# Turn Data into Insights with Data Lakes and Analytics on AWS

Osemeke Isibor

Partner Solutions Architect, AWS.

29th June 2020

# Data is a strategic asset for every organization

> " The world's most valuable resource is no longer oil, but data. *

# Data Analytics Workflow



Generation → Collection & Storage → Analytics & Computation → Consume

aws

# Traditionally data Analytics revolved around data warehouses

# Customers want more value from their data

Growing
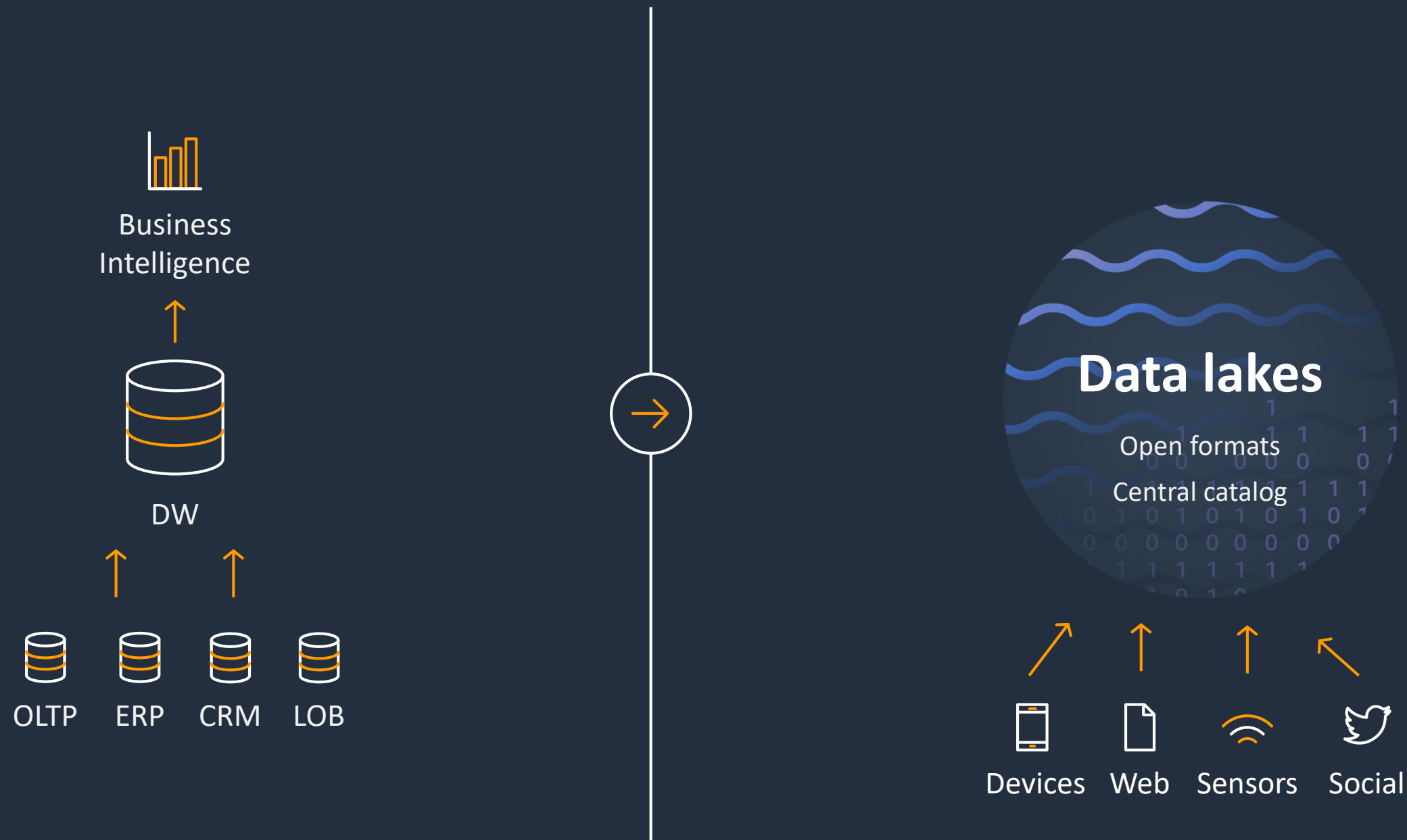exponentially

From new
sources

Increasingly
diverse

Used by
many people

Analyzed by
many applications

aws

# Customers moving to data lake architectures

Business
Intelligence

DW

OLTP    ERP    CRM    LOB

**Data lakes**

Open formats

Central catalog

Devices    Web    Sensors    Social

aws

# Bringing together the best of both worlds

Data Warehousing

Analytics

Machine Learning

Data lake

Extends or evolves DW architectures

Store any data in any format

Durable, available, and exabyte scale

Secure, compliant, auditable

Run any type of analytics from DW to Predictive

aws

# Why choose AWS for data lakes and analytics?

# 1. Easiest to build data lakes and analytics

**Amazon S3**: Highly durable, secure and scalable single storage layer for all analytics and ML

**AWS Lake Formation:** A service to build secure data lakes in days instead of months.

Deep integration across analytics & infrastructure(including federated queries)

---

The fastest way to go from zero to insights, covering all data for all users

aws

# 2. Most secure infrastructure for analytics
## Services for security and governance

## Security

- Amazon GuardDuty
- AWS Shield
- AWS WAF
- Amazon Macie
- VPC

## Identity

- AWS IAM
- AWS SSO
- Amazon Cloud Directory
- AWS Directory Service
- AWS Organizations

## Encryption

- AWS Certification Manager
- AWS Key Management Service
- Encryption at rest
- Encryption in transit
- Bring your own keys, HSM support

## Compliance

- AWS Artifact
- Amazon Inspector
- Amazon Cloud HSM
- Amazon Cognito
- AWS CloudTrail

aws

# 2. Most secure infrastructure: certifications

## Global

**CSA**
Cloud Security Alliance Controls

**ISO 9001**
Global Quality Standard

