AWS

# SUMMIT

# Building Your Data Lake on AWS

Ben Snively, Senior Solutions Architect, AWS

August 14, 2017

amazon
web services

# What is a Data Lake?

An **architectural approach** that allows you to store massive amounts of "**raw**" data into a central location

It's **readily available** to be **categorized, processed, analyzed, and consumed** by **diverse groups**

# Why use a Data Lake?

**Leverage all data within your organization**

Customer centricity
Business agility
Better predictions
Competitive advantage

**Leads to…**

# Legacy data architectures exist as isolated data silos



Hadoop cluster

Data warehouse appliance

SQL database

# Navigating the Data Lake…

Data Lake is a new and increasingly popular architecture to store and analyze massive volumes and heterogenous types of data in a centralized repository
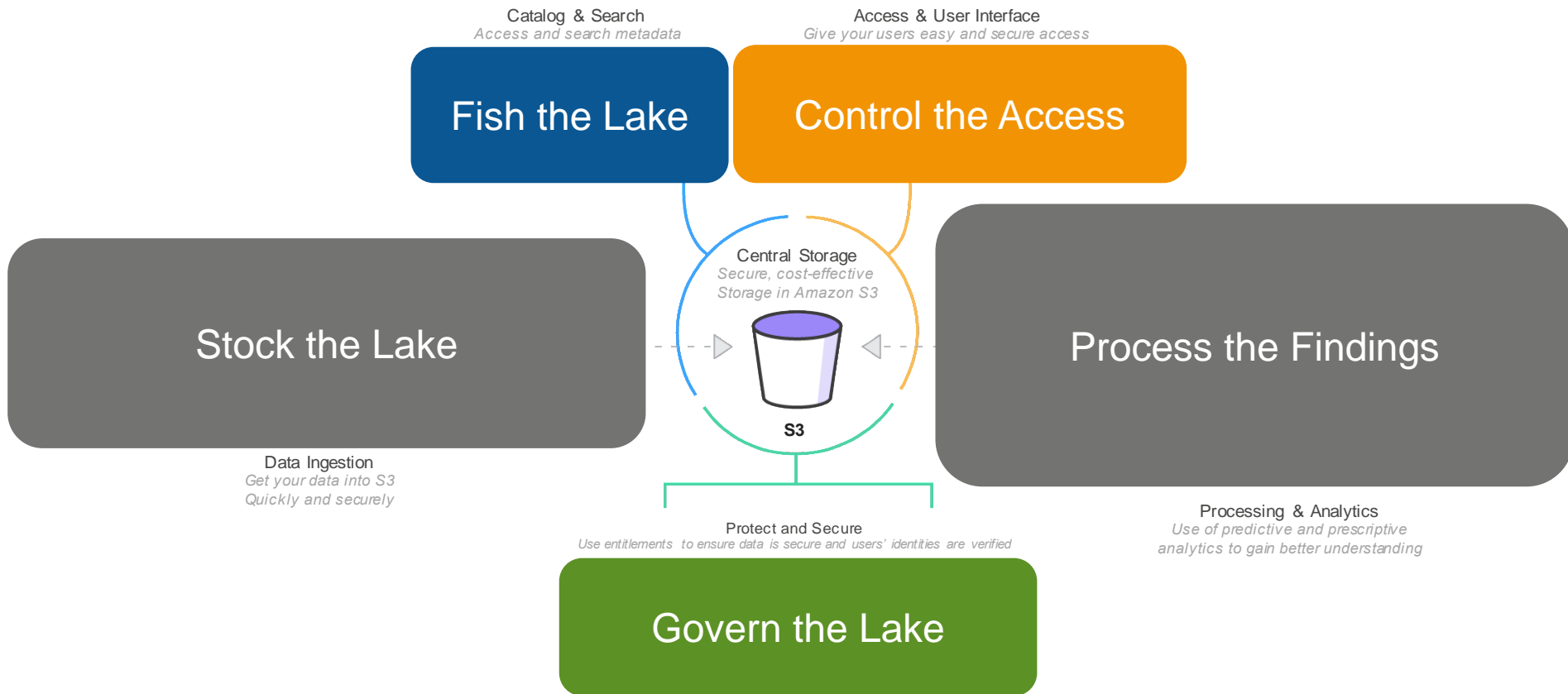
# Building a Data Lake on AWS

