# Getting Started with Amazon Redshift

Avinash Nidumbur
Sr Business Development Manager
Amazon Redshift
avinid@amazon.com
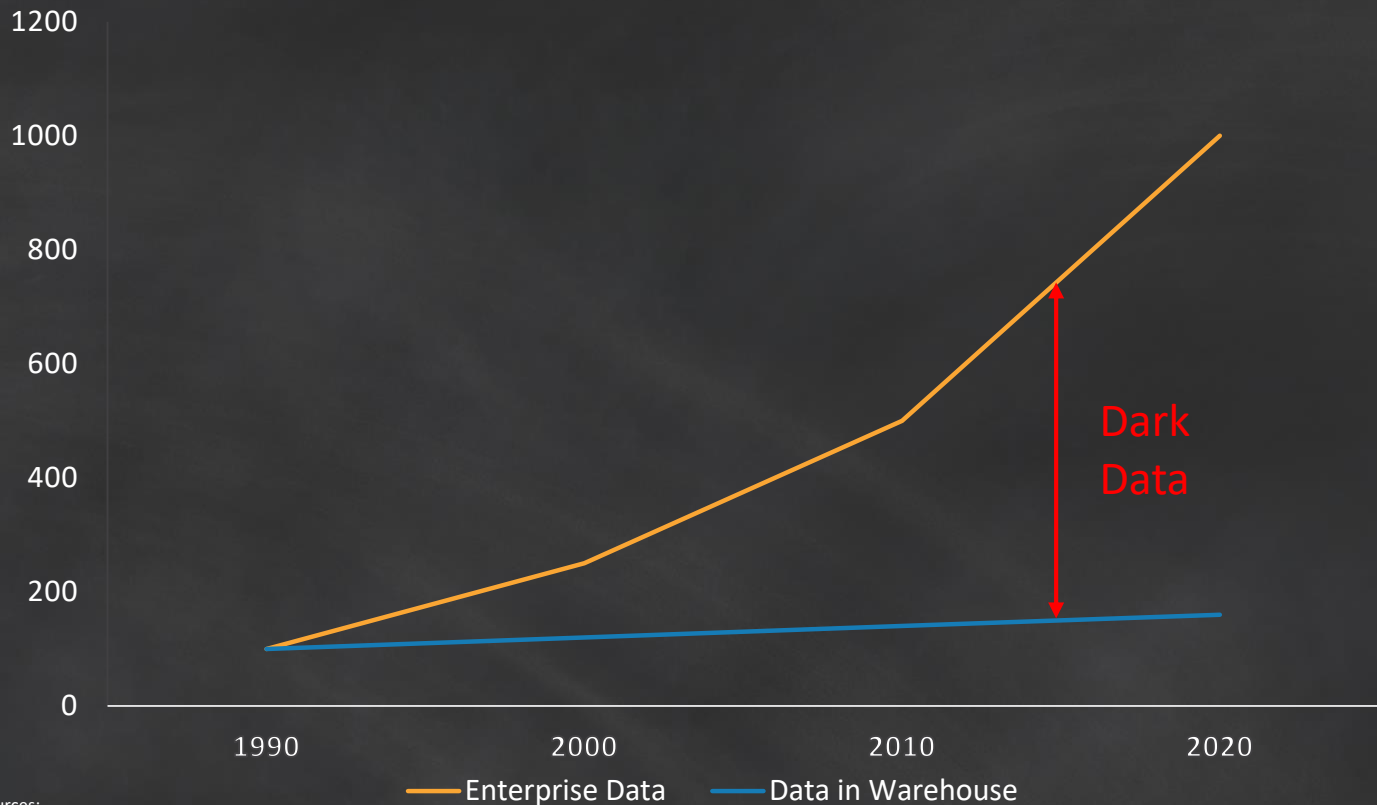
aws | Pop-up Loft

# Agenda

- Introduction
- Benefits
- Use cases
- Getting started
- Q&A

# AWS Big Data Portfolio

## Collect

Amazon Kinesis Firehose

AWS Direct Connect

Amazon Kinesis Streams

Amazon Snowball

## Store

Amazon S3

Amazon Glacier

Amazon Dynamo DB

Amazon RDS, Amazon Aurora

Amazon CloudSearch

Amazon Elasticsearch

## Analyze
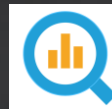
Amazon EMR

Amazon EC2

Amazon Redshift

Amazon Machine Learning

Amazon Kinesis Analytics

Amazon QuickSight

Amazon Athena

AWS Database Migration Service

AWS Data Pipeline

AWS Glue

# Amazon.com clickstream analytics
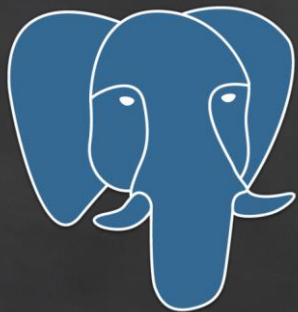
- Web log analysis for Amazon.com
  - PBs workload, 2TB/day@67% YoY
  - Largest table: 400 TB

- Understand customer behavior

- Previous solution
  - Legacy DW (Oracle)—query across 1 week/hr
  - Hadoop—query across 1 month/hr

aws

# Results with Amazon Redshift

- Query 15 months **in 14 min**

- Load 5B rows in **10 min**

- 21B w/ 10B rows: **3 days to 2 hrs**
(Hive → Redshift)

- Load pipeline: **90 hrs to 8 hrs**
(Oracle → Redshift)