**ISO 27001**
Security Management Controls

**ISO 27017**
Cloud Specific Controls

**ISO 27018**
Personal Data Protection

**PCI DSS Level 1**
Payment Card Standards

**SOC 1**
Audit Controls Report

**SOC 2**
Security, Availability, & Confidentiality Report

**SOC 3**
General Controls Report

## United States

**CJIS**
Criminal Justice Information Services

**DoD SRG**
DoD Data Processing

**FedRAMP**
Government Data Standards

**FERPA**
Educational Privacy Act

**ISO FFIEC**
Financial Institutions Regulation

**FIPS**
Government Security Standards

**FISMA**
Federal Information Security Management

**GxP**
Quality Guidelines and Regulations

**HIPPA**
Protected Health Information

**ITAR**
International Arms Regulations

**MPAA**
Protected Media Content

**NIST**
National Institute of Standards and Technology

**SEC Rule 17a-4(f)**
Financial Data Standards

**VPAT/Section 508**
Accountability Standards

## Asia Pacific

**FISC [Japan]**
Financial Industry Information Systems

**IRAP [Australia]**
Australian Security Standards

**K-ISMS [Korea]**
Korean Information Security

**MTCS Tier 3 [Singapore]**
Multi-Tier Cloud Security Standard

**My Number Act [Japan]**
Personal Information Protection

## Europe

**C5 [Germany]**
Operational Security Attestation

**Cyber Essentials Plus [UK]**
Cyber Threat Protection

**G-Cloud [UK]**
UK Government Standards

**IT-Grundschutz [Germany]**
Baseline Protection Methodology

aws

# 3. Most comprehensive and open

**Data, visualization, engagement, & machine learning**

Data

Dashboards

Digital User Engagement

Predictive Analytics

**Analytics**

Data Warehousing

Big Data Processing

Serverless Data processing

Interactive Query

Operational Analytics

Real time Analytics

**Data lake infrastructure & management**

Infrastructure

Security & Management

Data Catalog & ETL

**Data movement**

**Migration & Streaming Services**

aws

# 3. Most comprehensive and open

## Data, visualization, engagement, & machine learning

NEW

Data Exchange    QuickSight    Pinpoint    SageMaker    Comprehend    Lex    Polly    Rekognition    Translate

+ many more

## Analytics

Redshift    EMR (Spark & Hadoop)    AWS Glue (Spark & Python)    Athena    Elasticsearch Service    Kinesis Data Analytics