# Why AWS?

*Implementing a Data Lake architecture requires a broad set of tools and technologies to serve an increasingly diverse set of applications and use cases.*
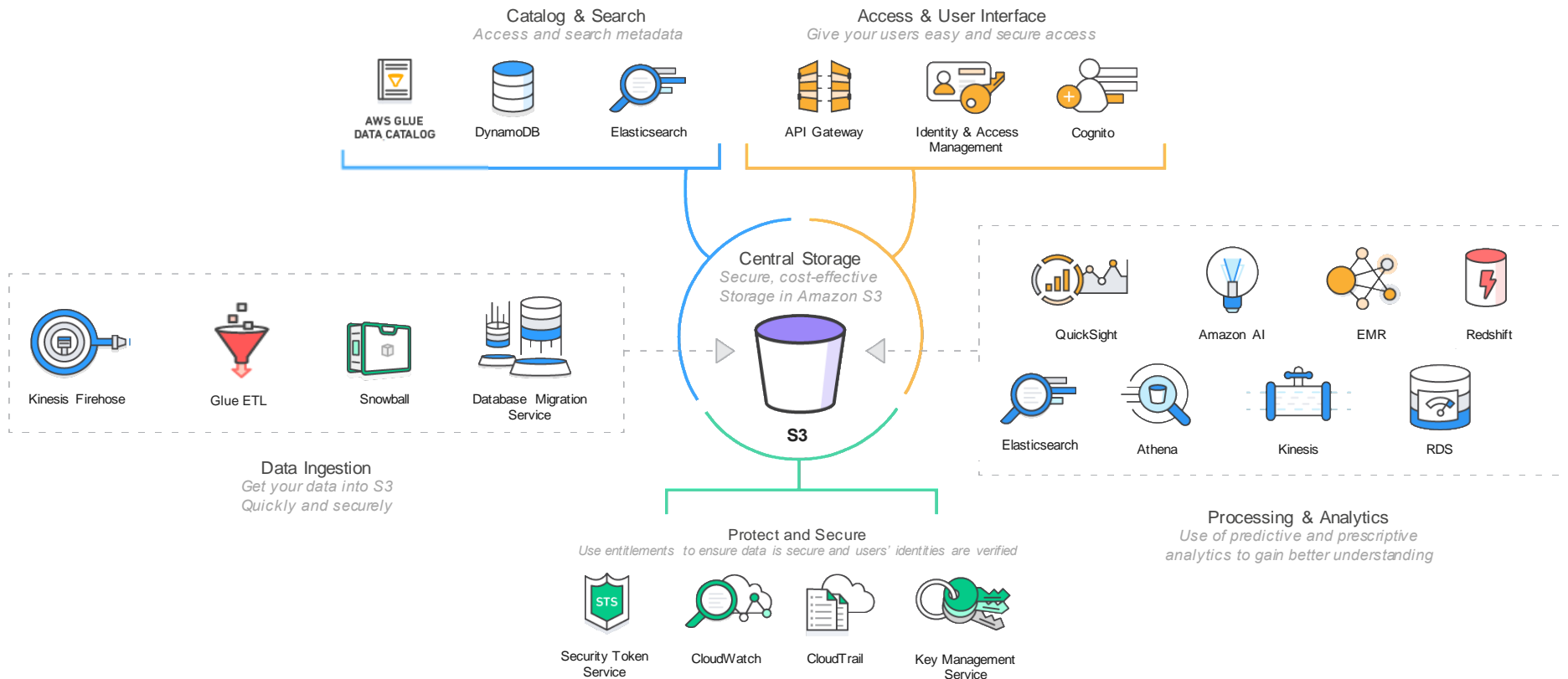
Boils down to:

Picking the right tool for the right job…on a consumption-based model…

# Data Lake reference architecture

# Data Lake reference architecture
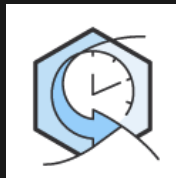
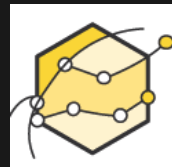# S3 – Center of the Data Lake

# Why Amazon S3 for Data Lake?

## Durable
Designed for 11 9s
of durability

## Available
Designed for
**99.99**% availability

## High performance
- Multiple Upload
- Range GET

## Easy to use
- Simple REST API
- AWS SDKs
- Read-after-create consistency
- Event notification
- Lifecycle policies

## Scalable
- Store as much as you need
- Scale storage and compute independently
- No minimum usage commitments

## Integrated
- Amazon EMR
- Amazon Redshift
- Amazon DynamoDB
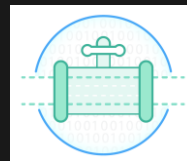
# Stock the Lake – Data Ingestion



**Amazon Glue ETL**

- Serverless ETL engine generates Python code that is entirely customizable, reusable, and portable.