- 100 node DS2.8XL clusters

- Easy resizing

- Managed backups and restore

- Failure tolerance and recovery

- 20% time of one DBA

- Increased productivity

**Amazon Redshift**

Relational data warehouse

Massively parallel

Fully managed

HDD and SSD platforms

$1,000/TB/year; starts at $0.25/hour

*a lot faster*
*a lot simpler*
*a lot cheaper*

aws

# Selected Amazon Redshift Customers

# Use Case: Traditional Data Warehousing

**Business Reporting**

**Advanced pipelines and queries**

**Secure and Compliant**

**Bulk Loads and Updates**

**Easy Migration** – Point & Click using AWS Database Migration Service

**Secure & Compliant** – End-to-End Encryption. SOC 1/2/3, PCI-DSS, HIPAA and FedRAMP compliant

**Large Ecosystem** – Variety of cloud and on-premises BI and ETL tools

NTT docomo

Japanese Mobile Phone Provider

SCHOLASTIC

World's Largest Children's Book Publisher

Nasdaq

Powering 100 marketplaces in 50 countries

aws

# Use Case: Log Analysis

**Log & Machine IOT Data**

**Clickstream Events Data**

**Time-Series Data**

**Cheap** – Analyze large volumes of data cost-effectively

**Fast** – Massively Parallel Processing (MPP) and columnar architecture for fast queries and parallel loads

**Near real-time** – Micro-batch loading and Amazon Kinesis Firehose for near-real time analytics

Interactive data analysis and recommendation engine

Ride analytics for pricing and product development

Ad prediction and on-demand analytics

# Use Case: Business Applications

**Multi-Tenant BI Applications**

**Back-end services**

**Analytics as a Service**

**Fully Managed** – Provisioning, backups, upgrades, security, compression all come built-in so you can focus on your business applications

**Ease of Chargeback** – Pay as you go, add clusters as needed. A few big common clusters & data marts

**Service Oriented Architecture** – Integrated with other AWS services. Easy to plug into your pipeline

Infosys®

Infosys Information Platform (IIP)

accenture >

Analytics-as-a-Service

Amplitude

Product and Consumer Analytics

aws

# Amazon Redshift architecture

- **Leader node**
  - Simple SQL endpoint
  - Stores metadata
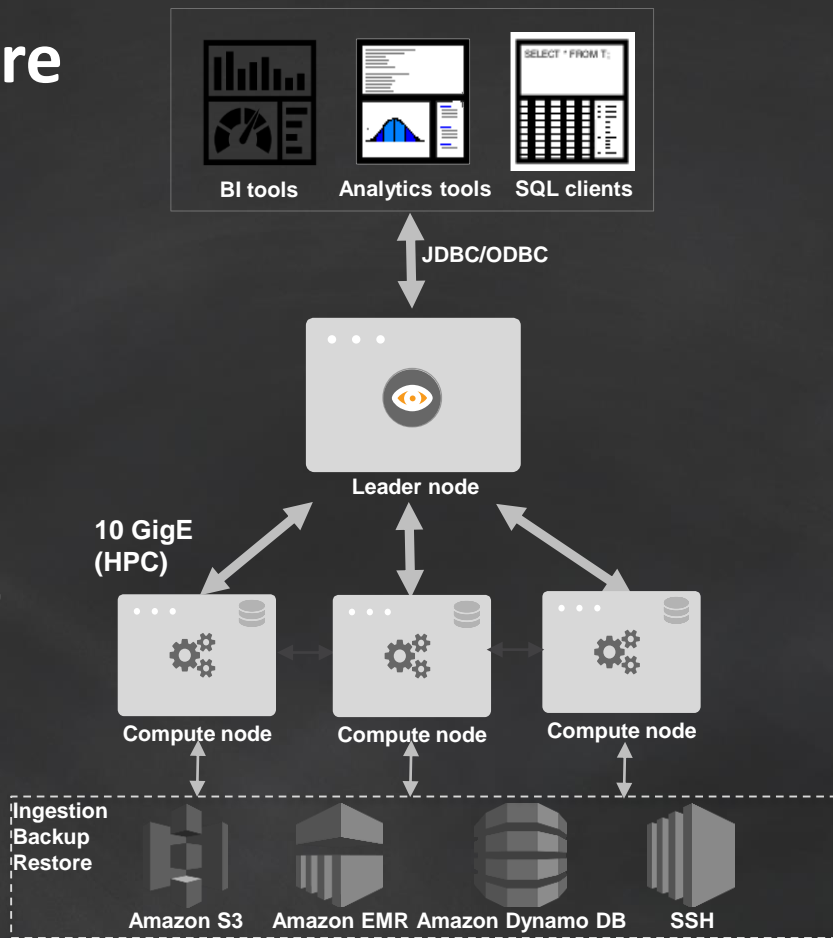  - Optimizes query plan
  - Coordinates query execution

- **Compute nodes**
  - Local columnar storage
  - Parallel/distributed execution of all queries, loads, backups, restores, resizes

- **Start at just $0.25/hour, grow to 2 PB (compressed)**
  - DC1/DC2: SSD; scale from 160 GB to 326 TB
  - DS2: HDD; scale from 2 TB to 2 PB

BI tools  Analytics tools  SQL clients

JDBC/ODBC

Leader node

10 GigE
(HPC)

Compute node    Compute node    Compute node

Ingestion
Backup
Restore

Amazon S3  Amazon EMR  Amazon Dynamo DB  SSH

aws

# Benefit #1: Amazon Redshift is fast

- **Dramatically less I/O**

  – Column storage

  – Data compression

  – Zone maps

  – Direct-attached storage

  – Large data block sizes

```
analyze compression listing;

  Table    |     Column      | Encoding
----------+-----------------+----------
 listing  | listid          | delta
 listing  | sellerid        | delta32k
 listing  | eventid         | delta32k
 listing  | dateid          | bytedict
 listing  | numtickets      | bytedict
 listing  | priceperticket  | delta32k
 listing  | totalprice      | mostly32
 listing  | listtime        | raw
```

| 10 | 10 \| 13 \| 14 \| 26 \|... |
|---|---|
| 324 | ... \| 100 \| 245 \| 324 |
| 375 | 375 \| 393 \| 417... |
| 623 | ... 512 \| 549 \| 623 |
| 637 | 637 \| 712 \| 809 ... |
| 959 | ... \| 834 \| 921 \| 959 |

| 1 | RFK | 900 Columbus | MOROCCO | MOROCCO | AFRICA | 25-989-741-2988 | BUILDING |
| 2 | JFK | 800 Washington | JORDAN | JORDAN | MIDDLE EAST | 23-768-687-3665 | AUTOMOBILE |
| 3 | LBJ | 700 Foxborough | ARGENTINA | ARGENTINA | AMERICA | 11-719-748-3364 | AUTOMOBILE |
| 4 | GWB | 600 Kansas | EGYPT | EGYPT | MIDDLE EAST | 14-128-190-5944 | MACHINERY |