## Data lake infrastructure & management

S3/Glacier    Lake Formation    AWS Glue

## Data movement

Database Migration Service | Snowball | Snowmobile | Kinesis Data Firehose | Kinesis Data Streams | Managed Streaming for Apache Kafka

aws

# 3. Open standards, formats, and Apache open source

| | | |
|---|---|---|
| Flink | Mahout | PyTorch |
| Ganglia | MapReduce | R |
| Hbase | MxNET | Scala |
| HCatalog | MySQL | Spark |
| HDFS | Oozie | Sqoop |
| Hive | ORC | SQL |
| Hudi | Parquet | TensorFlow |
| Java | Phoenix | Tez |
| JupyterHub | Pig | YARN |
| Kafka | Presto | Zeppelin |
| Livy | Python | Zookeeper |

aws

# 4. Most scalable, cost-effective, high-performance infrastructure for analytics

On-demand, Reserved, and Spot instances to reduce costs

100 Gbps bandwidth network interfaces for performance

Industry leading choice of 200+ instance types to meet workload needs

Five highly available storage tiers and intelligent tiering

aws

# 4. Most scalable, cost-effective infrastructure for analytics

Some examples of advanced cost saving capabilities in analytics services

## EMR

Autoscaling

57% less than on-premises per IDC report

## Redshift

Less than 1/10th of the cost of traditional, on-premises solutions

## Athena & QuickSight

Serverless pay only for what is used

Pricing per session for visualization

aws

# Get answers from your data in AWS Data Lake

# Query Directly with Amazon Athena

# Analyze with Hadoop on Amazon EMR

# Create Visualizations with Amazon QuickSight

# Train ML Models with Amazon SageMaker

# Fully integrate with other AWS Services

aws

# More data lakes and analytics than anywhere else

## Tens of thousands of data lakes run on AWS across all industries

# Complemented by AWS Partner Network (APN) Solutions providers

### Collection & preparation

ATTUNITY

Informatica

MATILLION

Paxata

talend

TRIFACTA

### Governance

collibra

DATAGUISE

DATAGUISE

DATA REPUBLIC

Informatica

ZALONI
THE DATA LAKE COMPANY

### Visualization

SISENSE

MicroStrategy

looker

tableau
SOFTWARE

TIBCO

aws

# Get help from Data & analytics APN consulting partners

## GLOBAL

accenture
Deloitte.
wipro
Infosys®
Mindtree
Welcome to possible

## JAPAN

class method
cloud pack
TECHORUS

## CHINA

博思云为 bosicloud
eCloud valley

## LATAM

Morris & Opazo
Business Solutions

## NORTH AMERICA

1Strategy
47Lining
A Hitachi Vantara Company
8K Miles
AGILISIUM
BLUESENTRY
CAMBRIDGE TECHNOLOGY
clearscale
Cloudwick
CORECOMPETE
MACTORES
VERTICAL TRAIL SOLUTIONS
NorthBay
Provectus
ONICA
PARIVEDA SOLUTIONS
slalom
softserve
iOLAP
TEKsystems
Our people make IT possible.

## EMEA

BEEVA
claranet
CloudMas
D2SI by devoteam
CLOUDREACH
EMET E&M COMPUTING
IPPON Discovery to Delivery
KCOM
Reply storm
SOLITA
dnm SEE CLEARER. ACT SMARTER.
tecRACER Cloud Enabling Your Business
DoiT INTERNATIONAL
keepler

## APAC

blazeclan Cloud IT Better
Powerup Cloud
TO THE NEW
bluepi

aws

# Learn analytics with AWS Training and Certification

**Resources created by the experts at AWS to help you build and validate data analytics skills**

New free digital course: "Data Analytics Fundamentals"

Classroom offerings, including "Big Data on AWS", feature AWS expert instructors and hands-on labs

Validate expertise with the "AWS Certified Big Data—Specialty" exam or the new "AWS Certified Data Analytics—Specialty" beta exam

Visit aws.amazon.com/training/paths-specialty/

aws

**Customer Success Powered by AWS.**

**FINRA oversees > 3,000 securities firms doing business in the United States.**

**Challenge:**

FINRA's legacy system did not scale well

- Up to 75 billion events per day

- Run complex surveillance queries over 20+ PB of data

**Solution:**

- Migrated their big data appliance to a S3 Data Lake and used EMR for ingestion and processing

- Migrated to RDS and testing Aurora

# FINRA uses S3 to Build Data Lake with EMR



- Required fast access across trillions of trade records (20PB+)

- Migrated from on-premises system

- Use Apache HBase on Amazon EMR to store and serve this data

- Use EMR engines— Spark, Presto, and Hive to process data

- Lower costs by 60% over on-premises system

aws

Nasdaq operates financial exchanges around the world, and processes
large volumes of data.