**Database Migration Service**

Migrate your RDBMS into S3 (as well as other targets)

**Amazon Kinesis Streams**

- Build your own custom applications that process or analyze streaming data

**Amazon Kinesis Firehose**

- Easily load massive volumes of streaming data into S3, Amazon Redshift, and Amazon Elasticsearch Service

# Stock the Lake – Data Ingestion

# Demonstration

What is the speed of the data?
What is the source of the data?

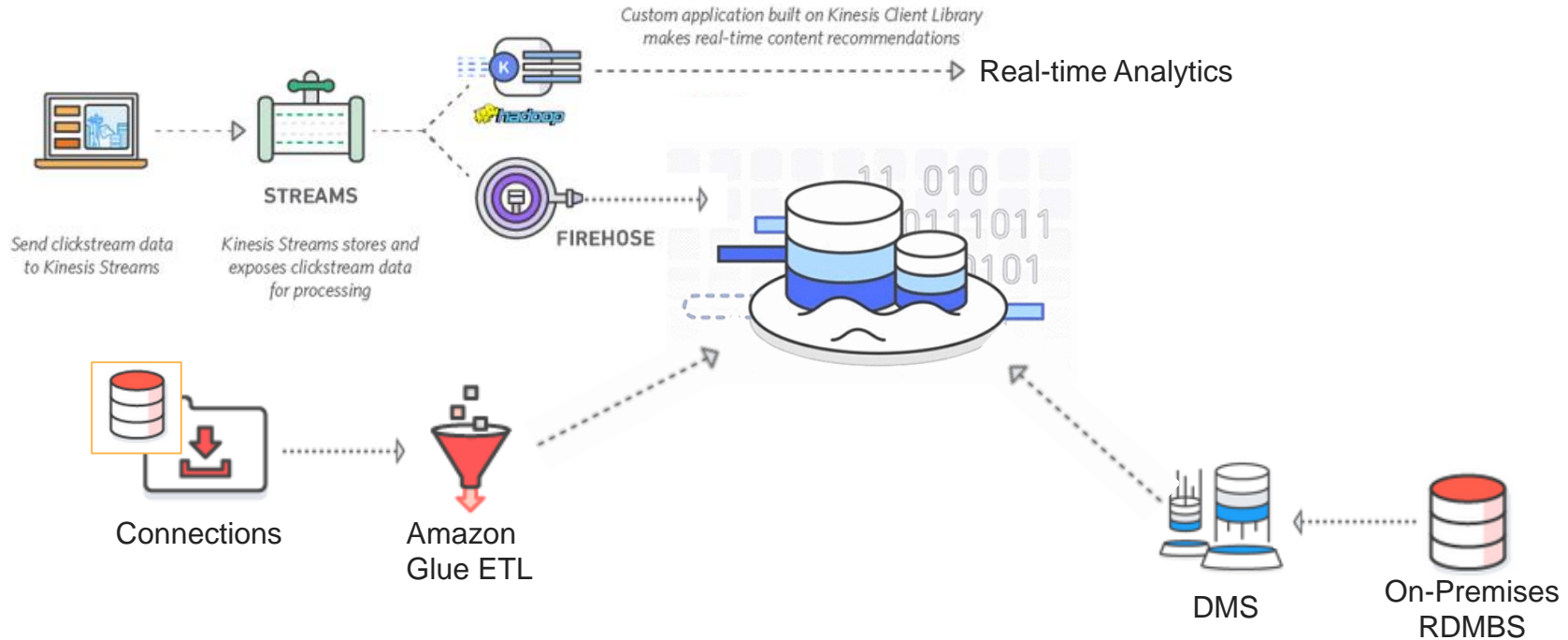# Fishing the Lake – Catalog/Search



Glue Data Catalog

Structural and Operational metadata

Amazon DynamoDB

Metadata/Tag Lookup

Search Content

### Diagram

Catalog & Search
*Access and search metadata*

Access & User Interface
*Give your users easy and secure access*

**Fish the Lake**

**Control the Access**

**Stock the Lake**

Central Storage
*Secure, cost-effective Storage in Amazon S3*

S3

**Process the Findings**

Data Ingestion
*Get your data into S3 Quickly and securely*

Protect and Secure
*Use entitlements to ensure data is secure and users' identities are verified*

Processing & Analytics
*Use of predictive and prescriptive analytics to gain better understanding*

**Govern the Lake**

# Demonstration

What type of data?

How is the data being queried?