Column 0    Column 1    Column 2

| 1,2,3,4 | RFK,JFK,LBJ,GWB | 900 Columbus,800 Washington, 700 Foxborough,600 Kansas |

# Benefit #2: Amazon Redshift is inexpensive

| DS2 (HDD) | Price per hour for DS2.XL single node | Effective annual price per TB compressed |
|---|---|---|
| On-demand | $ 0.850 | $ 3,725 |
| 1 year reservation | $ 0.500 | $ 2,190 |
| 3 year reservation | $ 0.228 | $   999 |

| DC1 (SSD) | Price per hour for DC1.L single node | Effective annual price per TB compressed |
|---|---|---|
| On-demand | $ 0.250 | $ 13,690 |
| 1 year reservation | $ 0.161 | $  8,795 |
| 3 year reservation | $ 0.100 | $  5,500 |

<u>Pricing is simple</u>
Number of nodes x price/hour
No charge for leader node
No upfront costs
Pay as you go

aws

# Benefit #3 : Amazon Redshift is easy to use

Provisioning in minutes

Automatic patching

SQL - Data loading

Backups are built-in

Security is built-in

Compression is built-in

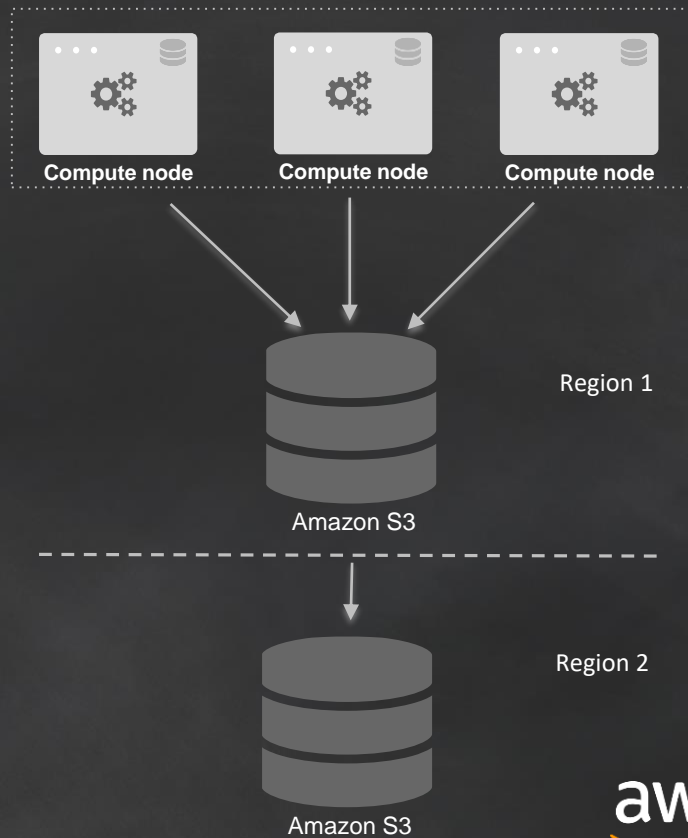# Benefit #4: Amazon Redshift is fully managed

## Continuous/incremental backups

Multiple copies within cluster
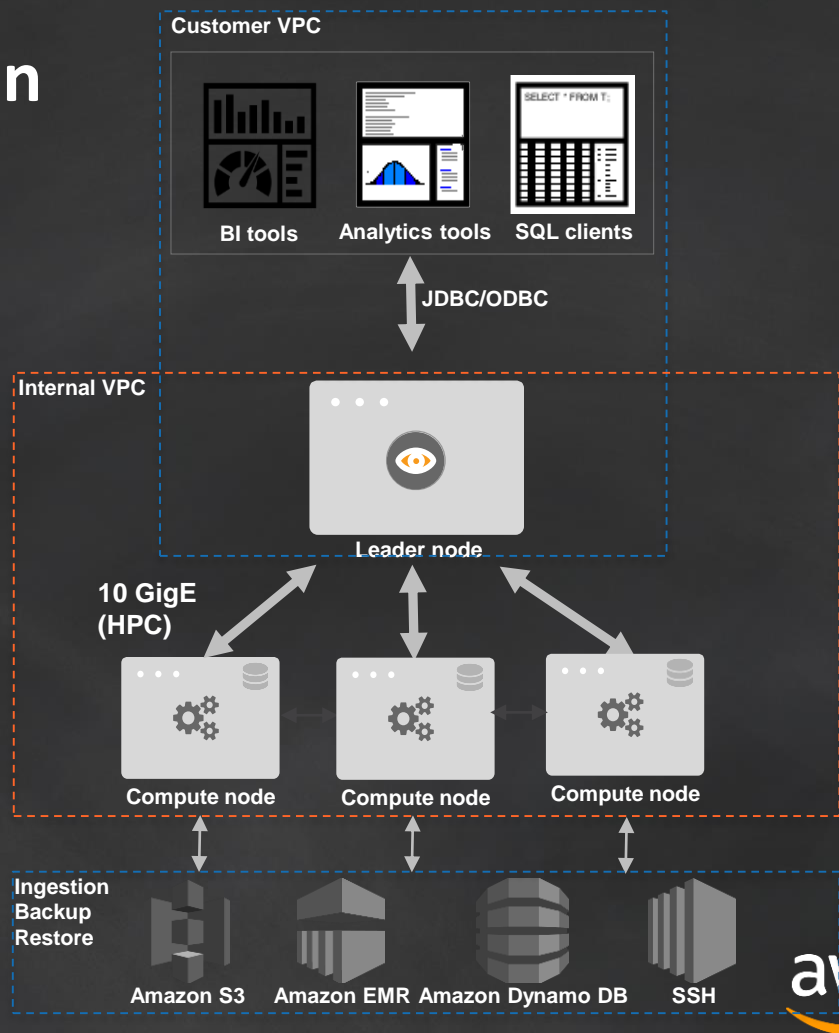
Continuous and incremental backups to Amazon S3

