



# Cloud Data Lake DevDay

# Agenda

<b>9:00-9:30 AM</b>	Opening Remarks - Enabling Cloud Data Lakes for Analytics
<b>9:30-9:45 AM</b>	Customer Use Cases
<b>9:45-9:55 AM</b>	Break
<b>9:55-10:40 AM</b>	Architecture - Incorporating Roles; Scaling, Launching and Managing Clusters; Managing Pools and Containers
<b>10:40-11:25 AM</b>	Notebooks - Data Loading using Scala, Machine Learning Analytics using Python, and Redshift integration for data delivery to analysts
<b>11:25-11:35 PM</b>	Q&A

# Current state of analytics



We have more data than ever before, often locked in systems and formats where we can't access it for analytics

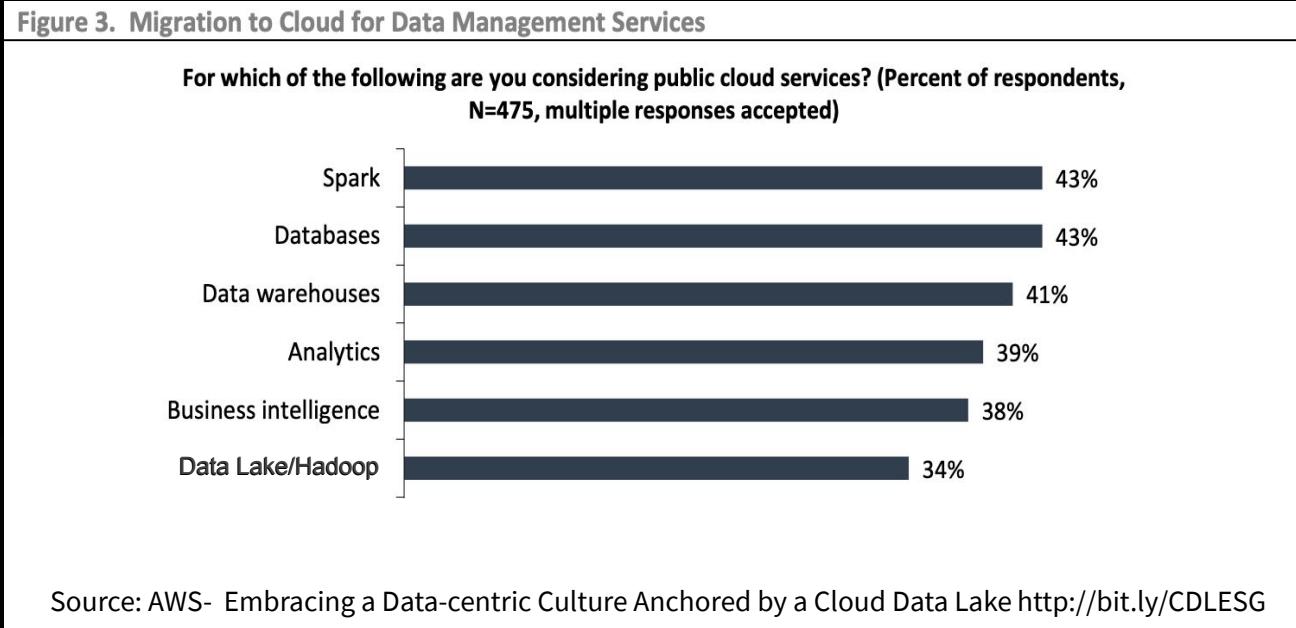


On premises systems require excess capacity to process peak needs



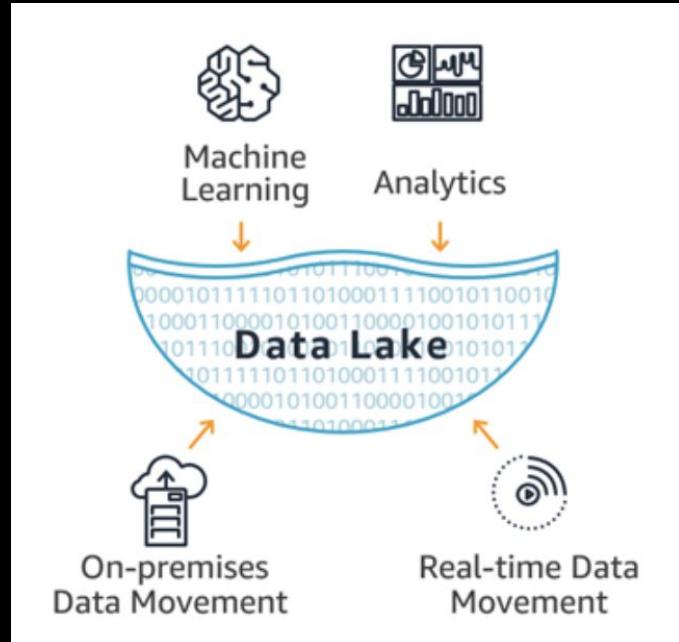
Data security is costly and experts are hard to find

# Analytic platforms are moving to the cloud



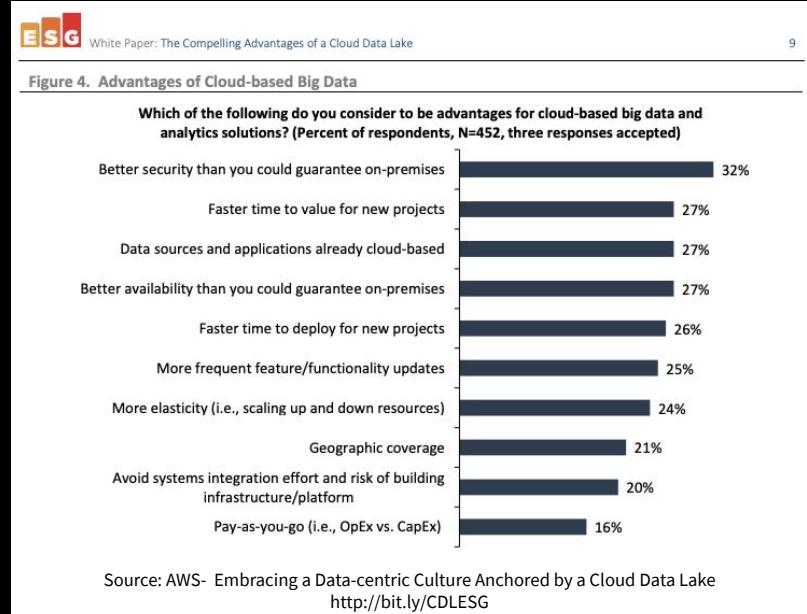
# What is a data lake?

A data lake is a centralized repository that allows you to store all your structured and unstructured data at any scale. Organizations see this as their golden repository.



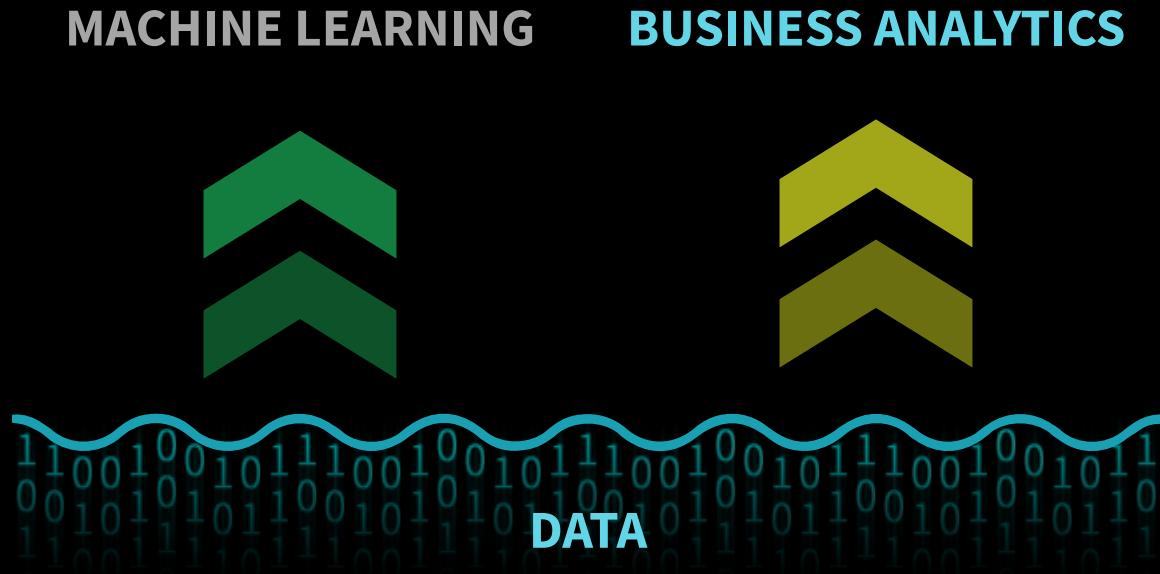
# Benefits of cloud data lakes

- Pervasive security features
- Performance and scalability
- Reliability and availability
- Economics
- Integration
- Agility



For organizations that are running analytics and machine learning workloads, the cloud makes the most sense, providing the balance of infinite scale and pay by use.