# Govern the Lake



Catalog & Search
*Access and search metadata*

Access & User Interface
*Give your users easy and secure access*

**Fish the Lake**

**Control the Access**

Central Storage
*Secure, cost-effective Storage in Amazon S3*

**Stock the Lake**

S3

**Process/Analyze the Catch**

Data Ingestion
*Get your data into S3 Quickly and securely*

Protect and Secure
*Use entitlements to ensure data is secure and users' identities are verified*

Processing & Analytics
*Use of predictive and prescriptive analytics to gain better understanding*

**Govern the Lake**

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

Security Token Service

CloudWatch

CloudTrail

Key Management Service

**Temporary Tokens**

**Performance**

**Auditing**

**Encryption**

# Control the Access



Fish the Lake

Control the Access

Catalog & Search
*Access and search metadata*

Access & User Interface
*Give your users easy and secure access*

Central Storage
*Secure, cost-effective
Storage in Amazon S3*

S3

Stock the Lake

Process the Findings

Data Ingestion
*Get your data into S3
Quickly and securely*

Protect and Secure
*Use entitlements to ensure data is secure and users' identities are verified*

Processing & Analytics
*Use of predictive and prescriptive
analytics to gain better understanding*

Govern the Lake

API Gateway

Identity & Access
Management
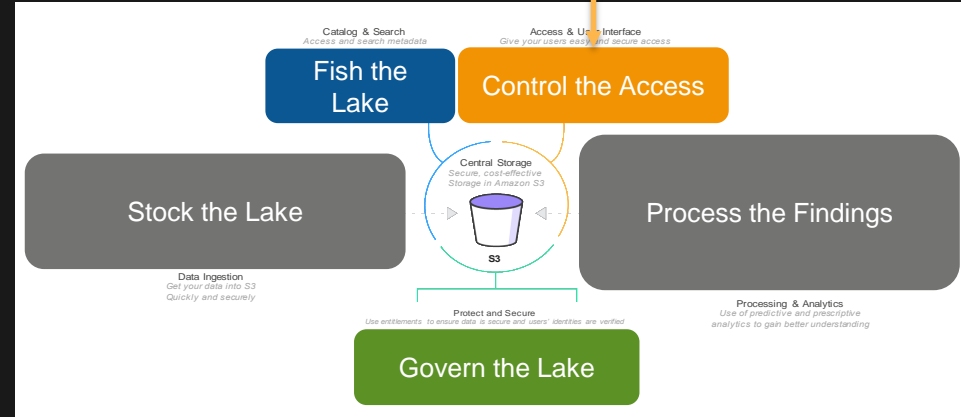
Cognito

Interfaces          Identity/Access          User Authentication

# Process/Analyze the Catch

## Processing & Analytics

### Real-time

- ElastiSearch Service
- Kinesis Analytics, Kinesis Streams
- Spark Streaming on EMR
- Apache Flink on EMR
- AWS Lambda
- Apache Storm on EMR