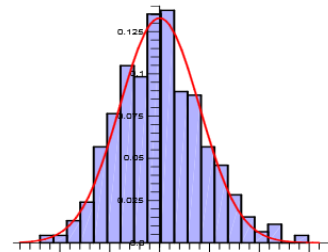Continuous and incremental backups across regions

Streaming restore



Compute node   Compute node   Compute node

Region 1

Amazon S3

Region 2

Amazon S3

aws

# Benefit #5: Security is built-in

- Load encrypted from S3

- SSL to secure data in transit
  - ECDHE perfect forward secrecy

- Amazon VPC for network isolation

- Encryption to secure data at rest
  - All blocks on disks and in S3 encrypted
  - Block key, cluster key, master key (AES-256)
  - On-premises HSM & AWS CloudHSM support

- Audit logging and AWS CloudTrail integration

- SOC 1/2/3, PCI-DSS, FedRAMP, BAA

Customer VPC

BI tools    Analytics tools    SQL clients

JDBC/ODBC

Internal VPC

Leader node

10 GigE
(HPC)

Compute node    Compute node    Compute node

Ingestion
Backup
Restore

Amazon S3    Amazon EMR    Amazon Dynamo DB    SSH

aws

# Benefit #6: Amazon Redshift is powerful

- Approximate functions

- User defined functions

- Machine learning

- Data science



*HyperLogLog: analysis of a near-optimal cardinality algorithm*

# Benefit #7: Amazon Redshift has a large ecosystem

# Amazon Redshift Spectrum

# Amazon Redshift Spectrum

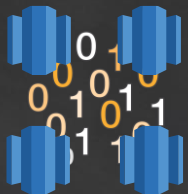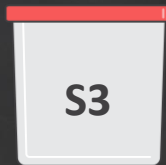**Run SQL queries directly against data in S3 using thousands of nodes**

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
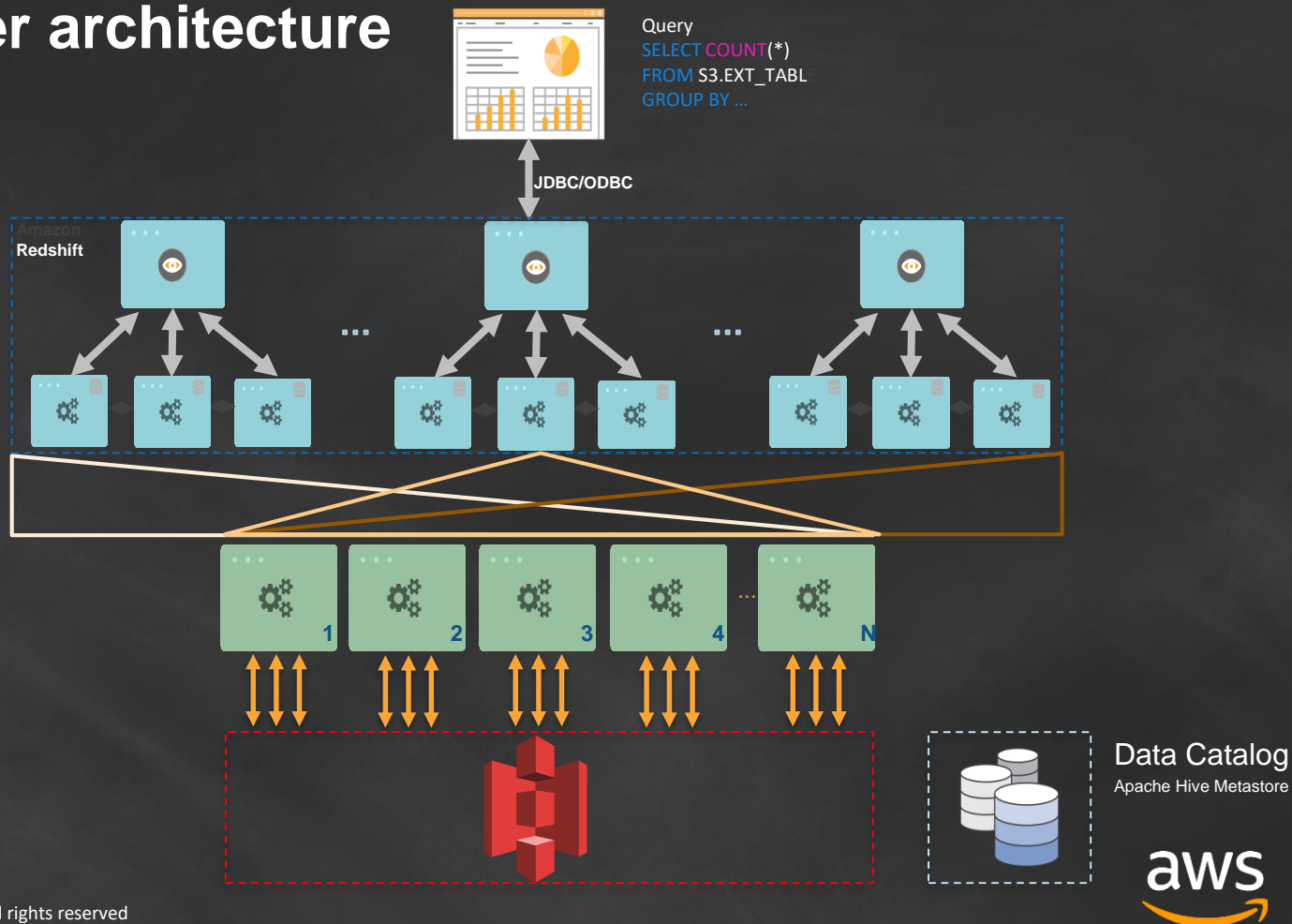
S3