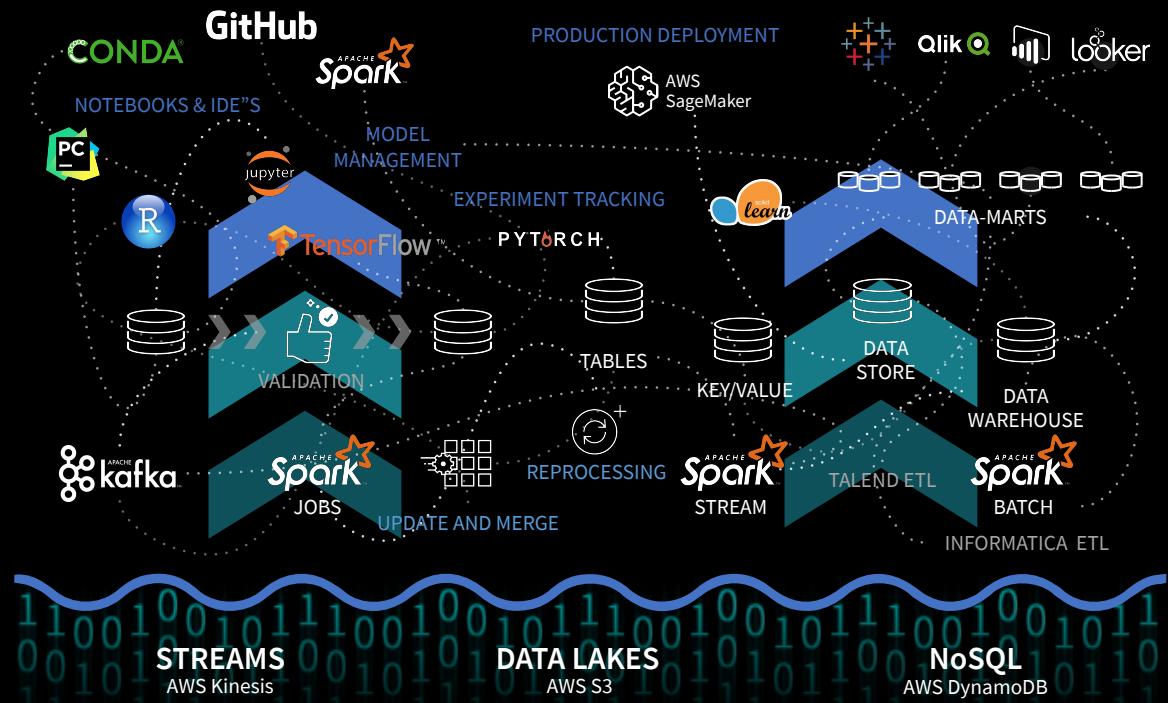
Organizations want to combine their massive **data lake** with **ML** and **BI** capabilities to unlock business value



## DATA SCIENCE AND MACHINE LEARNING

## BUSINESS ANALYTICS AND REPORTING

... but organizations fail to unlock business value due to data, technology and people silos



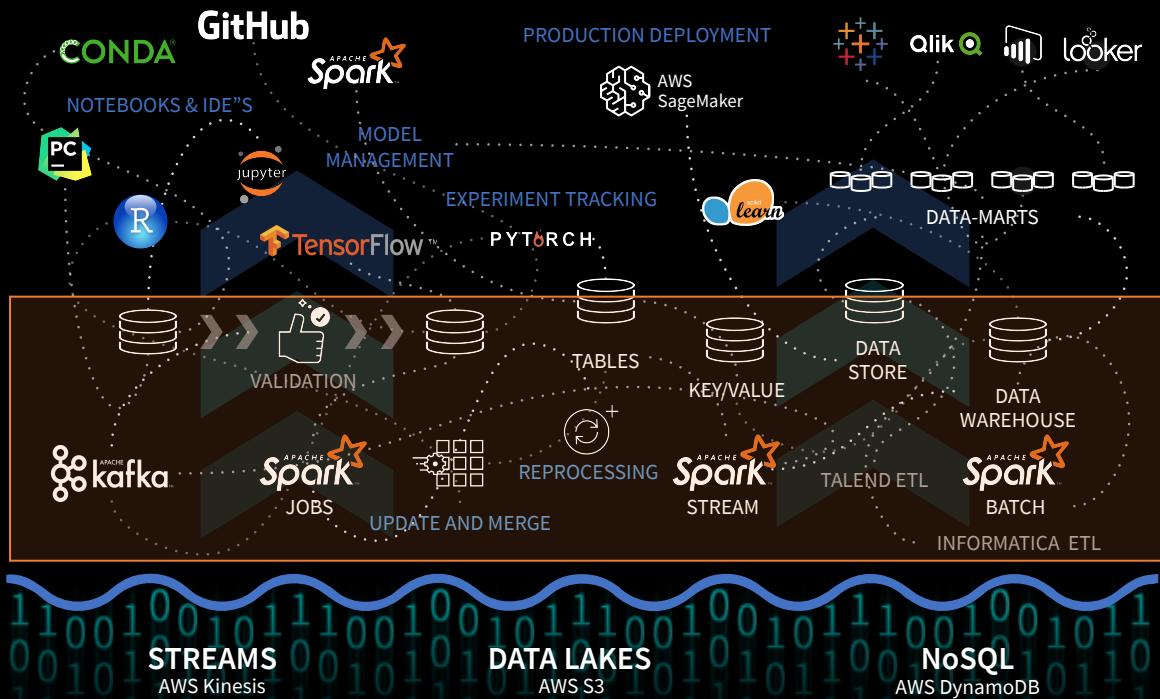
# Why do organizations fail to unlock business value?

## Data quality and reliability

Multiple copies of inconsistent data built with unreliable data pipelines. Organizations want to guarantee data quality while retaining flexibility to ingest all data.

## DATA SCIENCE AND MACHINE LEARNING

## BUSINESS ANALYTICS AND REPORTING



# Why do organizations fail to unlock business value?

## Disparate technologies

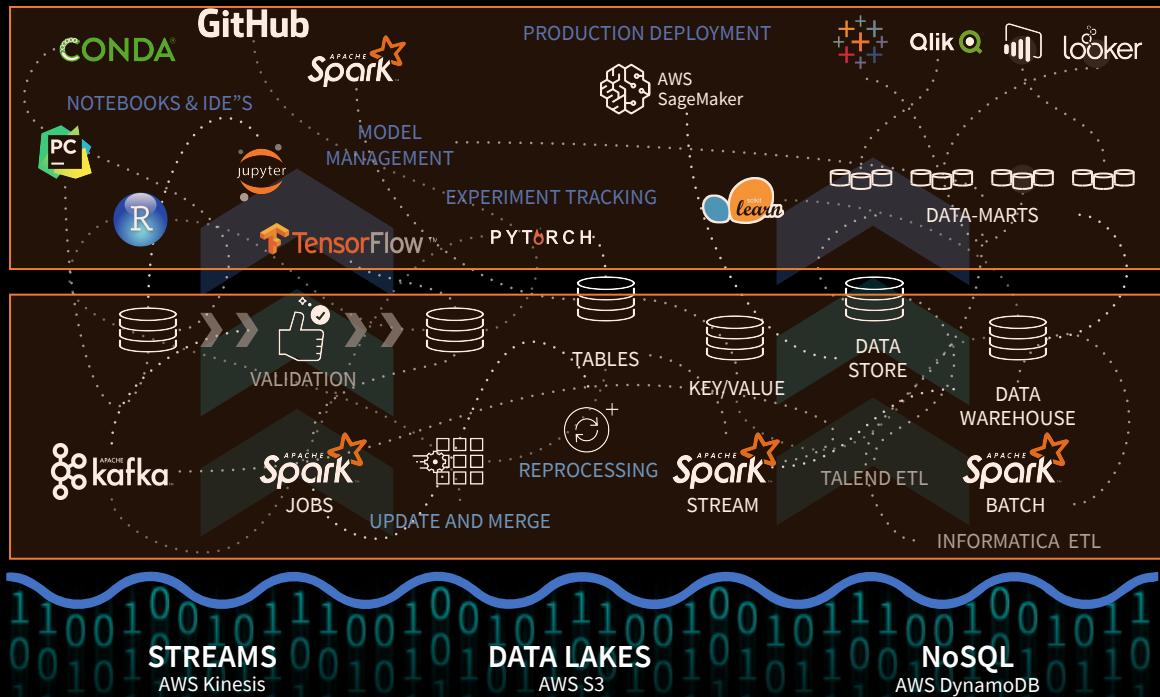
Need to stitch different IDE's, ML frameworks and deployment tools together