### Batch

- EMR Hadoop, Spark, Presto
- Redshift Data Warehouse
- Athena Query Service

### AI & Predictive

- Amazon Lex
- Amazon Polly
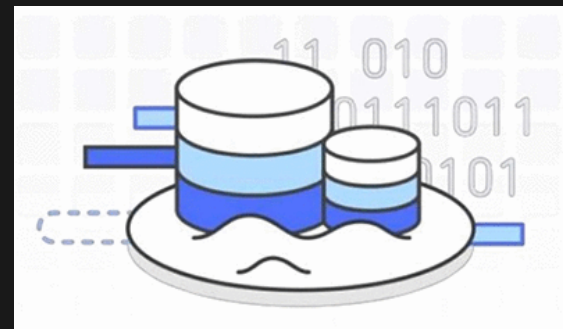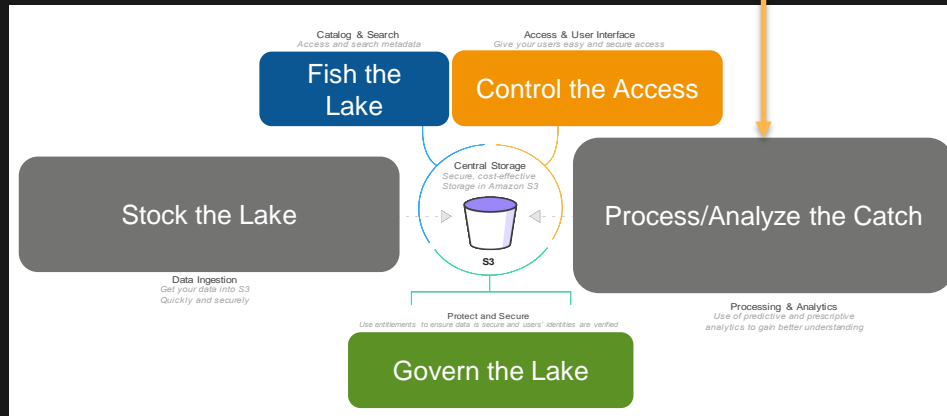- Amazon Rekognition
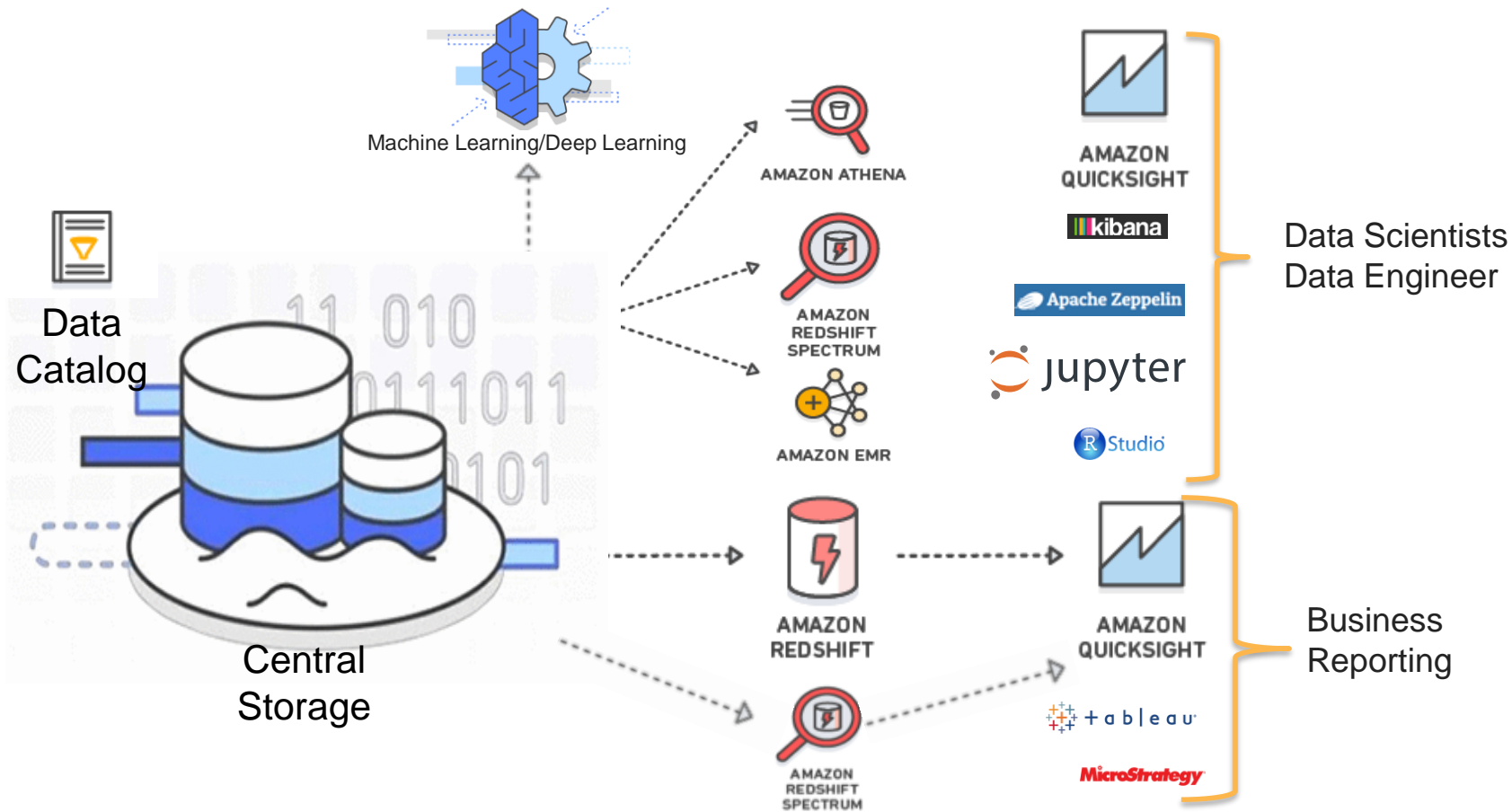- Machine Learning

### Transactional & RDBMS

- DynamoDB, NoSQL DB
- Aurora Relational Database

### BI & Data Visualization

---

Catalog & Search
*Access and search metadata*

Access & User Interface
*Give your users easy and secure access*

**Fish the Lake**

**Control the Access**

Central Storage
*Secure, cost-effective Storage in Amazon S3*

S3

**Stock the Lake**

**Process/Analyze the Catch**

Data Ingestion
*Get your data into S3 Quickly and securely*

Protect and Secure
*Use entitlements to ensure data is secure and users' identities are verified*

