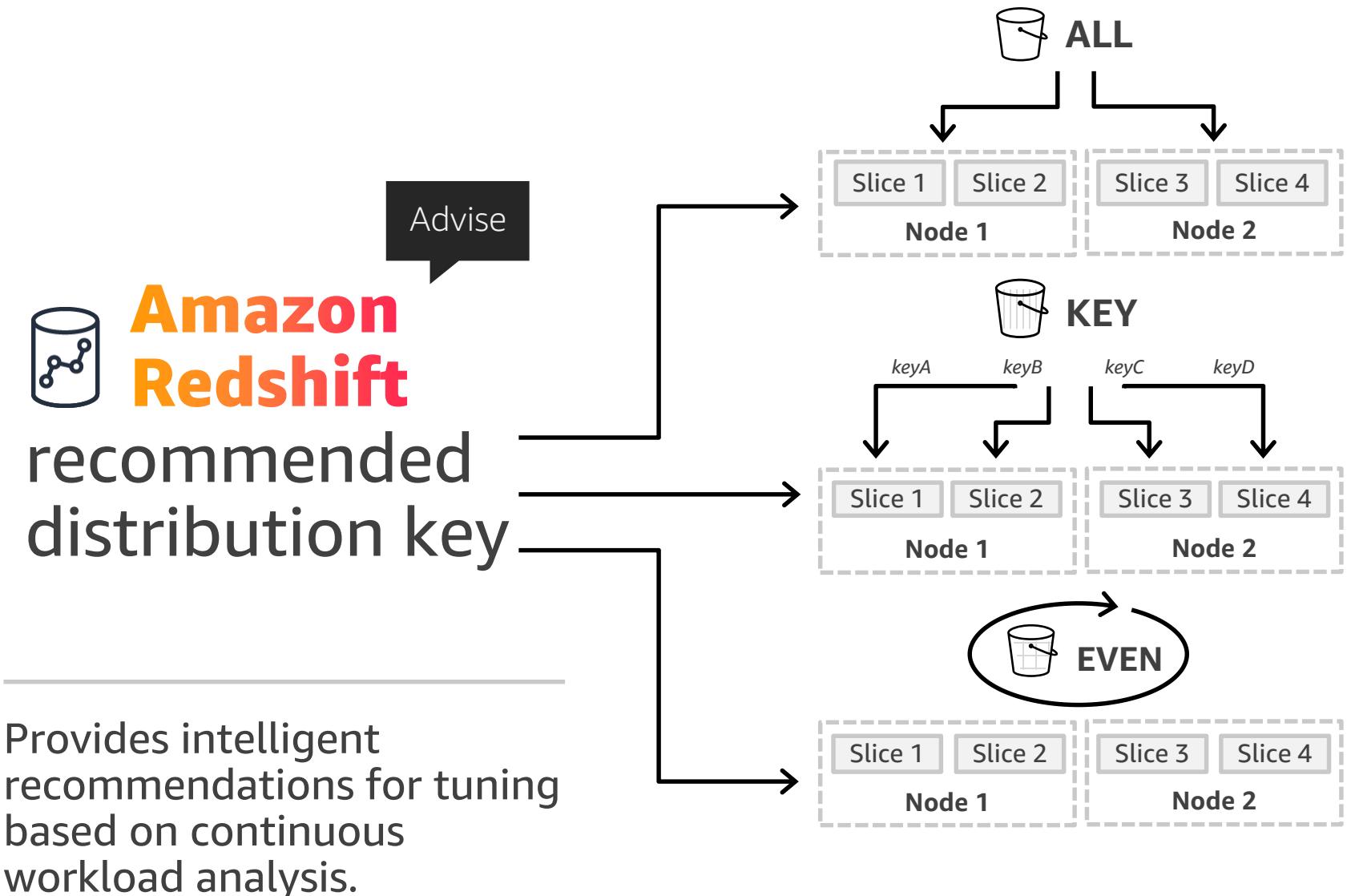
Amazon Redshift intelligent administration



Automates data distribution in tables for improved performance and disk space utilization.

No more messing with distkeys!

Coming Soon!