## Data quality and reliability

Multiple copies of inconsistent data built with unreliable data pipelines

## DATA SCIENCE AND MACHINE LEARNING

## BUSINESS ANALYTICS AND REPORTING



# Why do organizations fail to unlock business value?

## Disparate technologies

Need to stitch different IDE's, ML frameworks and deployment tools together

## Data quality and reliability

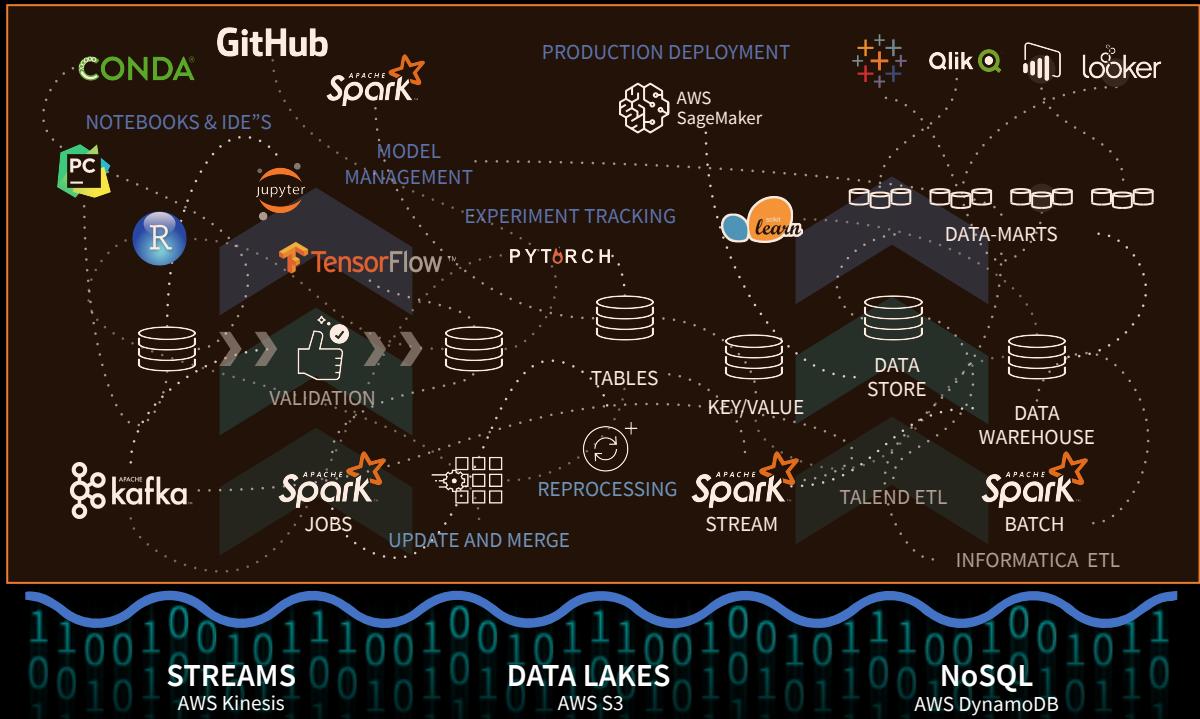
Multiple copies of inconsistent data built with unreliable data pipelines

## Fragmented security

End-to-end security with enterprise SLA's is a nightmare

## DATA SCIENCE AND MACHINE LEARNING

## BUSINESS ANALYTICS AND REPORTING

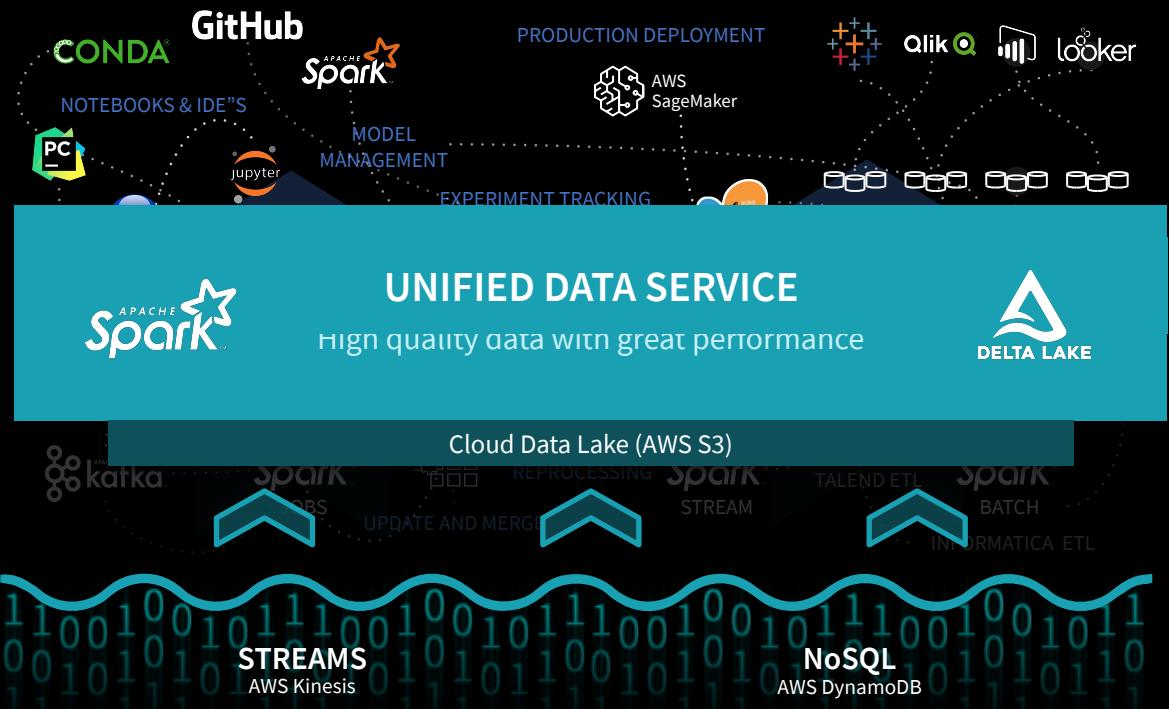


Databricks accelerates  
data-driven innovation

## DATA SCIENCE AND MACHINE LEARNING

## BUSINESS ANALYTICS AND REPORTING

Scalable, reliable, and  
fast data - built on your  
existing data lake



Databricks accelerates  
data-driven innovation

Collaborative workspace  
for data teams across the  
full lifecycle

Scalable, reliable, and  
fast data - built on your  
existing data lake

## DATA SCIENCE AND MACHINE LEARNING

## BUSINESS ANALYTICS AND REPORTING

### DATA SCIENCE WORKSPACE



**mlflow**

ML Runtime

P Y T H O N TensorFlow PyTorch



### BI INTEGRATIONS



Qlik



looker

Access all your data

### UNIFIED DATA SERVICE



High quality data with great performance



DELTA LAKE

Cloud Data Lake ( AWS S3 )