I ♥ SQL

aws

# Extend Redshift Queries To Amazon S3

**Amazon Redshift**

**Amazon S3**

Redshift Query Engine

Redshift Data

Data Lake

Query S3 directly or join data across Redshift and S3

Scale Redshift compute and storage separately

Support for CSV, Parquet, ORC, Grok, Avro and more formats

aws

# Life of a query

**1** Query
SELECT COUNT(*)
FROM S3.EXT_TABLE
GROUP BY…

**JDBC/ODBC**

**Amazon Redshift**

**9** Result is sent back to client

**2** Query is optimized and compiled at the leader node. Determine what gets run locally and what goes to Amazon Redshift Spectrum

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

**3** Query plan is sent to all compute nodes

**4** Compute nodes obtain partition info from Data Catalog; dynamically prune partitions

**7** Amazon Redshift Spectrum projects, filters, joins and aggregates

**5** Each compute node issues multiple requests to the Amazon Redshift Spectrum layer

**6** Amazon Redshift Spectrum nodes scan your S3 data

## Amazon S3
Exabyte-scale object storage

## Data Catalog
Glue / Apache Hive Metastore

# Use cases

- Extend data warehousing to data lake in Amazon S3

- Cost reduction for less frequently accessed data

- Better query performance and minimal time to insight at large scale

- Query many open formats and large data sets that can't be loaded into cluster

- Improved concurrency with multiple Redshift clusters querying common data

- Elastically scale compute resources separately from the storage layer in S3

- Use familiar SQL to cleanse, transform and load data from S3 to Redshift

aws

# Multiple cluster architecture



Query
SELECT COUNT(*)
FROM S3.EXT_TABL
GROUP BY ...

- Use multiple Redshift clusters to query same copy of data in s3

Data Catalog
Apache Hive Metastore

# Simplify and accelerate your ETL pipelines

| "Dirty" Logs | Clean Logs | CSV | COPY → | Staging Table<br>Amazon Redshift | Final Table<br>Amazon Redshift |

| "Dirty" Logs | `CREATE TABLE AS`<br>`SELECT C.. FROM S3.xxx WHERE …` → | Amazon Redshift |

aws

# Customer Use Cases

# Nasdaq: powering 100 marketplaces in 50 countries

Orders, quotes, trade executions, market "tick" data from 7 exchanges

7 billion rows/day

Analyze market share, client activity, surveillance, billing, and so on

Microsoft SQL Server on-premises

Expensive legacy DW ($1.16 M/yr.)

Limited capacity (1 yr. of data online)

Needed lower TCO

Must satisfy multiple security and regulatory requirements

Similar performance

# Nasdaq: powering 100 marketplaces in 50 countries



23 node DS2.8XL cluster

828 vCPUs, 5 TB RAM

368 TB compressed

2.7 T rows, 900 B derived

8 tables with 100 B rows

7 man-month migration

¼ the cost, 2x storage, room to grow

Faster performance, very secure

# Customers love Amazon Redshift Spectrum

**Time Inc.**

"Redshift Spectrum enables us to directly operate on our data in its native format in Amazon S3 with no preprocessing or transformation."

**NTT docomo**

"Redshift Spectrum will let us expand the universe of the data we analyze to 100s of petabytes over time. This is truly a game changer, and we can think of no other system in the world that can get us there."

**edmunds**

"Redshift Spectrum's fast performance across massive data sets is unprecedented."

**REDFIN**

"Our data science team using Amazon EMR can now collaborate with our marketing and product teams using Redshift Spectrum to analyze the same Amazon S3 data sets."

**yelp**

"Multiple teams can now query the same Amazon S3 data sets using both Redshift and EMR."

**RECRUIT**
Recruit Technologies Co.,Ltd.

"Redshift Spectrum will help us scale yet further while also lowering our costs."

**aws**

# Getting Started

# Provisioning

# Enter cluster details

# Select node configuration

# Select security settings and provision

# Point-and-click resize

# Data Modeling

# Zone maps

- Single column

- Compound

- Interleaved

# Single Column

- Table is sorted by 1 column
[ SORTKEY ( date ) ]

| Date | Region | Country |
|------|--------|---------|
| 2-JUN-2015 | Oceania | New Zealand |
| 2-JUN-2015 | Asia | Singapore |
| 2-JUN-2015 | Africa | Zaire |
| 2-JUN-2015 | Asia | Hong Kong |
| 3-JUN-2015 | Europe | Germany |
| 3-JUN-2015 | Asia | Korea |

- Best for:
  - Queries that use 1st column (i.e. *date)* as primary filter
  - Can speed up joins and group bys
  - Quickest to VACUUM

aws

# Compound

- Table is sorted by 1st column , then 2nd column etc.

`[ SORTKEY COMPOUND ( date, region, country) ]`

| Date | Region | Country |
|---|---|---|
| 2-JUN-2015 | Africa | Zaire |
| 2-JUN-2015 | Asia | Korea |
| 2-JUN-2015 | Asia | Singapore |
| 2-JUN-2015 | Europe | Germany |
| 3-JUN-2015 | Asia | Hong Kong |
| 3-JUN-2015 | Asia | Korea |

- Best for:
  - Queries that use 1st column as primary filter, then other cols
  - Can speed up joins and group bys
  - Slower to VACUUM

aws

- EVEN

- KEY

- ALL

Distribution

| ID | Gender | Name |
|----|--------|------|
| 101 | M | John Smith |
| 139 | M | Peter Black |
| 446 | M | Pat Partridge |
| 164 | M | Brian Snail |
| 209 | M | James White |

| ID | Gender | Name |
|----|--------|------|
| 101 | M | John Smith |
| 292 | F | Jane Jones |
| 139 | M | Peter Black |
| 446 | M | Pat Partridge |
| 658 | F | Sarah Cyan |
| 164 | M | Brian Snail |
| 209 | M | James White |
| 306 | F | Lisa Green |