Amazon Redshift intelligent maintenance

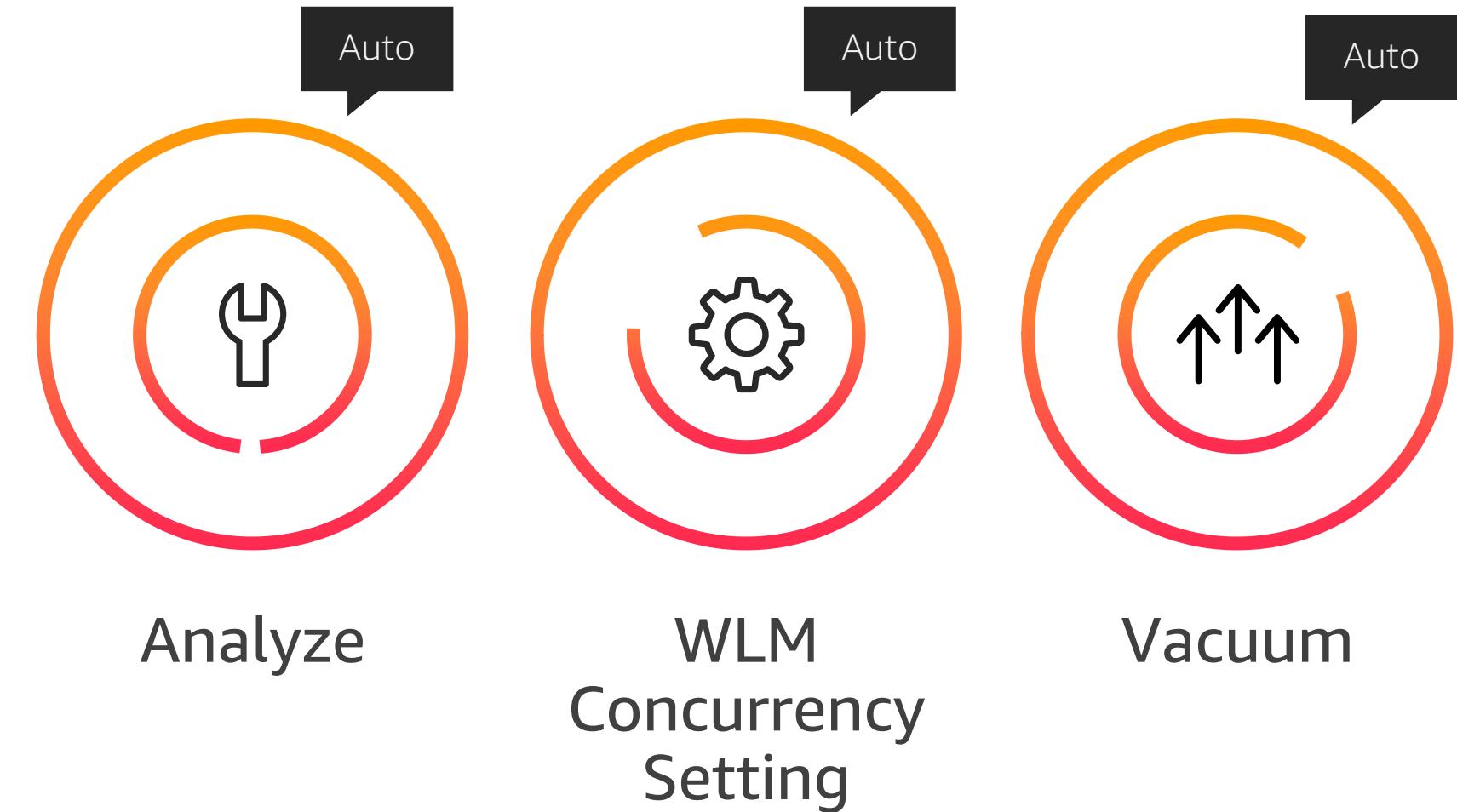
Coming Soon!



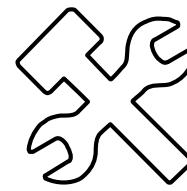
Maintenance processes like vacuum and analyze will automatically run in the background.

Moving towards zero-maintenance.

Redshift will automatically adjust the WLM concurrency setting to deliver optimal throughput.



Run stored procedures in Amazon Redshift



Bring your existing **Stored Procedure** and run in Redshift.