Databricks accelerates  
data-driven innovation

Collaborative workspace  
for data teams across the  
full lifecycle

Scalable, reliable, and  
fast data - built on your  
existing data lake

One fully-integrated  
security model for  
production infrastructure

## DATA SCIENCE AND MACHINE LEARNING

## BUSINESS ANALYTICS AND REPORTING

### DATA SCIENCE WORKSPACE



**mlflow**

ML Runtime

P Y T H O N TensorFlow PyTorch

### BI INTEGRATIONS



Qlik



looker

Access all your data



### UNIFIED DATA SERVICE

High quality data with great performance



DELTA LAKE

### ENTERPRISE CLOUD SERVICE



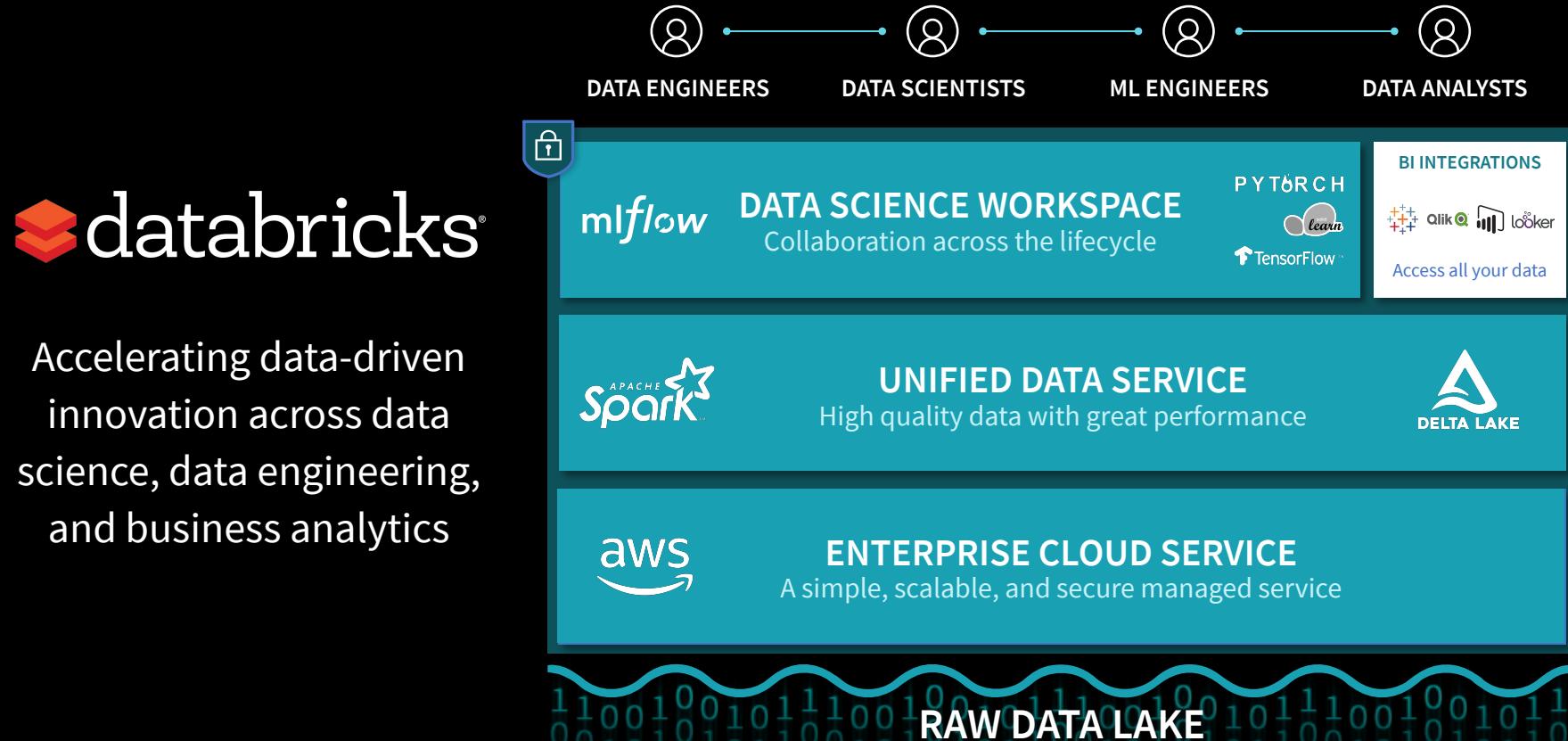
Enterprise  
Security

Simple  
Administration

Production  
Scale



# UNIFIED DATA ANALYTICS PLATFORM



DEV DAY

aws | databricks

# Successful cloud analytics implementations

Cloud Migration \$100M Fin Srv Company	Pipeline Automation \$1B Marketing Svcs Org	Retail ML \$36B Retail Company	Fraud Detection \$27B Fin Srv Company
<ul style="list-style-type: none"><li>Previously running Hadoop on premises</li><li>Maintenance <b>cost \$2M/yr</b></li><li>Migrated to Databricks on AWS</li><li>Estimated costs savings <b>of almost \$1M.</b></li><li>Much more capable platform to drive their business</li></ul>	<ul style="list-style-type: none"><li>Previous projects were <b>costing \$65-70 per run</b> with Hadoop on premises</li><li>Databricks ran faster and provided expertise to optimize these jobs</li><li>Reduced cost to <b>under \$1 per run</b></li></ul>	<ul style="list-style-type: none"><li>Enabled a DevOps approach to data science and productionized machine learning models</li><li>Now used by every data science team in the organization</li><li>Annual <b>savings of \$15M</b> from a unified data analytics platform</li></ul>	<ul style="list-style-type: none"><li>Established statistically-based fraud detection; running <b>successfully in production for several years</b></li><li>Extending to detecting phone motion and key entry patterns to detect bots and other automated applications</li></ul>

Estimated cost savings up to \$1M/year

Created automated pipelines that enable hundreds of machine learning models

Now powering every team and domain with data science and machine learning

Developing innovative AI applications with huge business implications

# Today you will learn

- Keys to implement a cloud data lake
- How to prepare data for analytics
- How to create automated data pipelines
- How to create an infinite storage capacity that scales economically

These lessons will enable you to have a huge impact on your organization



# The Role of Data Lakes

# Cloud data lakes are great for data storage

Data Lake is a file system that supports

- High-velocity data
- Open storage
- Separate storage from compute
- Many structures and even unstructured data
- Lower-cost storage



RAW DATA LAKE

# Organizations want to operationalize

To operationalize data lakes, you need features you expect on a database

- Transactions
- Recovery
- Snapshots
- Indexes

But with the benefits of a data lake:

- No schema, flexible schema or fixed schema
- Integrated stream and batch processing
- Open format so your data is not locked and subject to vendor tax.



RAW DATA LAKE

# Cloud data lake blueprint

- 1 Manage your ingest process - handle batch and streaming data
- 2 Transform to open formats like Parquet
- 3 Utilize ACID transactions to prevent partial writes
- 4 Curate your data and refine it in a series of steps
- 5 Setup a data catalog for your data
- 6 Define centralized security, governance and audit procedures
- 7 Distribute results to LOB units for action



RAW DATA LAKE



# Introducing Delta Lake

# A new standard for building data lakes



Open Format Based on Parquet

ACID Transactions

Apache Spark API's

# Data reliability challenges with data lakes



**Lack of schema enforcement** creates inconsistent and low quality data

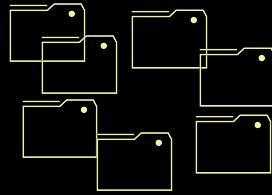


**Failed production jobs** leave data in corrupt state requiring tedious recovery

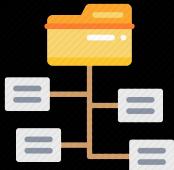


**Lack of consistency** makes it almost impossible to mix appends and reads, batch and streaming

# Performance challenges with data lakes



**Too many small or very big files** - more time opening & closing files rather than reading contents (worse with streaming).



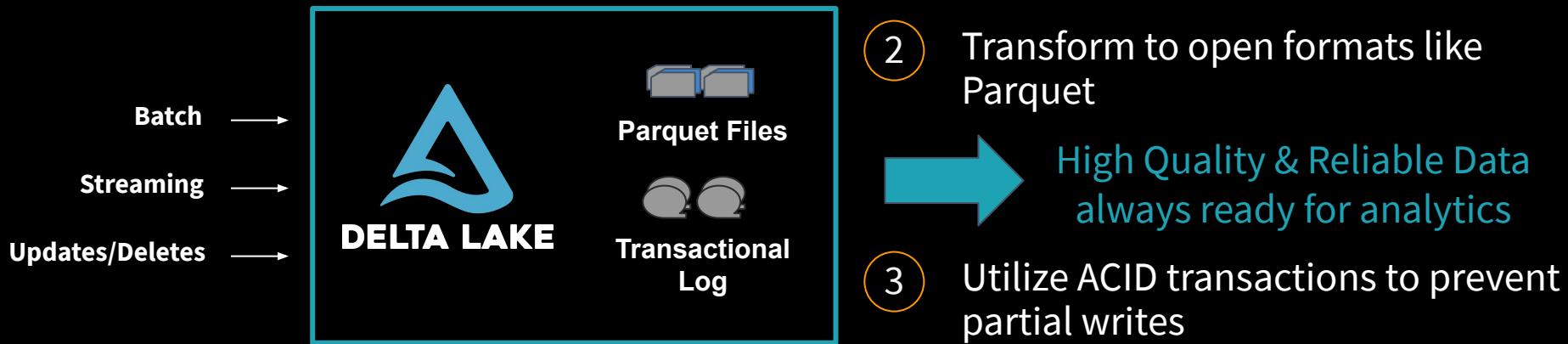
**Partitioning aka “poor man’s indexing”**- breaks down if you picked the wrong fields or when data has many dimensions, high cardinality columns.



**No caching** - storage throughput is low (storage is 20-50MB/s/core vs 300MB/s/core for local SSDs).

# Delta Lake ensures data reliability

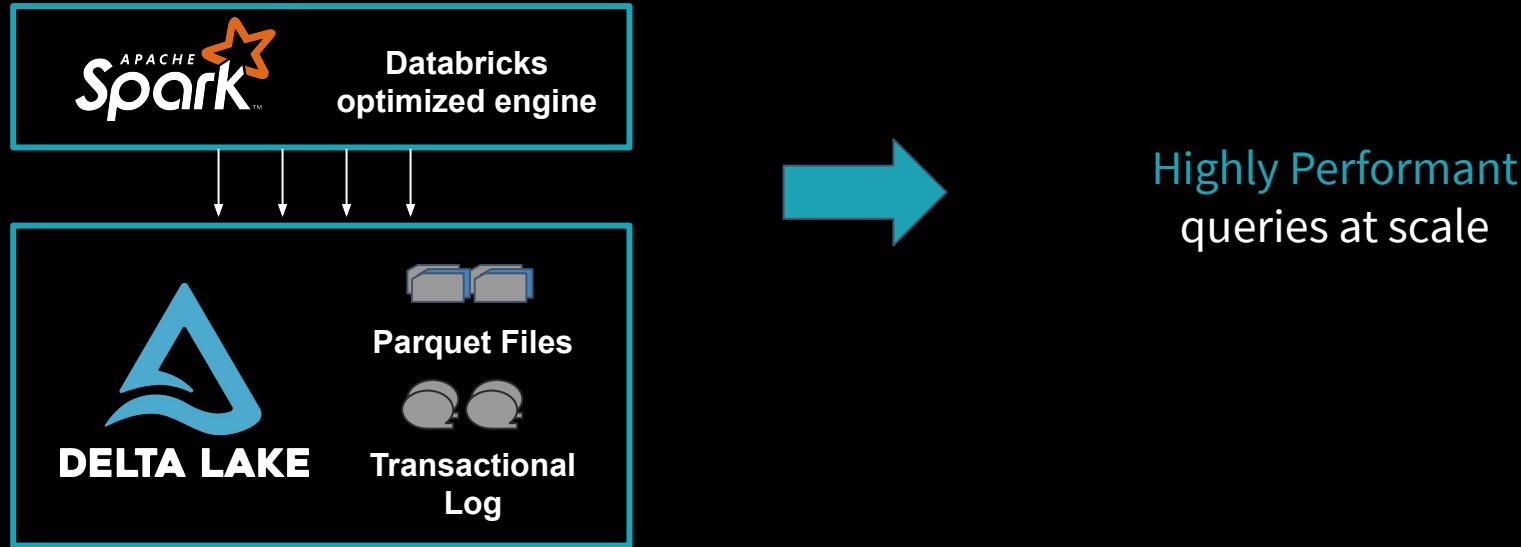
- 1 Manage your ingest process - handle batch and streaming data



