



Unifying Data Warehousing with Data Lakes

Ali Ghodsi, Co-Founder & CEO

Oct 25, 2017



Many enterprises are undergoing
a data transformation

Databricks Customers Across Industries

Financial Services



JPMORGAN
CHASE & CO.



Healthcare & Pharma



Media & Entertainment



Data & Analytics Services



Technology



Public Sector



Retail & CPG



Consumer Services



Marketing & AdTech



Energy & Industrial IoT



Health care AI cloud dataset use case

Financial Services



Healthcare & Pharma



Media & Entertainment



Data & Analytics Services



Technology



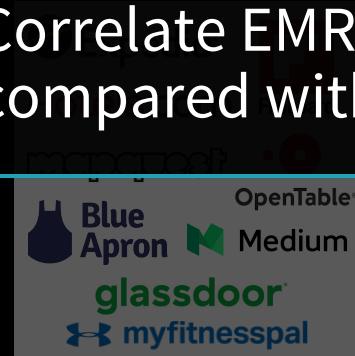
Public Sector



Retail & CPG



Consumer Services



Marketing & Ad Tech



Energy & Industrial IoT



Correlate EMR of 50,000 patients compared with their DNA

Enterprise AI use case

Financial Services



Public Sector

Healthcare & Pharma



Media & Entertainment



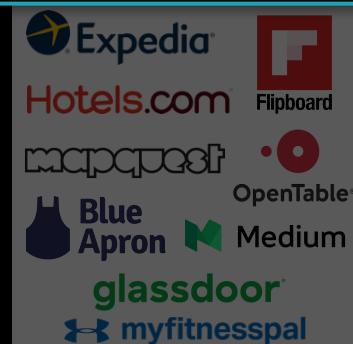
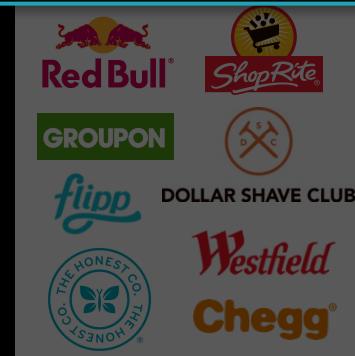
Data & Analytics Services



Technology



Provide recommendations to sales
using NLP and deep learning



Real-time AI use-case

Financial Services



JPMORGAN
CHASE & CO.

Nasdaq

Public Sector



Healthcare & Pharma



Retail & CPG



Media & Entertainment



Data & Analytics Services



Curb abusive behavior
across gamers globally

Technology



Big Data was the Missing Link for AI

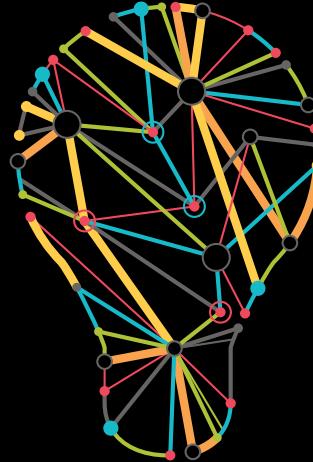
BIG DATA



Customer Data
Emails/Web pages
Click Streams
Sensor data (IoT)
Video/Speech
...