Processing & Analytics
*Use of predictive and prescriptive analytics to gain better understanding*

**Govern the Lake**

# Process/Analyze the Catch



Machine Learning/Deep Learning

Data Catalog

Central Storage

AMAZON ATHENA

AMAZON REDSHIFT SPECTRUM

AMAZON EMR

AMAZON REDSHIFT

AMAZON REDSHIFT SPECTRUM

AMAZON QUICKSIGHT

kibana

Apache Zeppelin

jupyter

R Studio

AMAZON QUICKSIGHT

tableau

MicroStrategy

Data Scientists
Data Engineer

Business Reporting

# Interactive query service



Amazon
Athena

- **Query directly from Amazon S3**
- **Use ANSI SQL**
- **Serverless**
- **Multiple data formats**
- **Pay per query**

Amazon
Elastic
MapReduce

**Hadoop/HDFS clusters**

**Hive, Pig, Impala, Hbase, Spark, Presto**

**Easy to use, fully managed**

**On-demand, reserved instance, and**

**Spot pricing**

**Tight integration with Amazon S3,**

**DynamoDB, and Kinesis**

# Resizable clusters

Easy to add and remove compute capacity on your cluster.

# ON A SINGLE MACHINE

COST: 4h x $1.06 = **$4.24**
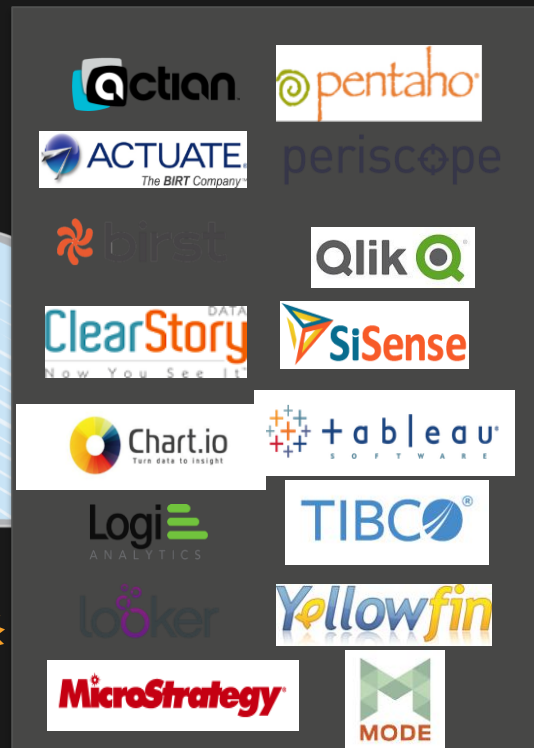PROCESSING TIME: **4h**

# ON MULTIPLE MACHINES

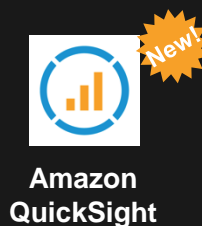COST: 4 x 1h x $1.06 = **$4.24**
PROCESSING TIME: **1h**

# Amazon Redshift works with third-party analysis tools

JDBC/ODBC

**Amazon Redshift**

actian
ACTUATE. The BIRT Company
birst
ClearStory DATA Now You See It
Chart.io Turn data to insight
Logi ANALYTICS
looker
MicroStrategy
pentaho
periscope
Qlik Q
SiSense
tableau SOFTWARE
TIBCO
Yellowfin
MODE

New!

**Amazon QuickSight**

# Amazon Redshift has security built in

SSL to secure data in transit

Encryption to secure data at rest
- AES-256; hardware accelerated
- All blocks on disks and in Amazon S3 encrypted
- HSM support