## Key Features

- ACID Transactions
- Schema Enforcement
- Unified Batch & Streaming
- Time Travel/Data Snapshots

# Delta Lake optimizes performance



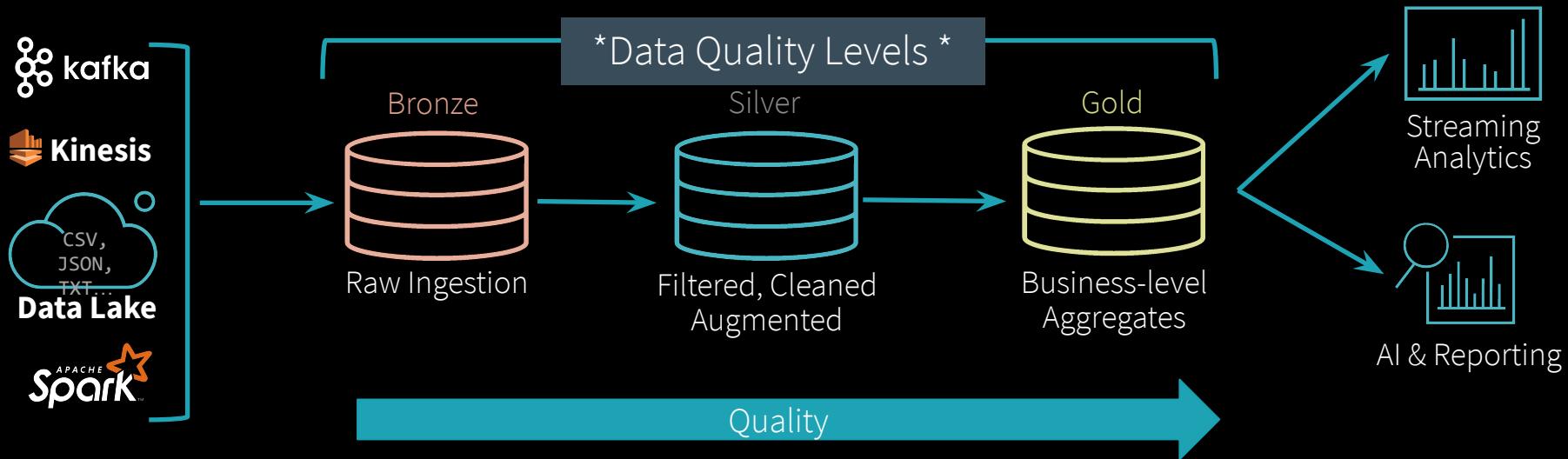
## Key Features

- Indexing
- Compaction
- Data skipping
- Caching

# The DELTA LAKE

4

Curate your data and refine it in a series of steps



Delta Lake allows you to *incrementally* improve the quality of your data until it is ready for consumption.

# Delta gives you the features you need

The features you expect on a database

- Transactions
- Recovery
- Snapshots
- Indexes

With the benefits of a data lake:

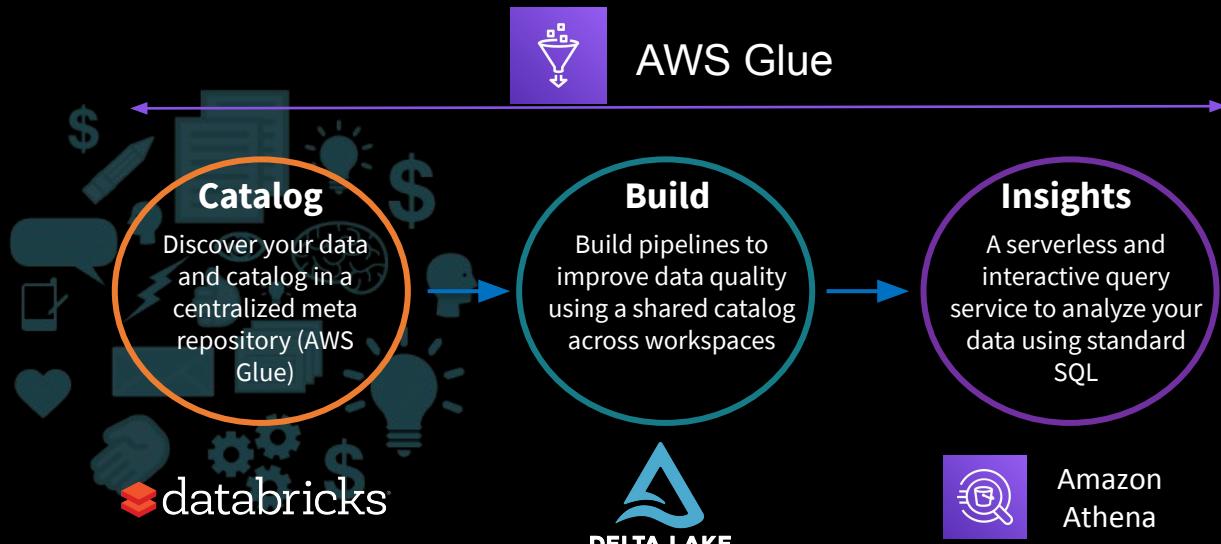
- No schema, flexible schema or fixed schema
- Integrated stream and batch processing
- Open format so your data is not locked and subject to vendor tax.



# Integration with Glue



## 5 Setup a data catalog for your data



### Glue integration

- Attach an IAM role to a Databricks cluster with the proper Glue Metastore permission
- Set the Spark Configuration, prior to launch the cluster
- Manage table definitions
- Simplifies manageability by using the same AWS Glue catalog across multiple workspace

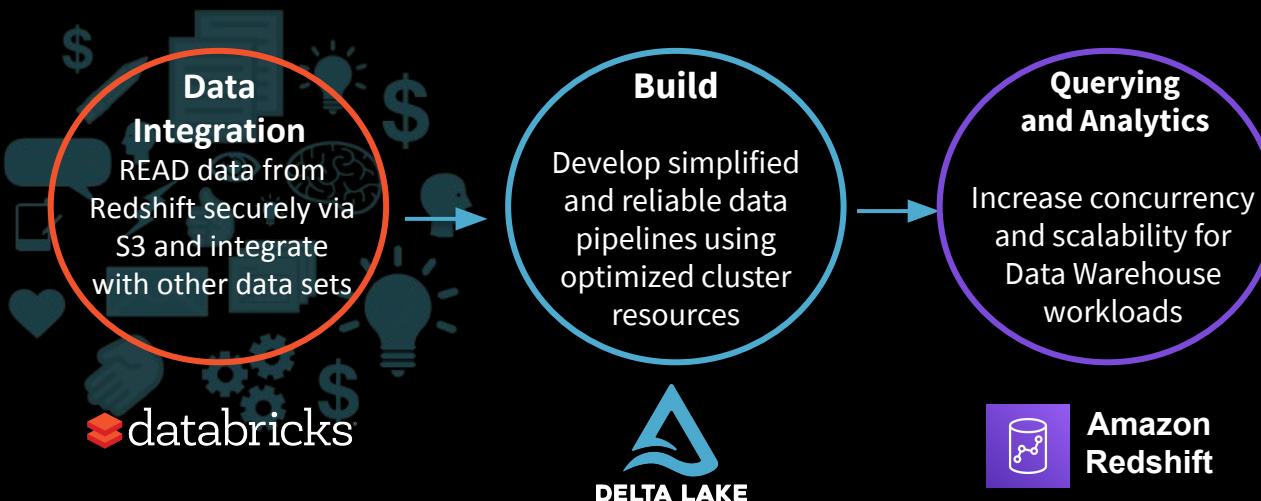
### Athena integration

- Create a manifest file and an external table definition
- Query the external table using Athena

# Integration with Redshift



## 7 Distribute results to LOB units for action



### Redshift integration

- Download and install the official Redshift JDBC driver on the cluster
- Authenticate to S3 and Redshift using IAM role
- S3 acts as an intermediary to store bulk data when reading from or writing to Redshift
- Secure your JDBC connection and encrypt the UNLOAD and COPY data operations following the Redshift documentation