Migrating to Redshift is even easier!

Redshift will support Stored Procedure in PL/pgSQL format, enabling you to bring your existing Stored Procedure to Redshift.

Coming Soon!

Support for stored procedure provides the ability to run code where the data is to efficiently run ETL, data validation, and custom business logic.



DATE data type

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

Query external tables during a resize operation

Specify the root of an S3 bucket as the source for an existing table

Spectrum support for JSON and ION

ALTER TABLE ADD/DROP COLUMN for external tables is now supported via standard JDBC calls

IN-list predicate processing in Spectrum scans

Spectrum queries with aggregations on partition columns

Table property to specify the file compression type for external tables

Map datatypes in Spectrum to contain arrays

Spectrum support for nested data

Renaming external table columns

Retrieving metadata for late-binding views
Support for Enhanced VPC Routing

Push the LENGTH() string function to Spectrum

New Spectrum regions

Arrays of arrays and arrays of maps

*Since re:Invent 2017

Amazon Redshift Spectrum

Extend the data warehouse to exabytes of data in S3 Data Lake

No loading required

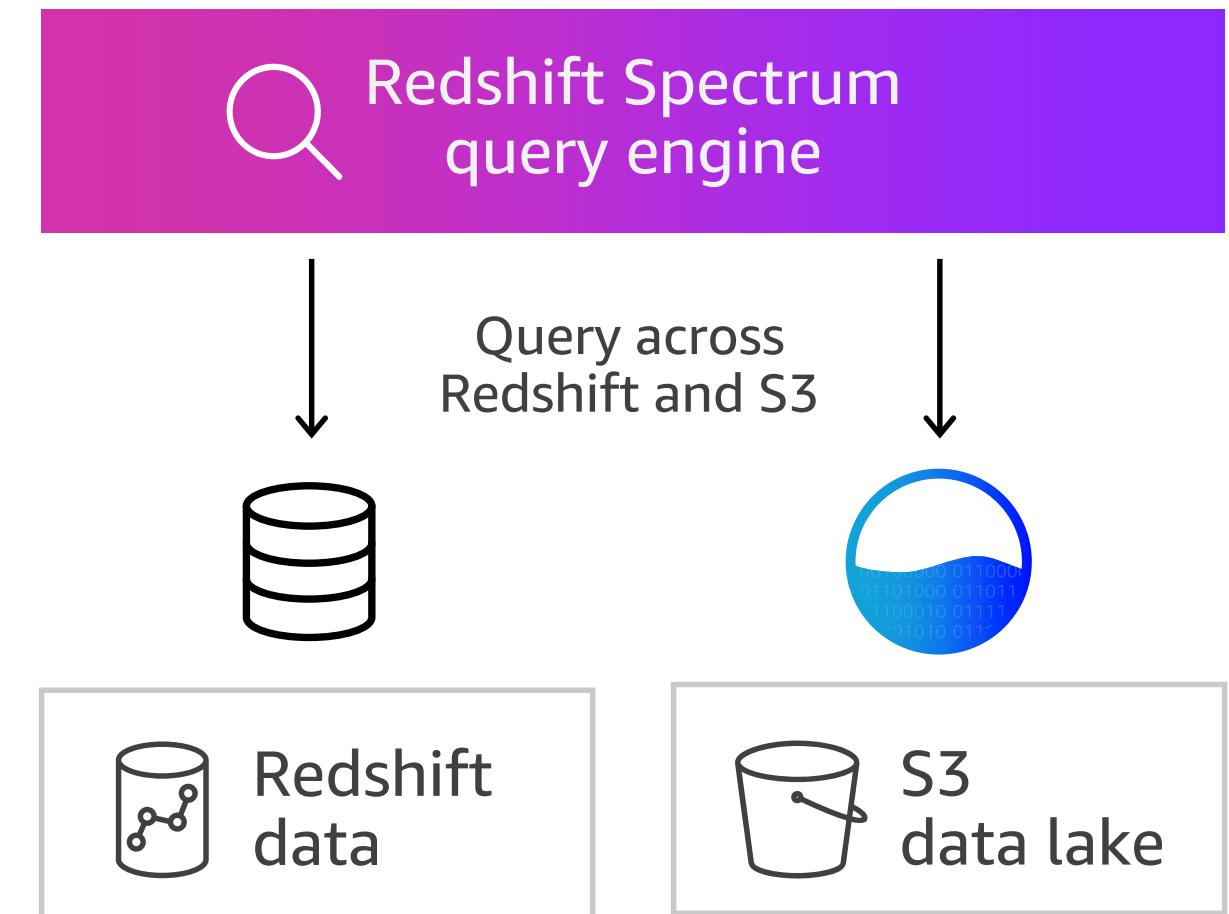
Scale compute and storage separately

Directly query data stored in S3

Parquet, ORC, Avro, Grok, and CSV data formats

- Unload to Parquet
- Spectrum Request Accelerator

Coming Soon!



Amazon Redshift is Scalable

Redshift Spectrum: Exabyte data lake query in under three minutes

Imagine you are the manager at a Seattle book store. An author released her 8th book in a popular series, and you need to figure out how many copies to order.

You have historical book sales data in Amazon S3

Roughly 140 terabytes of customer item order detail records for each day over the past 20 years

190 million files across 15,000 partitions in S3

One partition per day for USA and rest of world

Total data size is over an exabyte

Redshift Spectrum
<3 minutes

Optimization

Compression	-----	5X
Columnar file format	-----	10X
Scanning with 2500 nodes	-----	2,500X
Static partition elimination	-----	2X
Dynamic partition elimination	-----	350X
Amazon Redshift query optimizer	-----	40X

* Query used a 20 node DC1.8XLarge Amazon Redshift cluster

* Not actual sales data—generated for this demo based on data format used by Amazon Retail.