KEY

1

2

3

4

| ID | Gender | Name |
|----|--------|------|
| 292 | F | Jane Jones |
| 658 | F | Sarah Cyan |
| 306 | F | Lisa Green |

# DISTSTYLE KEY

- EVEN
  - Tables with no joins or group by

- KEY
  - Large Fact tables
  - Large dimension tables

- ALL
  - Medium dimension tables (1K – 2M)
  - Small dimension tables

aws

# Loading Data

# Data loading options

Flat files

Amazon S3

Amazon
Redshift

AWS

Corporate Data center

AWS Cloud

# Data loading options



Source DBs

ETL

Corporate Data center

INFORMATICA

bryte Systems

ATTUNITY CloudBeam

hapyrus

snapLogic

talend*
integration at any scale

AWS

Amazon Redshift

AWS Cloud

# Data loading options

# ETL data into your data warehouse

# Instantly query your data lake on Amazon S3



**AMAZON S3** → **AWS GLUE CRAWLERS** → **AWS GLUE DATA CATALOG** → **AMAZON ATHENA** / **AMAZON EMR** / **AMAZON REDSHIFT SPECTRUM** → **BI TOOLS**

1. Crawlers scan your data sets and populate the Glue Data Catalog

2. The Glue Data Catalog serves as a central metadata repository

3. Once catalogued in Glue, your data is immediately available for analytics

# Redshift Spectrum: Defining External Schema and External Tables

1. Define an external schema in Amazon Redshift using the AWS Glue Data Catalog or your own Apache Hive Metastore

   ```
   CREATE EXTERNAL SCHEMA <schema_name>
   FROM { [ DATA CATALOG ] | HIVE METASTORE }
   DATABASE 'database_name'
   IAM_ROLE 'iam-role-arn'
   ```

2. Register external tables using Amazon Athena, your Hive Metastore client, or from Amazon Redshift

   ```
   CREATE EXTERNAL TABLE <schema_name>.<table_name>
   [PARTITIONED BY <column_name, data_type, …>]
   STORED AS file_format
   LOCATION s3_location
   [TABLE PROPERTIES property_name=property_value, …];
   ```