# Cloud native enterprise solution

## 6 Define centralized security, governance and audit procedures

Consistent security  
and roles enable easy  
enterprise roll-outs



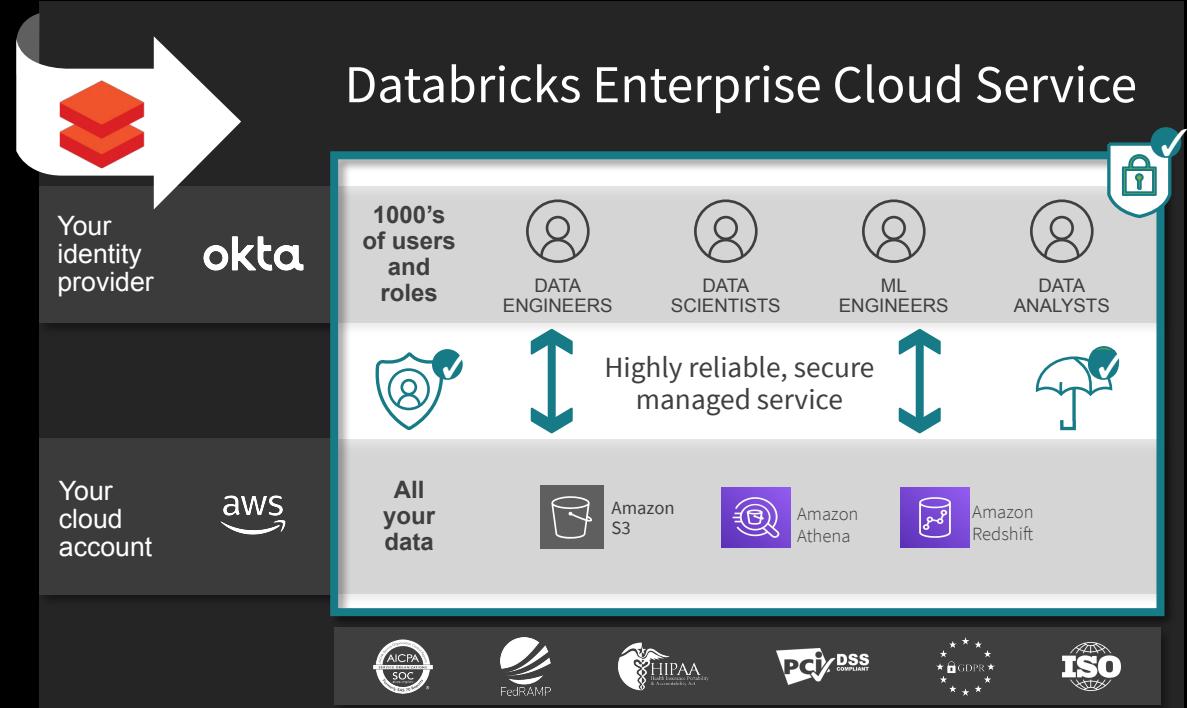
Adopt roles from SSO pass through



Consolidated security



Consistent governance



# Best practices for building a cloud data lake

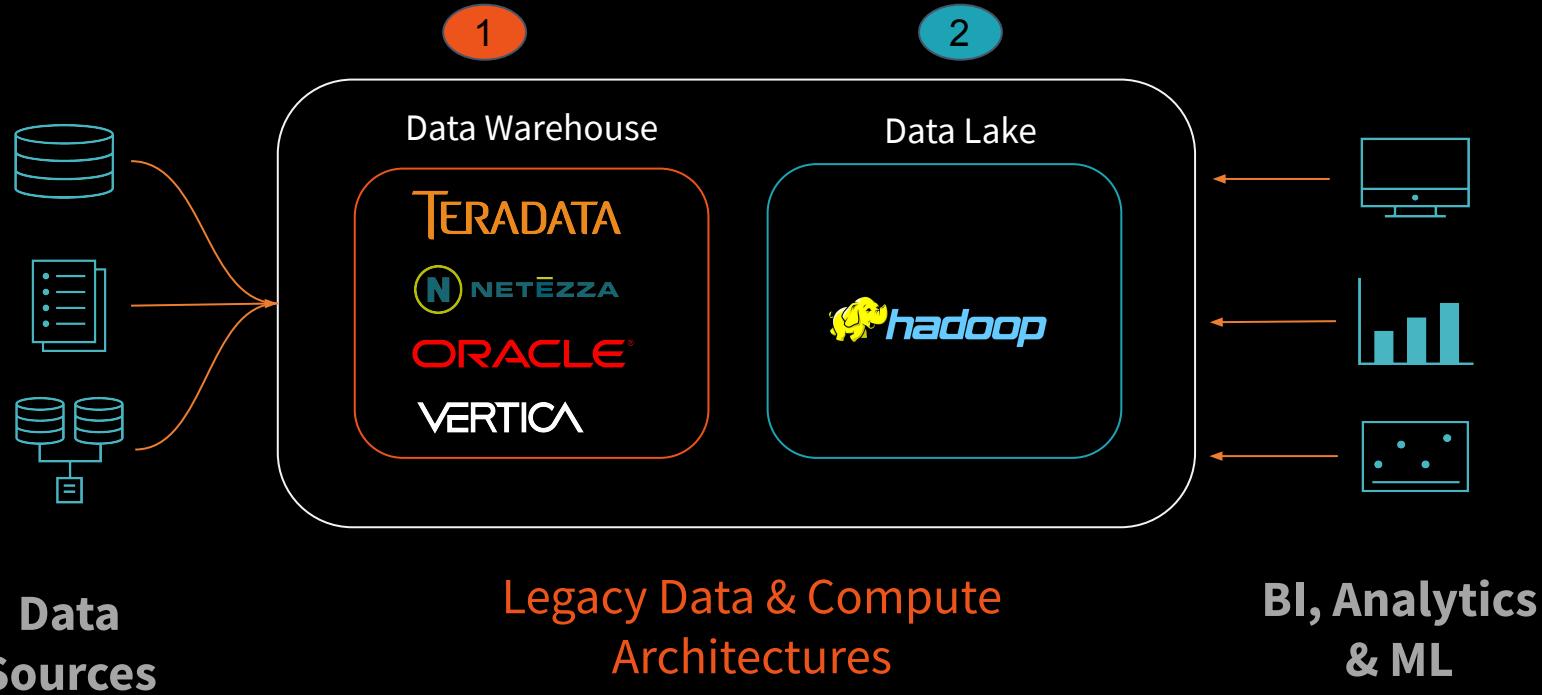
- 1 Manage your ingest process - handle batch and streaming data [using Delta Lake](#)
- 2 Transform to open formats like Parquet [using Delta Lake](#)
- 3 Utilize ACID transactions to prevent partial writes [using Delta Lake](#)
- 4 Curate your data and refine it in a series of steps [using Delta Lake](#)
- 5 Setup a data catalog for your data [using Glue](#)
- 6 Define centralized security, governance and audit procedures [using AWS security](#)
- 7 Distribute results to LOB units for action [using Delta Lake or Redshift](#)



RAW DATA LAKE

# On-premises analytics to cloud data lakes

## Common Patterns



# Hadoop: Complex, costly, risk-prone

1

Data reliability &  
performance  
Issues

Hung jobs, corrupted-data

Inconsistent schemas

No roll-backs

2

Rigid, In-flexible  
clusters - cannot  
scale up/ down

Unresponsive to business needs

Cluster configurations rigid

Hardware procurement cycles

3

Devops  
Nightmare; Low  
productivity

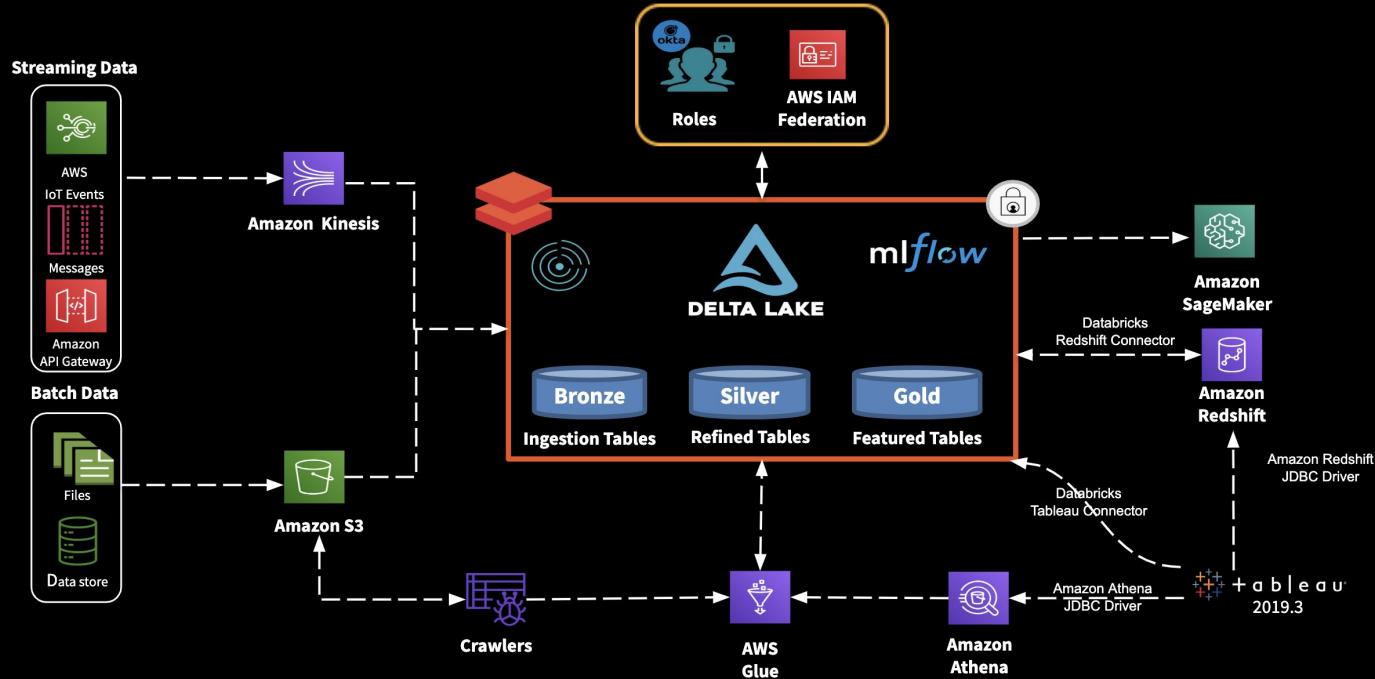
Silo'd AI/ Big Data environments

Shared resource conflicts

Upgrades, Maintenance delays

# AWS data lake implementation

## The Databricks Unified Data Analytics Platform





# Kabbage migration

*Lower costs, more capabilities*

	<b>Data Platform 2.0</b>	<b>Data Platform 3.0</b>
Technology Stack - Major Components	MapR, Sql Server	Athena, Databricks, S3, Akka, Aurora
Technology Stack - Minor Components	Influx, Kafka, EMR, Druid, Nirvana, StreamSets, DCOS, Cloudera, Arrango, Vertica, S3	Influx, Kafka, EMR, Sql Server, Kubernetes
Cloud Environment	Hybrid	100% AWS
Maintainability	Labor Intensive	Auto Scalable, Auto Restart
Ease of Data Import/Export	Very Hard	Easy
Yearly Cost (approximate)	\$2,000,000	\$1,000,000 - \$1,400,000