**Encrypt your previously
unencrypted cluster with 1-click**

Default access
privileges

Encrypt unloaded data using S3
server-side encryption with AWS
KMS keys

Federated
authentication with
single sign-on

Cross-region backups for
KMS-encrypted clusters

Improvements to security

IAM roles with COPY
and UNLOAD
commands

Superusers to grant users
access to all rows in
selected system tables

Tag-based
permissions

Enhanced
VPC Routing

SAS integration
enhancements

***Since re:Invent 2015**

Security is **built-in**



End-to-end encryption



Integration with AWS Key Management Service

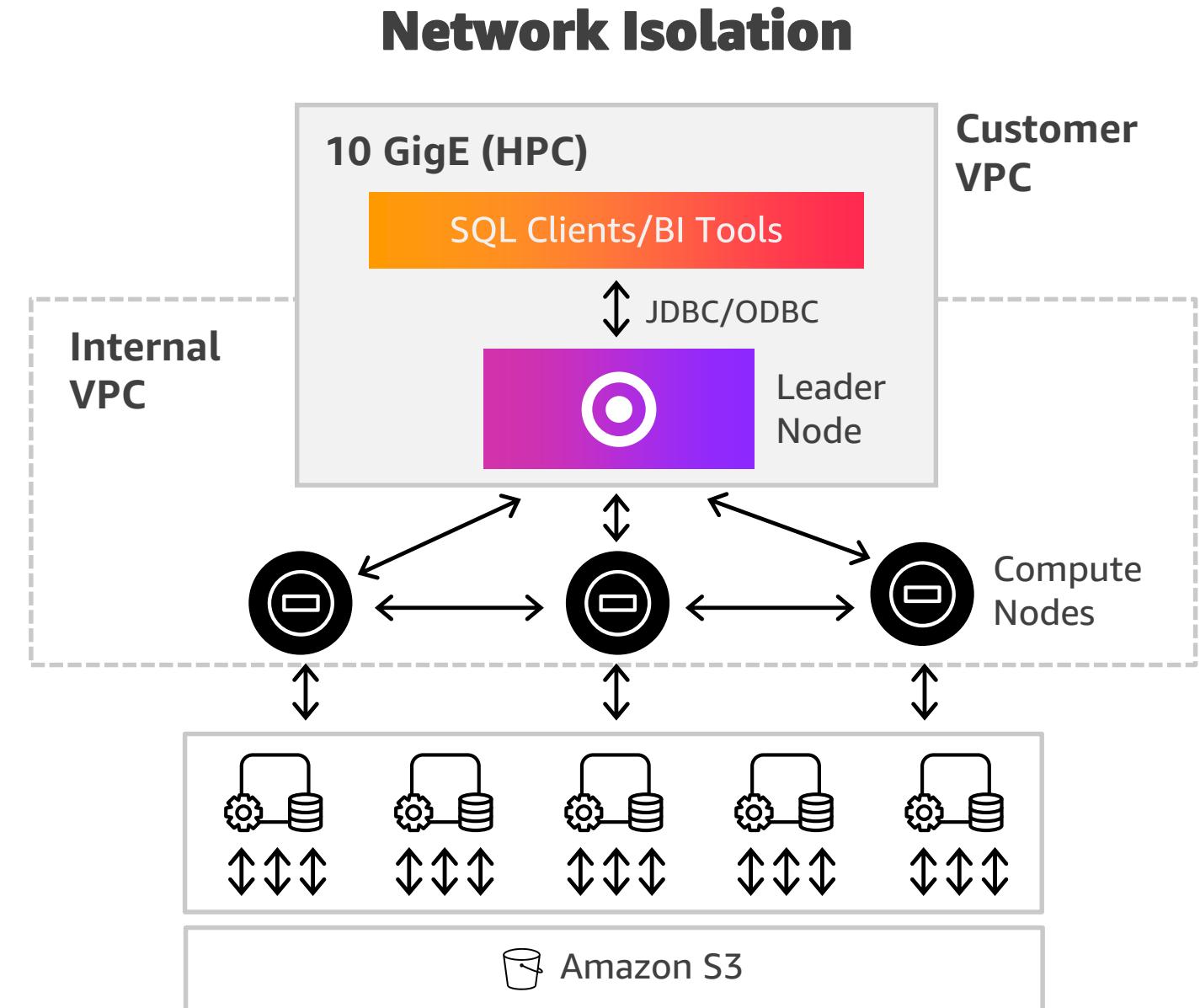
Select compliance certifications*



FedRAMP



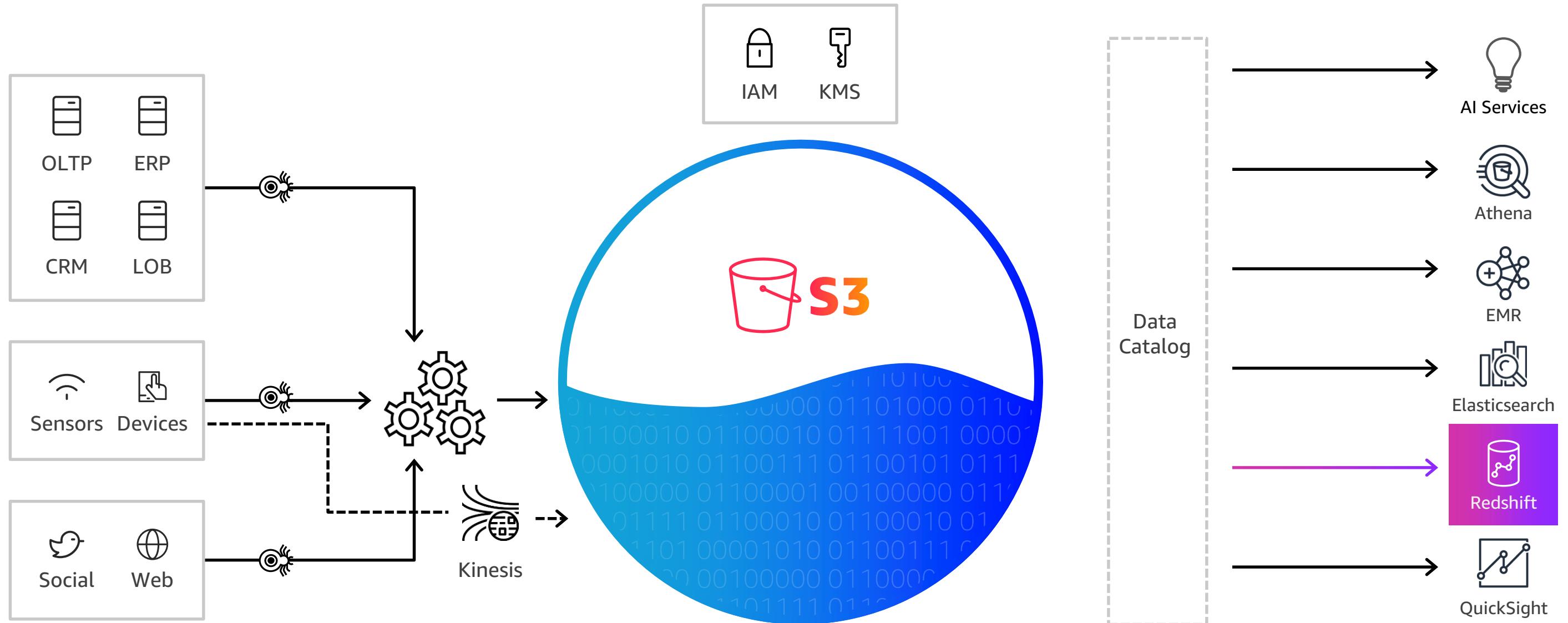
HIPAA
COMPLIANT



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

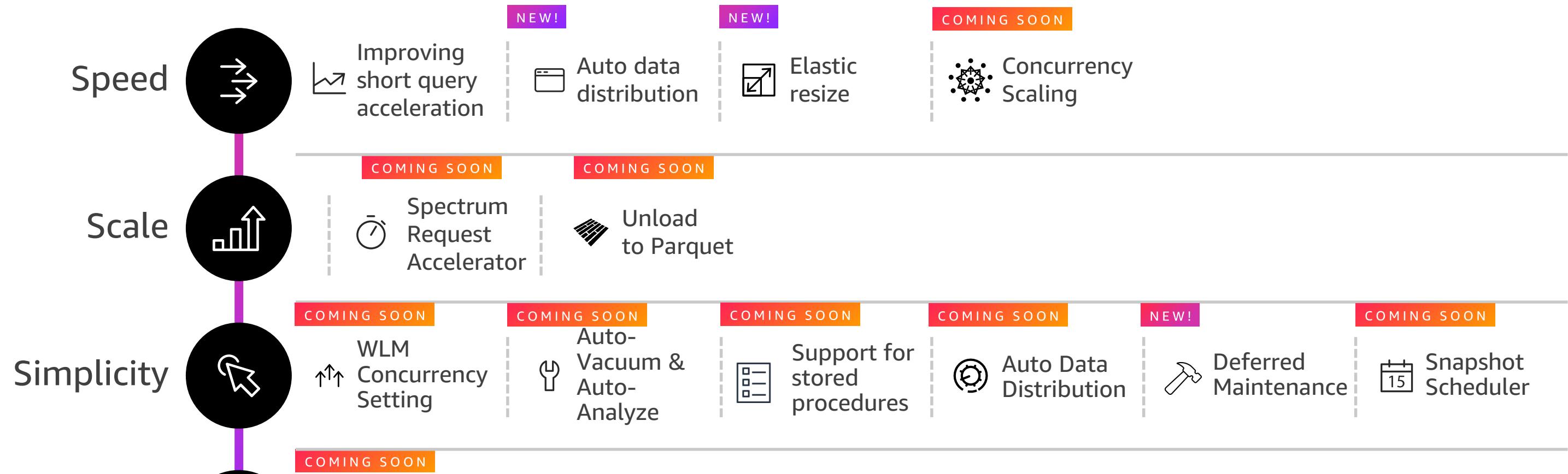
Integration with Amazon Lake Formation

Coming Soon!



Amazon Redshift

New features



More places to learn about Amazon Redshift

Try it out for yourself:
aws.amazon.com/redshift/



© 2018, Amazon Web Services, Inc. or its affiliates. All rights reserved.



**Modern Data
Warehousing**
on AWS ebook



**Blog on
performance
matters:**

Amazon Redshift is
now up to 3.5x
faster for real-
world workloads



Sign up for
**Concurrency
Scaling**



Amazon Redshift
and the **art of
performance
optimization**
in the cloud

by Werner Vogels



Amazon Redshift
**customer
use cases**



**Building
a Proof of
Concept**
for Amazon
Redshift



Thank you!

Vidhya Srinivasan
Vid@amazon.com



Please complete the session
survey in the mobile app.