No direct access to compute nodes

Audit logging, AWS CloudTrail, AWS KMS integration

Amazon VPC support

SOC 1/2/3, PCI-DSS Level 1, FedRAMP, HIPAA

Customer VPC

SQL Clients/BI Tools

JDBC/ODBC

Internal VPC

Leader Node

10 GigE (HPC)

Compute Node

Compute Node

Compute Node

Ingestion Backup Restore

Amazon S3/Amazon DynamoDB

# Redshift Spectrum

Leverages Amazon Redshift's advanced cost-based optimizer

Pushes down projections, filters, aggregations and join reduction
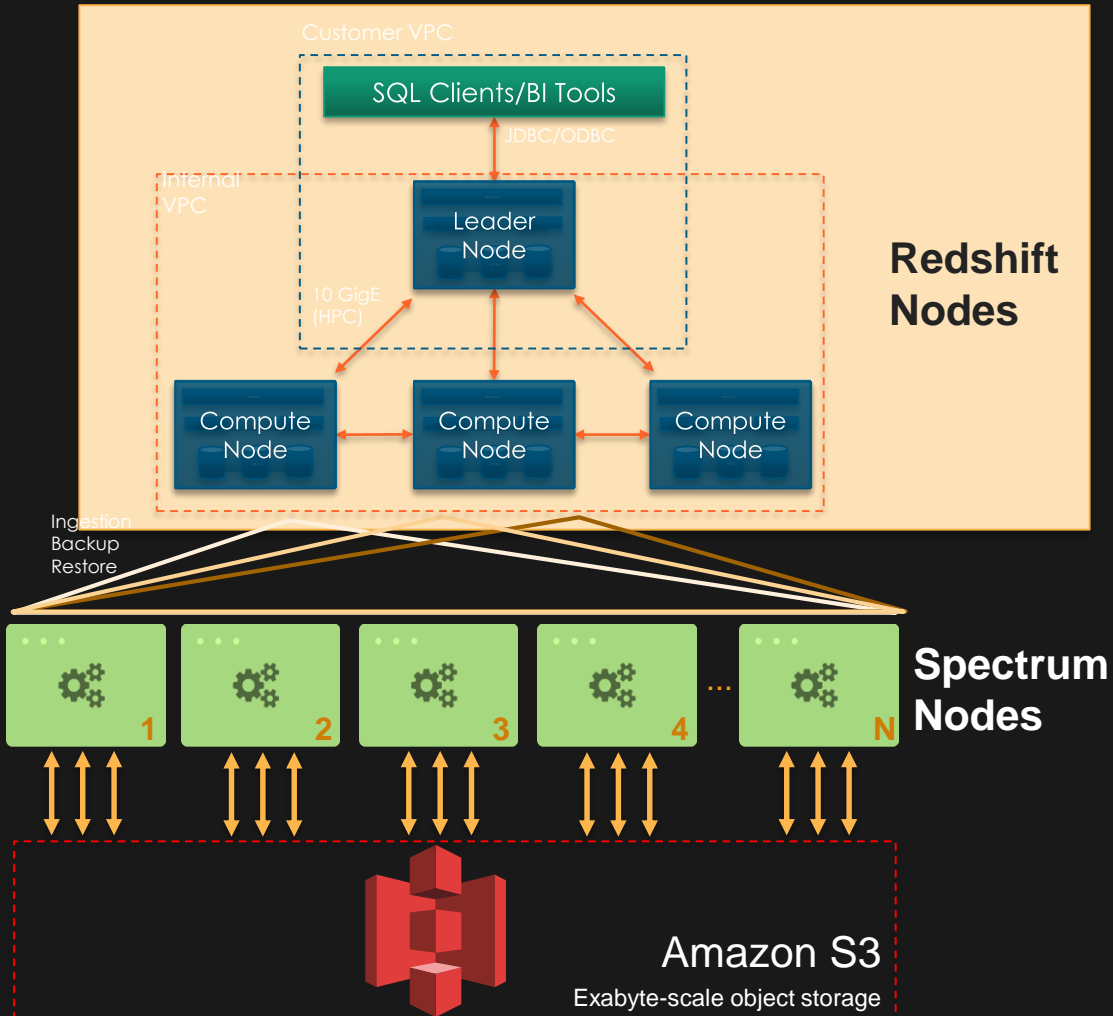
Dynamic partition pruning to minimize data processed

Automatic parallelization of query execution against Amazon S3 data

Efficient join processing within the Amazon Redshift cluster

Data Catalog
Apache Hive Metastore

Customer VPC

SQL Clients/BI Tools

JDBC/ODBC

Internal VPC

Leader Node

10 GigE (HPC)

Compute Node

Compute Node

Compute Node

Redshift Nodes

Ingestion Backup Restore

1  2  3  4  ...  N

Spectrum Nodes

Amazon S3
Exabyte-scale object storage

# Amazon AI
## Intelligent services powered by deep learning

# Proven customer success

The vast majority of big data use cases deployed in the cloud today run on AWS.