# Data lake pipeline modernization case study

AWS ML workloads managed by Databricks since 2017

Data lake modernization project Q4 2019: Modernized native S3 data lake with Databricks Delta

## Before: 72 person-days

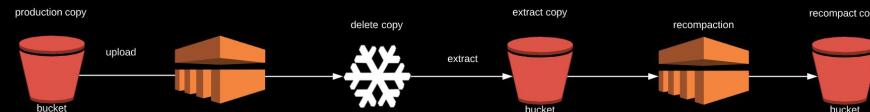
- Multiple (4) Copies of data; operational overhead, low data fidelity
- Manual Recompaction - operational overhead, costs
- Total = \$397,000 / Quarter (incl \$180K Snowflake + \$160K EMR costs/quarter)

## After: 6 person days

- No copies. Deletes occur in-place.
- New Cost Estimate \$23,000 per quarter

**Total Yearly Savings:  
\$1,496,000**

Before: Data lake architecture



# Databricks' migration solution

## On-premises



Other EDWs

Data  
Metastore  
Apps/ Jobs/ Queries  
Security Governance

## Migration

Reduce risk and accelerate outcomes  
Automate migration data + code  
Proven tools and methodologies  
Working with your preferred partners

## Databricks



## ISV Partners

**Migration:** WanDisco, StreamSets, Talend, Infoworks

**Security:** Privacera, Immuta  
**CDC:** Attunity, Syncsort

## Program for Hadoop Migration

- Assessment
- Data replication
- Frameworks

[www.databricks.com/migration](http://www.databricks.com/migration)

## Consulting & SI Partners

KnowledgeLens, Slalom, Clarity Insights, CapGemini, MindTree, Cognizant, TCS

# Now available on AWS Marketplace

- Simplified licensing
- One billing for Databricks usage

<http://bit.ly/DATABRICKS>

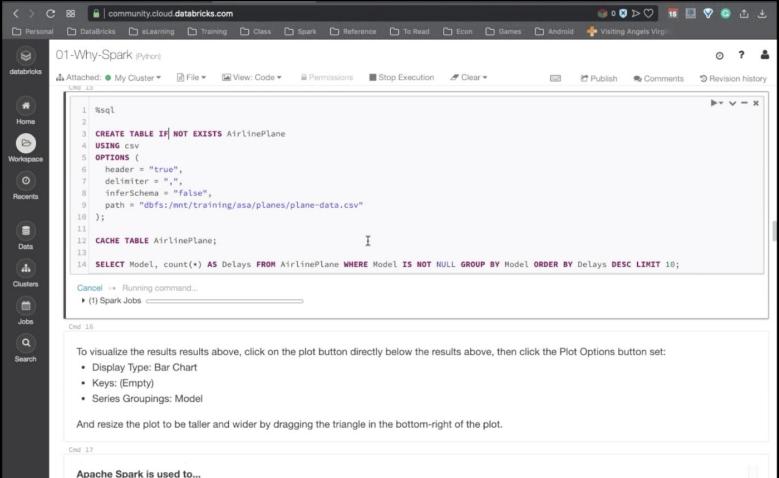
The screenshot shows the AWS Marketplace interface. At the top, there's a search bar and navigation links for Categories, Delivery Methods, Solutions, Migration Mapping Assistant, Your Saved List, Partners, and Sell in AWS Marketplace. Below the header, the Databricks logo is displayed next to the product name "Databricks Unified Analytics Platform". It's sold by "Databricks Inc.". A brief description states: "Databricks Unified Analytics Platform is a cloud-based service for running your analytics in one place - from highly reliable and performant data pipelines to state-of-the-art machine learning. From the original creators of Apache Spark and MLLib, it provides data science and engineering teams ready-to-use clusters with optimized Apache Spark and various ML frameworks (e.g., TensorFlow) coupled with powerful collaboration capabilities to improve productivity across the ML lifecycle." There are "Show more" and "Continue to Subscribe" buttons. The main content area has tabs for Overview (which is selected), Pricing, Usage, Support, and Reviews. The "Product Overview" section contains detailed information about the platform's capabilities, such as its use of DBUs and its support for various languages and frameworks. A "Highlights" sidebar lists several key features: reliable infrastructure, enterprise-grade security, buildable data pipelines, outstanding collaboration, and ML lifecycle support.

# Free Training

Three virtual 2-hour training sessions using Databricks on AWS:

- Getting Started with Apache Spark
- Data Engineering and Streaming Analytics
- Machine Learning

Sign up now:  
<http://bit.ly/TrainingAWS>



```
1 ssql
2
3 CREATE TABLE IF NOT EXISTS AirlinePlane
4 USING csv
5 OPTIONS (
6   header = "true",
7   inferSchema = "true",
8   inferschema = "false",
9   path = "dbfs:/mnt/training/asa/planes/plane-data.csv"
10 );
11
12 CATCH TABLE AirlinePlane;
13
14 SELECT Model, count(*) AS Delays FROM AirlinePlane WHERE Model IS NOT NULL GROUP BY Model ORDER BY Delays DESC LIMIT 10;
```

To visualize the results results above, click on the plot button directly below the results above, then click the Plot Options button set:

- Display Type: Bar Chart
- Keys: (Empty)
- Series Groupings: Model

And resize the plot to be taller and wider by dragging the triangle in the bottom-right of the plot.

Cmd 16

Cancel → Running command...  
► (1) Spark Jobs

Cmd 17

Apache Spark is used to...

# Customer Profile Slides

# VIACOM

## Use Cases: Quality of Service (QoS)

Across their 4+ billion subscribers, Viacom delivers tailored experiences and uses real-time network insights to optimize quality of service during live events like the MTV Music Awards.

## Why Databricks:

- Scalable platform that can handle 1.2PBs of daily data
- Single engine for processing ad hoc, batch, and real-time data, so data teams only need to monitor one system
- Data can be pushed to Tableau to support BI users

## Impact:

- **33% reduction** in video delays by predicting network issues
- **7x increase** in views per session
- Improved audience engagement and retention

A large black granite plaque on a city building features the white FINRA logo. The logo consists of the word "FINRA" in a bold, sans-serif font, with a stylized "V" shape composed of a grid of small triangles extending from the letter "R".

Analyzes 100 billion stock market events per day to identify fraud and wrong doing.

**With Databricks we have created an ecosystem that allows data scientists self-service access to data, a direct path from prototyping to production and access to a common data lake from a suite of analytic technologies**





“Our early large-scale data strategy has already paid off with a new drug target for chronic liver disease.”

#### Saving Lives by Scaling Data Science

- **Data Science on** 250K genomes and growing to 500K genomes (10TB)
- **Increase model effectiveness by reducing** full data set query runtimes from 30 minutes to 5 seconds
- **Increased productivity of** Bioinformaticians and clinical researchers by automating devops





**Construction Industry:** 98% of projects over on cost or time. And annual productivity has increased by 1% over the last 20 years

## Optimizing Construction Sequencing with TensorFlow and Apache Spark

- Millions of parts in 3-dimensional space
- Billions in budget over 10-15 years long projects
- Potential savings of 10's to 100's of millions

# SPARK+AI SUMMIT 2020

A large, diverse crowd of people is shown from the waist up, sitting in rows in a dark auditorium. They are all looking towards the left side of the frame, presumably at a stage or presentation area. The lighting is low, with a blue tint, creating a professional and focused atmosphere.