**Challenge:**

Nasdaq wanted to make their large historical data footprint available
to analyze as a single dataset.

**Solution:**

- Use Amazon Redshift for interactive querying

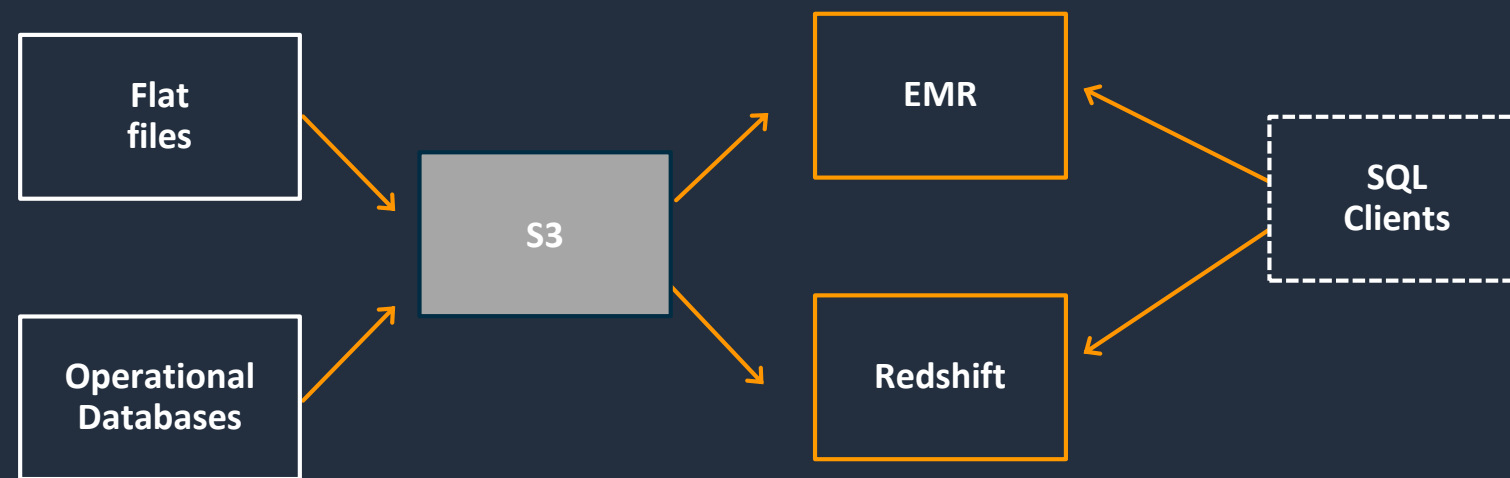- Use Amazon S3 as a Data Lake, and Presto on EMR to process historical data

# Nasdaq Uses AWS to Build a Data Lake



Data from all 7 exchanges operated by Nasdaq
(orders, quotes, trade executions)

- Migrate legacy on-premises warehouse to Amazon Redshift

- **4.8B rows** inserted per trading day (orders, trades, quotes)