# Case study: Re-architecting compliance

> *"For our market surveillance systems, we are looking at about 40% [savings with AWS], but the real benefits are the business benefits: We can do things that we physically weren't able to do before, and that is priceless."*
> - Steve Randich, CIO

**FINRA**

## What FINRA needed
- Infrastructure for its market surveillance platform
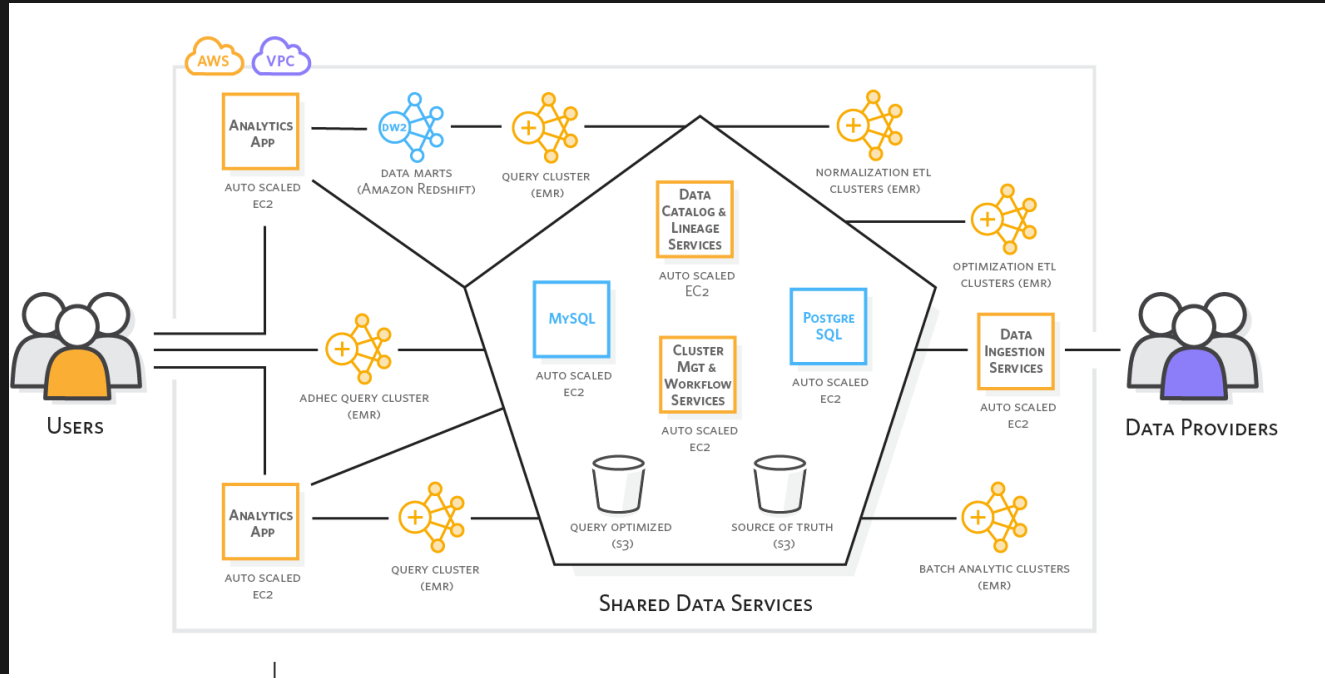- Support of analysis and storage of approximately 75 billion market events every day

## Why they chose AWS
- Fulfillment of FINRA's security requirements
- Ability to create a flexible platform using dynamic clusters (Hadoop, Hive, and HBase), Amazon EMR, and Amazon S3
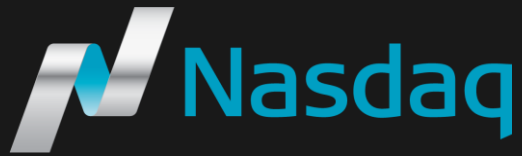
## Benefits realized
- Increased agility, speed, and cost savings
- Estimated savings of $10-20M annually by using AWS

# Fraud detection



FINRA uses Amazon EMR and Amazon S3 to process up to 75 billion trading events per day and securely store over 5 petabytes of data, attaining savings of $10-20M per year.

- Nasdaq implements an Amazon S3 data lake + Amazon Redshift data warehouse architecture
- Most recent two years of data is kept in the Amazon Redshift data warehouse and snapshotted into Amazon S3 for disaster recovery
- Data between two and five years old is kept in Amazon S3
- Presto on Amazon EMR is used to ad-hoc query data in Amazon S3
- Transitioned from an on-premises data warehouse to Amazon Redshift & Amazon S3 data lake architecture
- Over 1,000 tables migrated
- Average daily ingest of over 7B rows
- *Migrated* off legacy DW to AWS (start to finish) *in 7 man-months*
- AWS costs were *43%* of legacy budget for the same data set (~1100 tables)

# Building a Data Lake

An **architectural approach** that allows you to store massive amounts of "**raw**" data into a central location

It's **readily available** to be **categorized, processed, analyzed, and consumed** by **diverse groups**