JUNE 22-25, 2020 | SAN FRANCISCO

ORGANIZED BY  databricks

**SAVE \$450  
REGISTER BY MARCH 31**

**SAVE YOUR SPOT**

**EXPANDED  
TECHNICAL TRAINING**

**LEARN MORE**

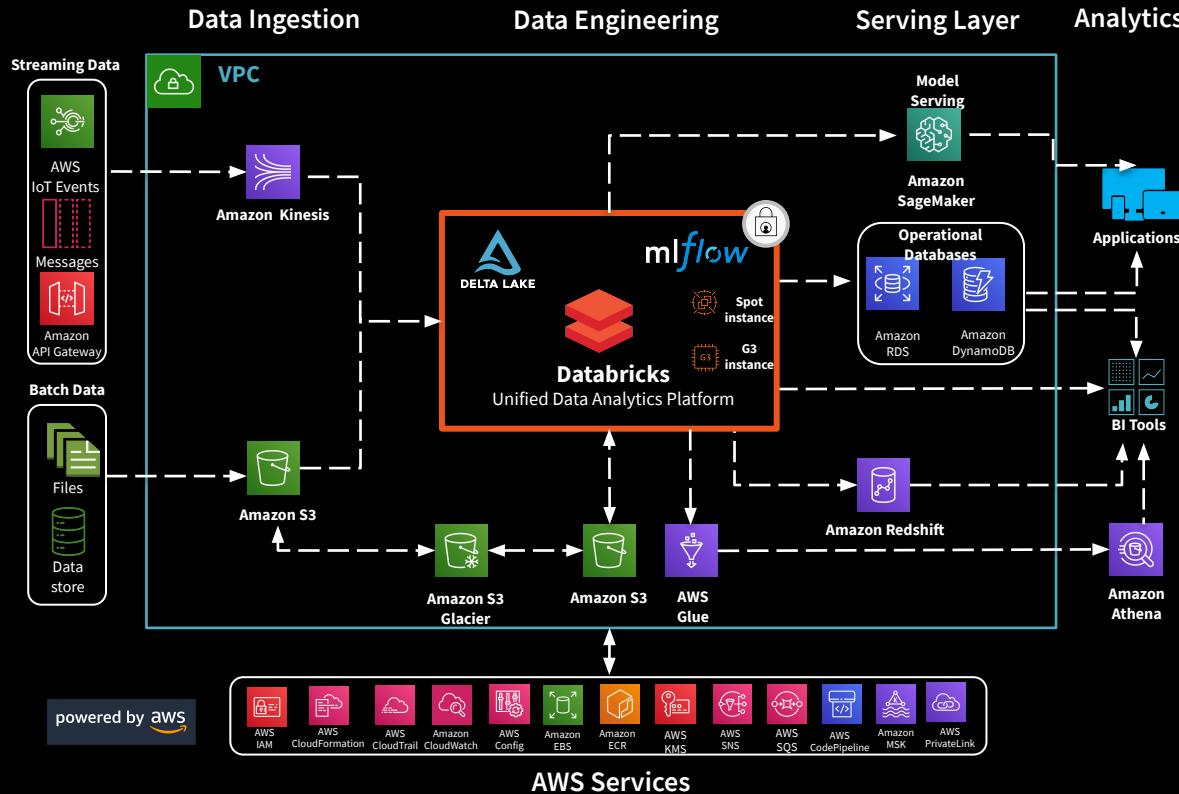
# Agenda

<b>9:00-9:30 AM</b>	Opening Remarks - Enabling Cloud Data Lakes for Analytics
<b>9:30-9:45 AM</b>	Customer Use Cases
<b>9:45-9:55 AM</b>	Break
<b>9:55-10:40 AM</b>	Architecture - Incorporating Roles; Scaling, Launching and Managing Clusters; Managing Pools and Containers
<b>10:40-11:25 AM</b>	Notebooks - Data Loading using Scala, Machine Learning Analytics using Python, and Redshift integration for data delivery to analysts
<b>11:25-11:35 PM</b>	Q&A



# Databricks Unified Data Analytics Technical Interactive Demo

# The Databricks Unified Data Analytics Platform on AWS

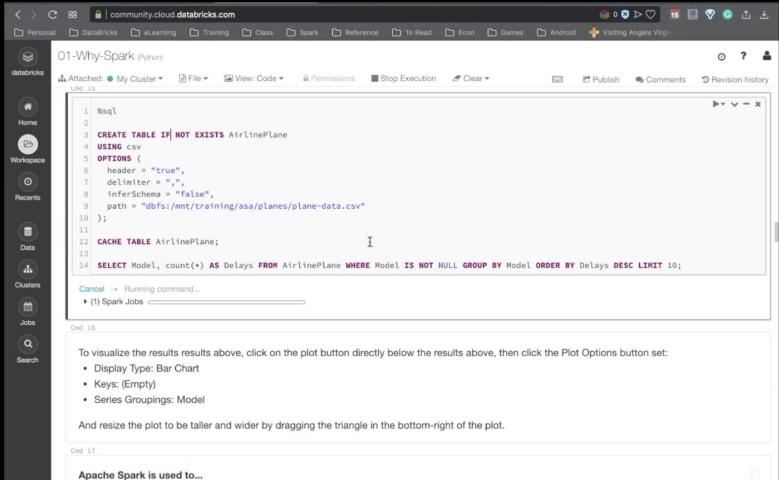


# Free Training

Three virtual 2-hour training sessions using Databricks on AWS:

- Getting Started with Apache Spark
- Data Engineering and Streaming Analytics
- Machine Learning

Sign up now:  
<http://bit.ly/TrainingAWS>



```
1 #sql
2
3 CREATE TABLE IF NOT EXISTS AirlinePlane
4 USING csv
5 OPTIONS (
6   header = "true",
7   inferSchema = "true",
8   inferschema = "true",
9   path = "dbfs:/mnt/training/asa/planes/plane-data.csv"
10 );
11
12 CATCH TABLE AirlinePlane;
13
14 SELECT Model, count(*) AS Delays FROM AirlinePlane WHERE Model IS NOT NULL GROUP BY Model ORDER BY Delays DESC LIMIT 10;
```

To visualize the results results above, click on the plot button directly below the results above, then click the Plot Options button set:

- Display Type: Bar Chart
- Keys: (Empty)
- Series Groupings: Model

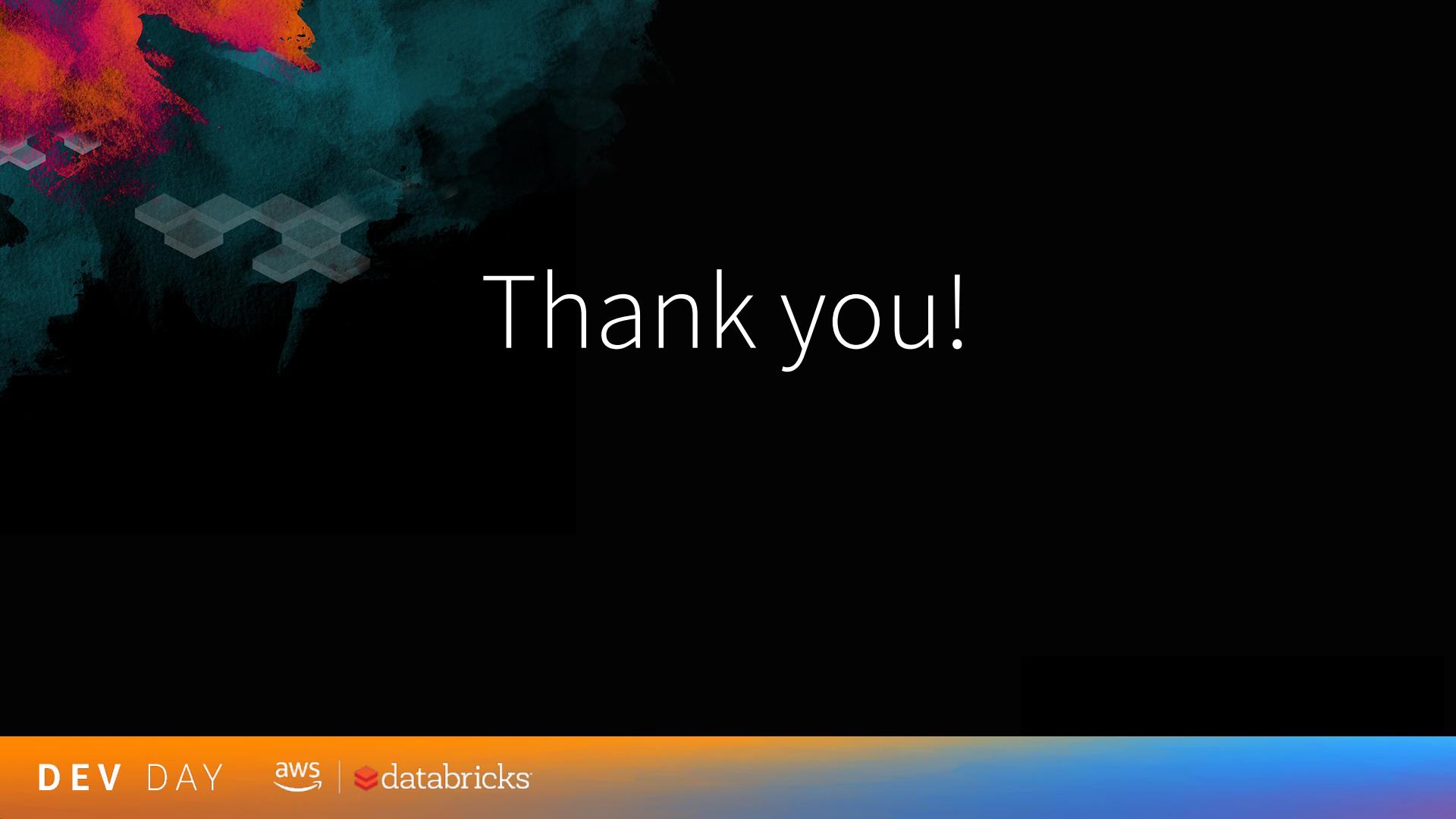
And resize the plot to be taller and wider by dragging the triangle in the bottom-right of the plot.

Cmd 16

Cancel → Running command...  
► (1) Spark Jobs

Cmd 17

Apache Spark is used to...



# Thank you!