- Ingest data from multiple sources, validates, and stages in S3

- **Redshift reads data out of S3** for fast queries

- Presto on EMR and S3 used for analysis of massive historical data set
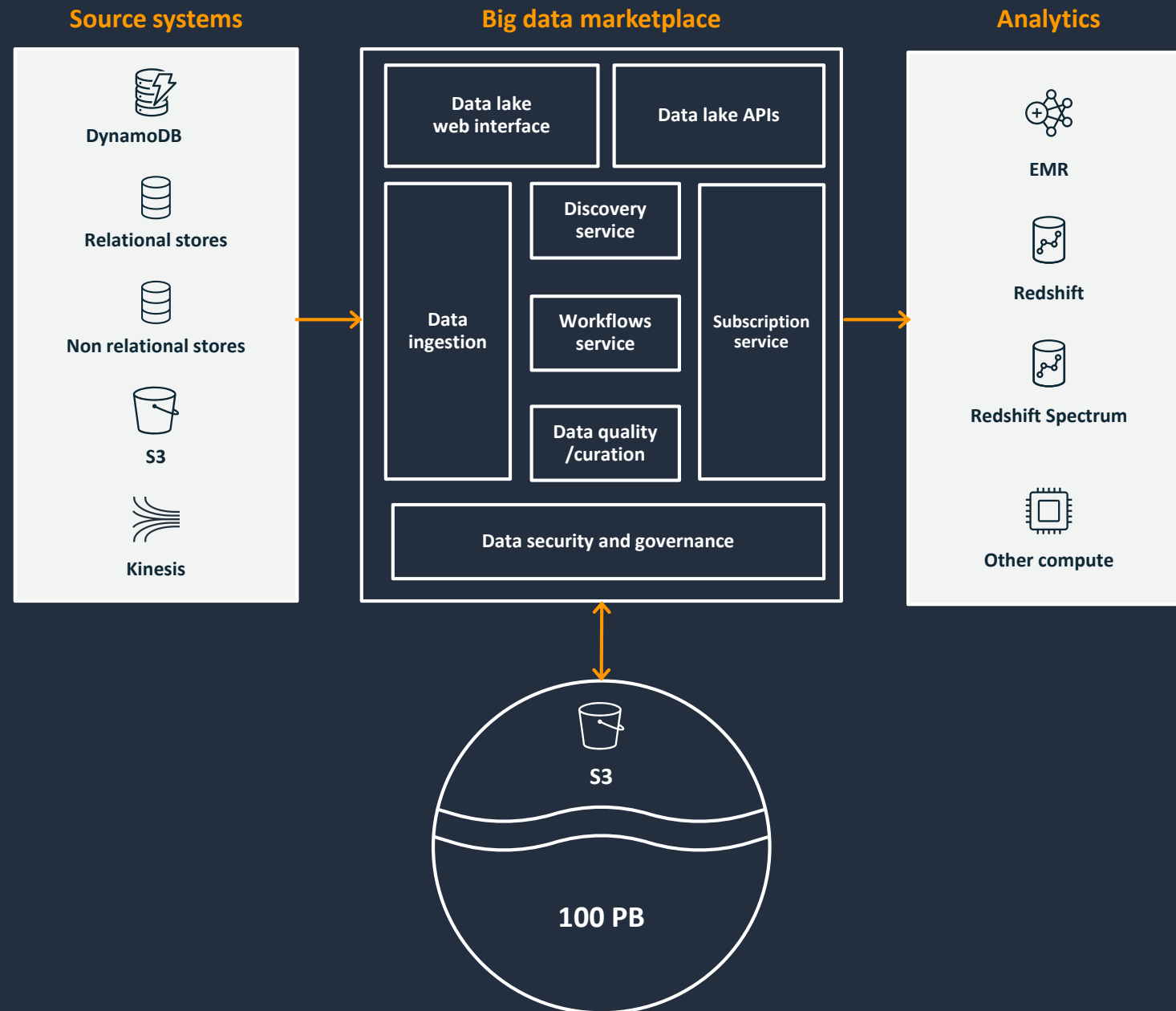
# amazon.com

## Challenge:

- Amazon needed to analyze a massive amount of data to find insights, identify opportunities, and evaluate business performance.

- The Oracle DW did not scale, was difficult to maintain, and costly.

## Solution:

- Amazon deployed a data lake with Amazon S3, and now runs analytics with Amazon Redshift, Redshift Spectrum, and Amazon EMR.

# EQUINOX

Equinox Fitness has a number of health and fitness brands

## Challenge:

- Their data warehouse had limited integration.

- They needed to reduce administration and costs, blend structured and semi-structured data for analytics, and evolve into a data lake strategy.

## Solution:

- Migrated to Amazon Redshift to combine data from disparate sources.

- They land data directly in an Amazon S3 data lake and perform analytics using Amazon Redshift, Redshift Spectrum, and Amazon EMR.