3. Query external tables

   ```
   SELECT … FROM <schema_name>.<table_name> …
   ```

aws

# Querying

# Amazon Redshift works with your existing BI tools



**Amazon Redshift**
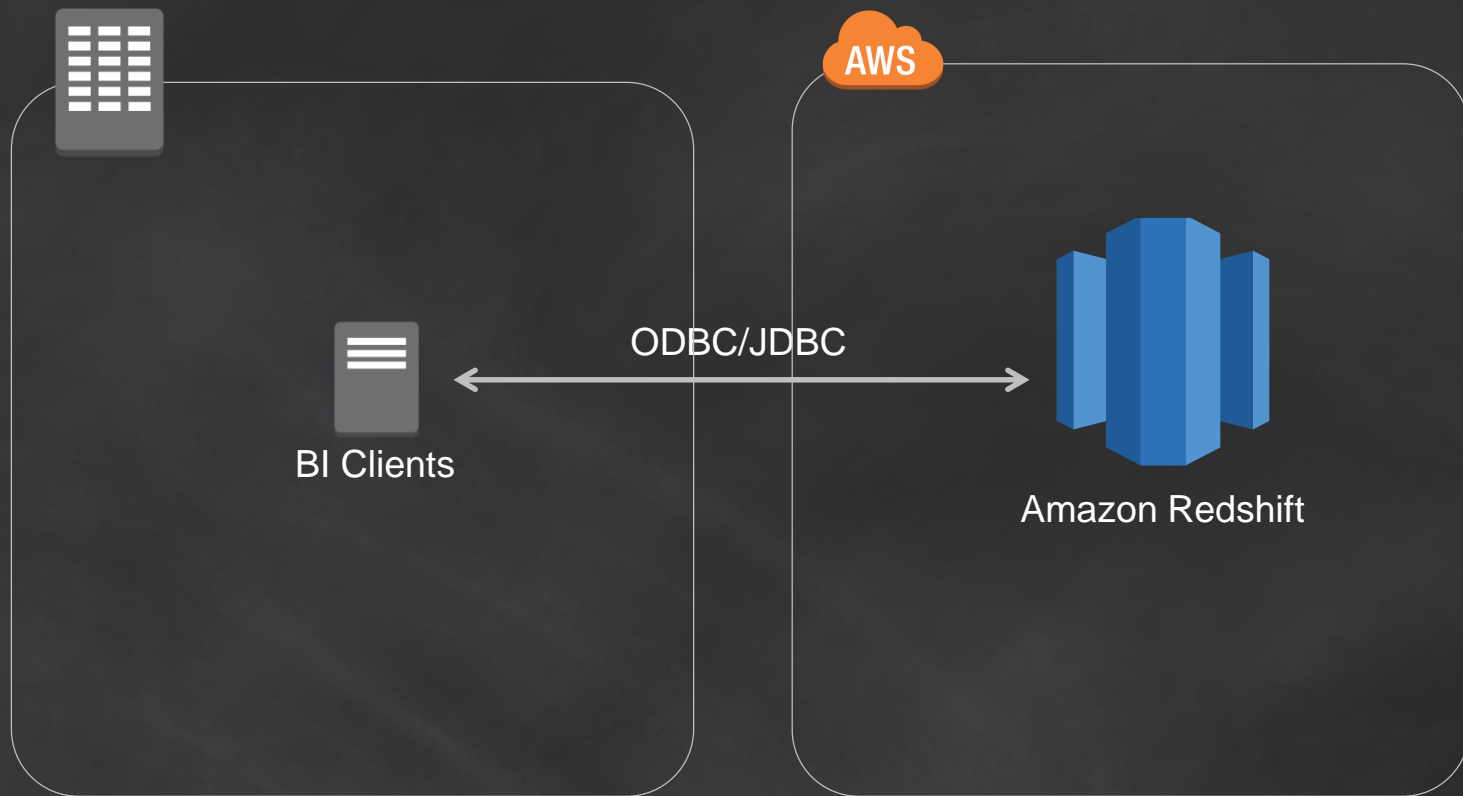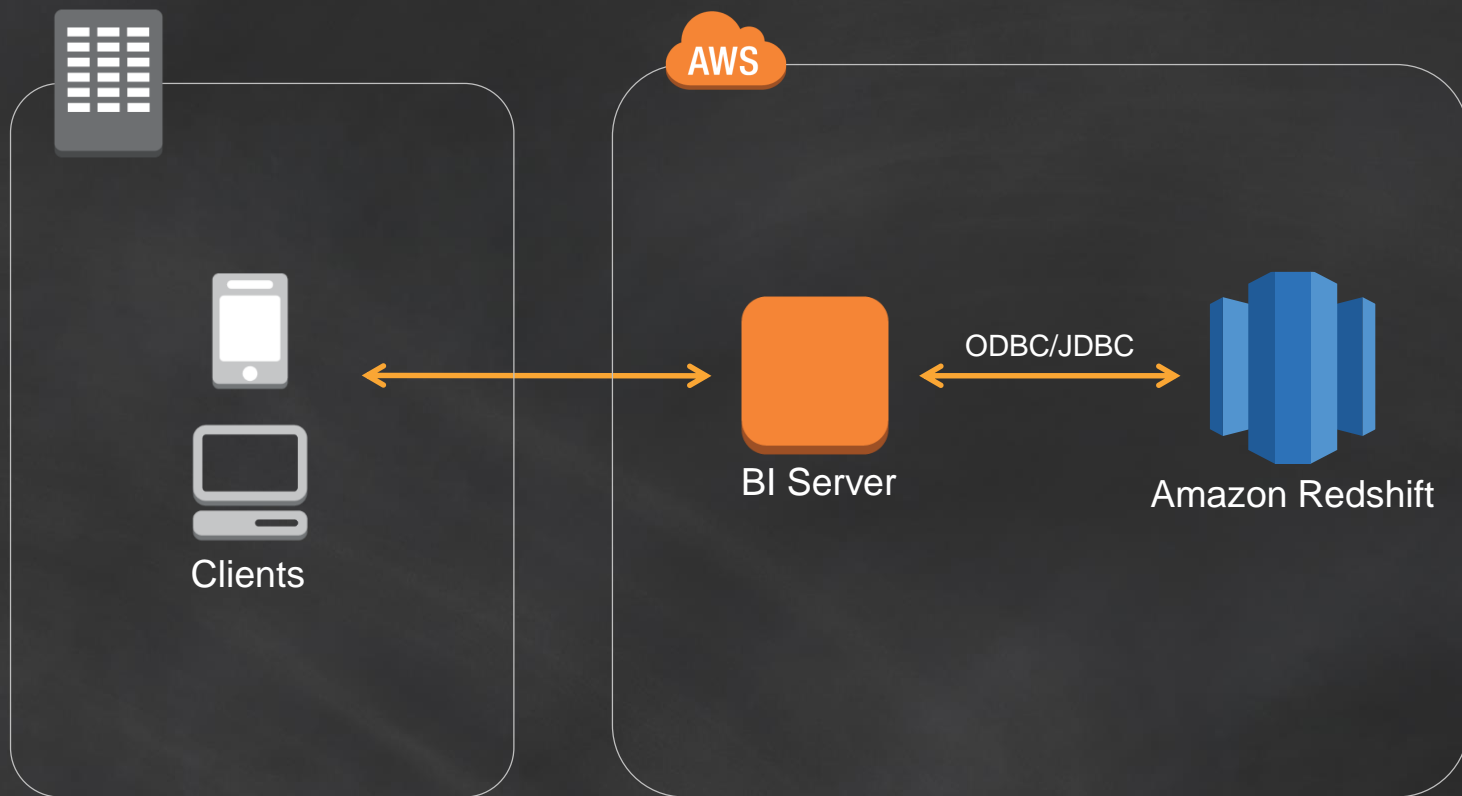
**JDBC/ODBC**

ODBC/JDBC

BI Clients

AWS

Amazon Redshift

aws

AWS

ODBC/JDBC

Clients

BI Server

Amazon Redshift

aws

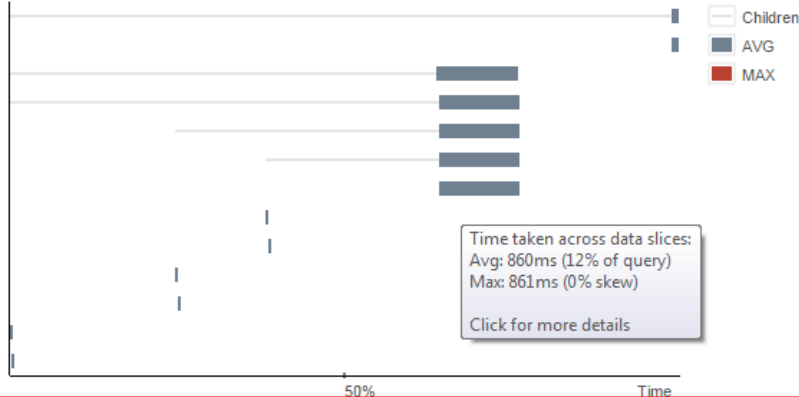# View explain plans

## Query Execution Details

**Plan** | Actual

```
XN Hash Join DS_BCAST_INNER  (cost=52602972108.91..68656434844.99 rows=23608 width=71)
  ->  XN Seq Scan on supplier s  (cost=0.00..10000.00 rows=1000000 width=22)
  ->  XN Hash  (cost=52602972049.89..52602972049.89 rows=23608 width=57)
        ->  XN Hash Join DS_BCAST_INNER  (cost=1752.83..52602972049.89 rows=23608 width=57)
              ->  XN Hash Join DS_BCAST_INNER  (cost=1718.72..52499650202.87 rows=70082 width=61)
                    ->  XN Hash  (cost=31.95..31.95 rows=861 width=4)
                          ->  XN Hash Join DS_BCAST_INNER  (cost=624.99..14000065191.80 rows=714354 wid
                          ->  XN Hash  (cost=874.99..874.99 rows=87499 width=33)
                                ->  XN Seq Scan on lineorder lo  (cost=0.00..149965.90 rows=14996590 wi
                                ->  XN Hash  (cost=499.99..499.99 rows=49999 width=20)
                                      ->  XN Seq Scan on part p  (cost=0.00..499.99 rows=49999 width=20
```