# 80% cost savings by migrating to Amazon Redshift

**EQUINOX**



- Migrated from Teradata data warehouse

- Built a DW with Redshift and data lake with S3
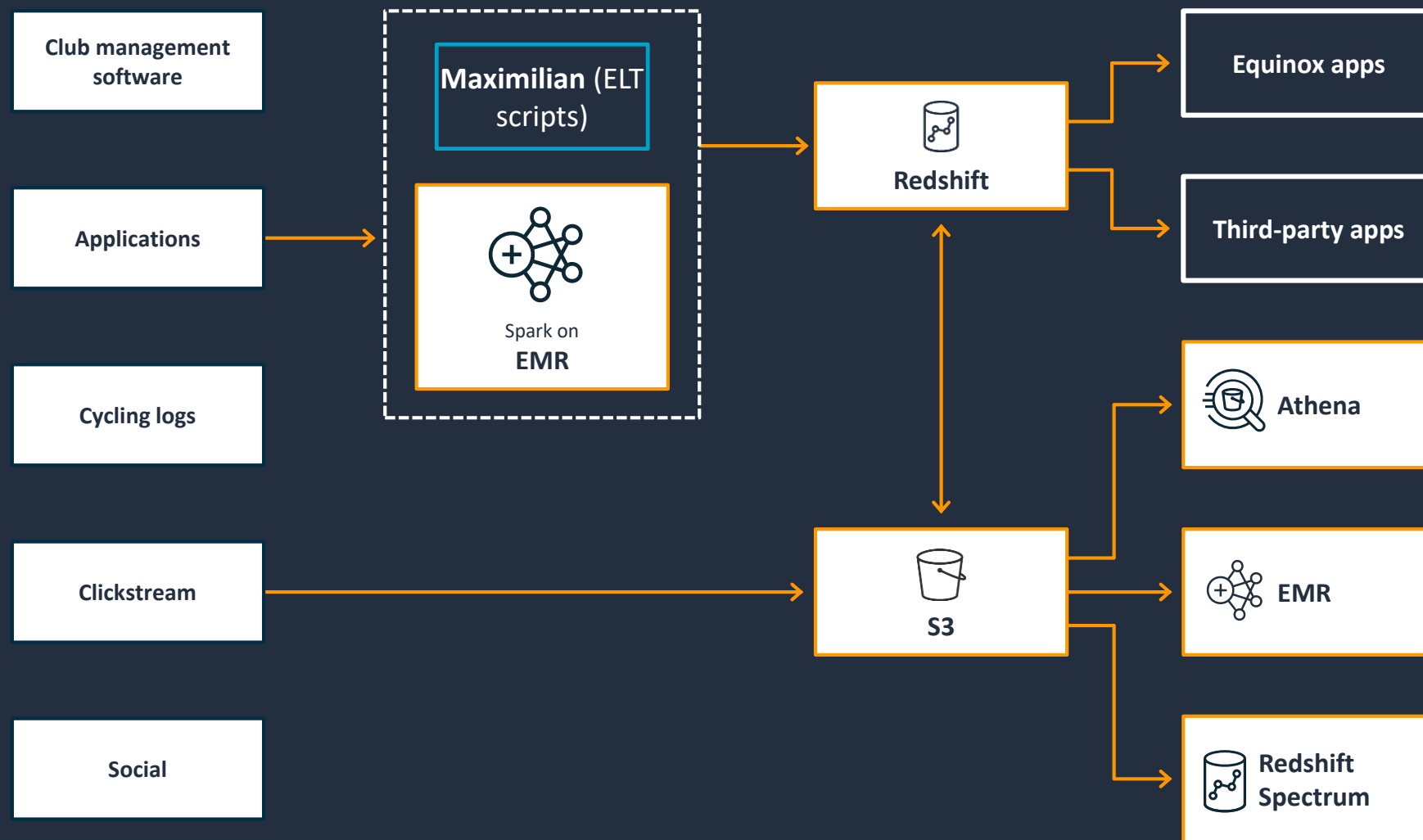
- Analytics on data lake with Amazon Athena, Amazon Redshift Spectrum, and Amazon EMR

- Increased user productivity to move faster

- Amazon Redshift costs ~20% of its original Teradata maintenance and support

- Report time reduced from months to days

aws

# Next steps...

**(1)** **Sign up for an AWS account**

Instantly get access
to the AWS Free Tier

**(2)** **Learn with 10-minute tutorials**

Explore and learn with
simple tutorials

**(3)** **Start building with AWS**

Begin building with step-by-step
guide to help you launch
your AWS project

aws

# AWS Data Exchange
## Easily find and subscribe to 3rd-party data in the cloud

### Quickly find diverse data in one place

>1,000 data products

>80 data providers including include Dow Jones, Change Healthcare, Foursquare, Dun & Bradstreet, Thomson Reuters, Pitney Bowes, Lexis Nexis, and Deloitte

### Easily analyze data

Download or copy data to S3

Combine, analyze, and model with existing data

Analyze data with EMR, Redshift, Athena, and AWS Glue

### Efficiently access 3rd party data

Simplifies access to data: No need to receive physical media, manage FTP credentials, or integrate with different APIs

Minimize legal reviews and negotiations

aws

# Feedback Survey

Please let us know what you think about the session.

Get help from AWS with your analytics project.



Scan QR code to complete survey

As a token of appreciation, we'll mail you an AWS souvenir

aws

# Thank you!

iosemeke@amazon.com

aws