Plan | **Actual**

| Expand Level | Collapse Level | Expand All |

- Hash Join (Broadcasted Inner Table)
  - Sequential Scan on supplier s
- Hash
  - Hash Join (Broadcasted Inner Table)
    - Hash Join (Broadcasted Inner Table
      - Hash Join (Broadcasted Inner Tab
        - Sequential Scan on lineorder lo
        - Hash
          - Sequential Scan on part p
      - Hash
        - Sequential Scan on customer c
  - Hash
    - Sequential Scan on dwdate d

Legend: Children / AVG / MAX

Time taken across data slices:
Avg: 860ms (12% of query)
Max: 861ms (0% skew)

Click for more details

50% ............ Time

```sql
select
    lo_orderkey,
    p_name,
    c_name,
    s_address,
    lo_quantity
from
    lineorder lo,
    part p,
    supplier s,
    customer c,
    dwdate d
where
    lo_custkey = c_custkey
    and     lo_partkey = p_partkey
    and     lo_suppkey = s_suppkey
    and     lo_orderdate = d_datekey
    and     d_sellingseason = 'Summer'
```
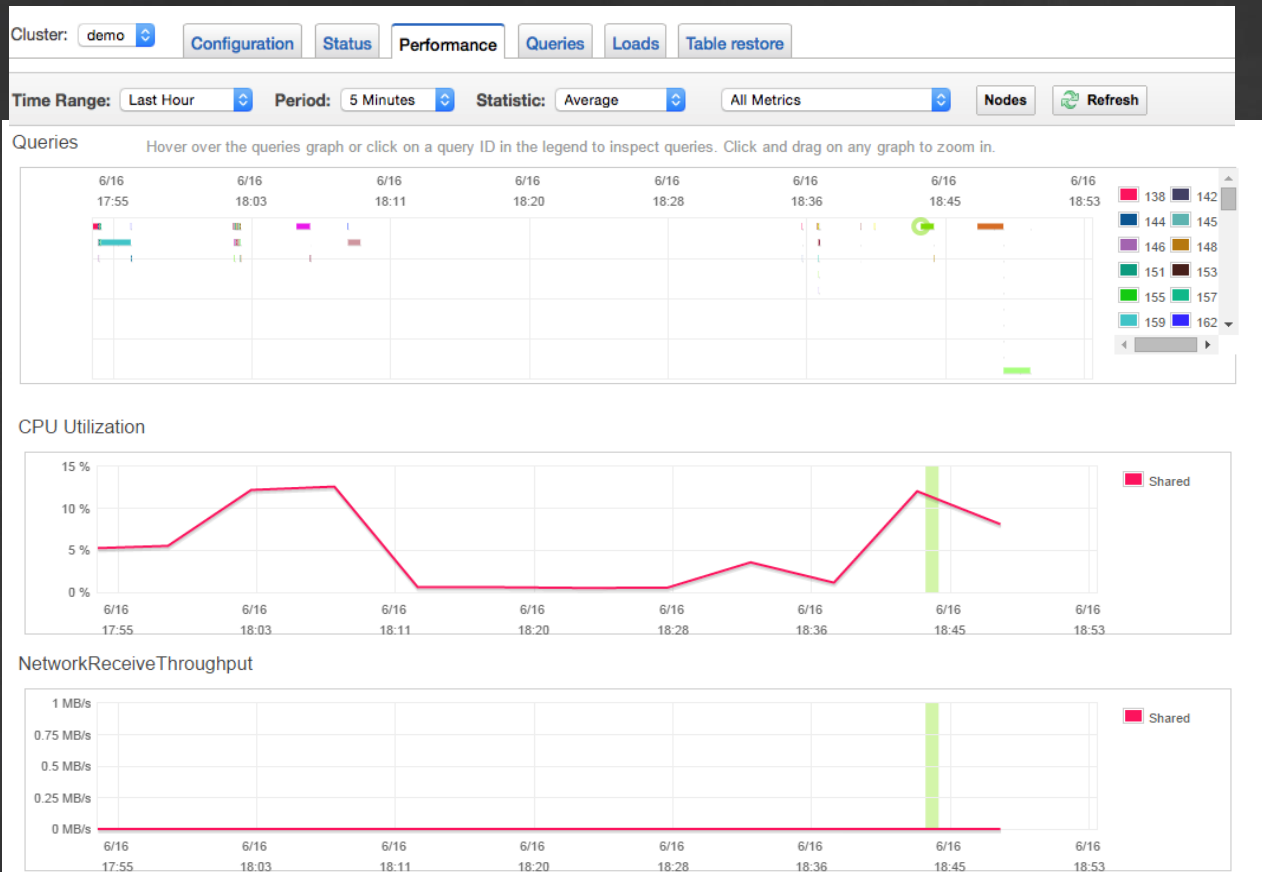
aws

# Monitor query performance

# Resources

- **Detail Pages**
  - http://aws.amazon.com/redshift
  - https://aws.amazon.com/redshift/spectrum
  - https://aws.amazon.com/marketplace/redshift/
  - https://aws.amazon.com/redshift/developer-resources/
  - Amazon Redshift Utilities - GitHub

- **Best Practices**
  - http://docs.aws.amazon.com/redshift/latest/dg/c_loading-data-best-practices.html
  - http://docs.aws.amazon.com/redshift/latest/dg/c_designing-tables-best-practices.html
  - http://docs.aws.amazon.com/redshift/latest/dg/c-optimizing-query-performance.html

aws

**Everything and Anything Startups
Need to Get Started on AWS**

aws.amazon.com/activate