GREAT RESULTS



Most companies are Struggling with Big Data

Hardest part of AI isn't AI

“Hidden Technical Debt in Machine Learning Systems”, Google NIPS 2015

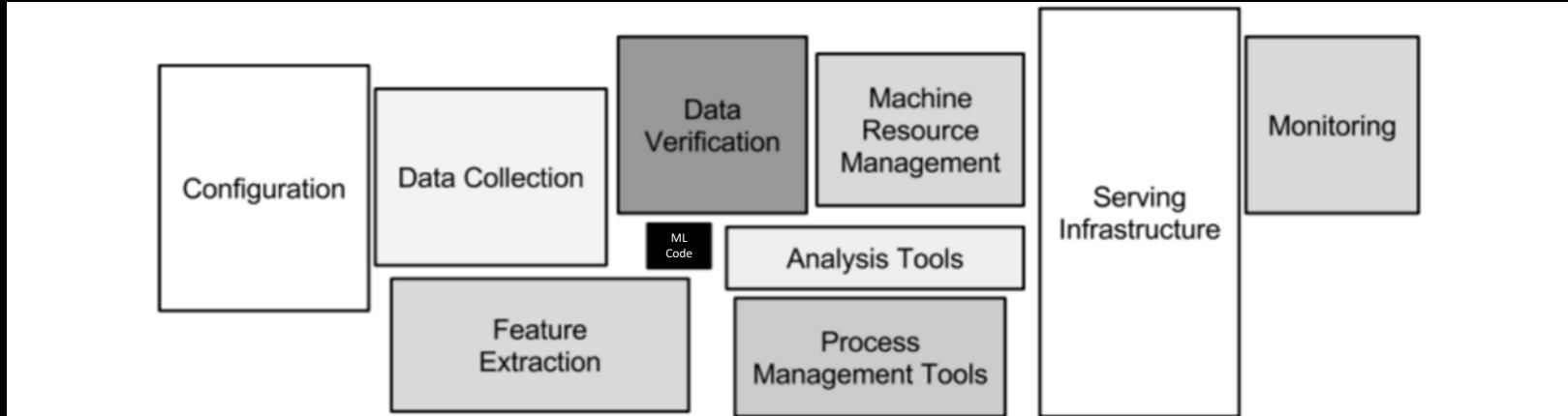


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

The hardest part of AI is Big Data

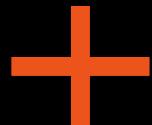
Building Predictive Applications
is really Hard!



databricks®

Unified Analytics Platform

UNIFIED
PROCESSING
ENGINE



UNIFIED
EXPERIENCE
ACROSS TEAMS

The Evolution of Big Data

The Era of the Data Warehouse

Data Warehouse (DW)

ETL important data to central DW and get Business Intelligence (BI)

THE GOOD

- Pristine Data
- Fast Queries
- Transactional

THE BAD

- Expensive to Scale, not Elastic
- Requires ETL, Stale Data, No Real-Time
- No Predictions, No ML
- Closed formats (lock in)

Not Future Proof – Missing Predictions, Real-time, Scale

The Era of the Data Lake

Hadoop Data Lake

ETL all data to central scalable open lake for all use cases

THE GOOD

- Massive scale
- Inexpensive Storage
- Open Formats (Parquet, ORC)
- Promise of ML & Real Time Streaming

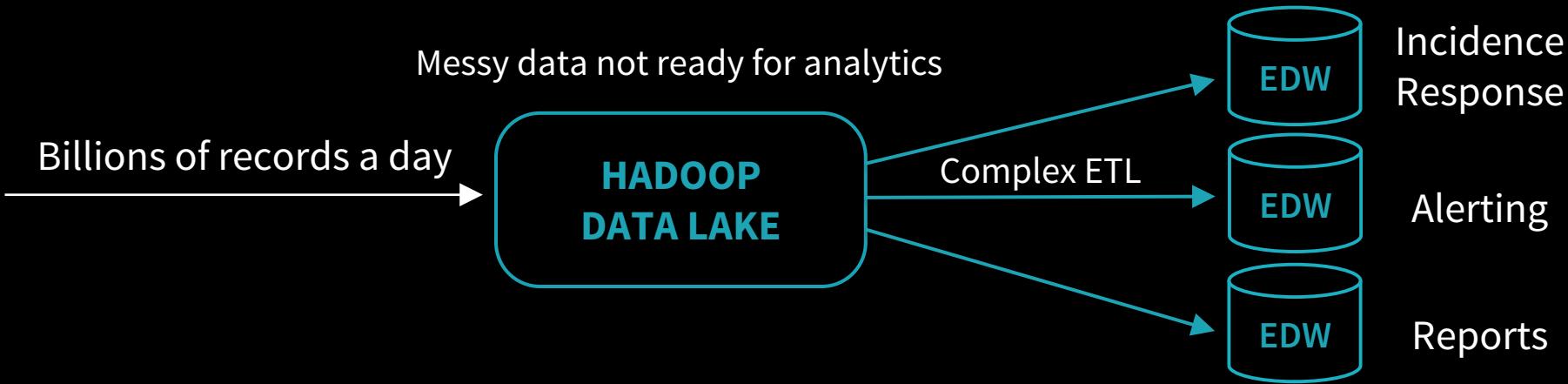
THE BAD

- Inconsistent Data
- Unreliable for Analytics
- Lack of Schema
- Poor Performance

Become a cheap messy data store with poor performance

The Current State of Data Platforms

Info Sec at a Fortune 100 Company



ENTERPRISE DATA WAREHOUSE

- Only 2 weeks of data
- Very expensive to scale
- Proprietary Formats
- No Predictions (ML)

DISADVANTAGES OF ARCHITECTURE

- Poor agility in responding to new threats
- Scale Limitations, no historical data
- 6 Months and twenty people to build

The Next Generation Data Platform

Announcing Databricks Delta

First **UNIFIED** data management system that delivers:



The
SCALE
of data lake

The
RELIABILITY & PERFORMANCE
of data warehouse

The
LOW-LATENCY
of streaming



The
SCALE
of data lake

The
**RELIABILITY &
PERFORMANCE**
of data warehouse

The
LOW-LATENCY
of streaming

Databricks Delta

THE GOOD OF DATA LAKES

- Massive scale on Amazon S3
- Open Formats (Parquet, ORC)
- Predictions (ML) & Real Time Streaming

THE GOOD OF DATA WAREHOUSES

- Pristine Data
- Transactional Reliability
- Fast Queries (10-100x)

Enables Predictions, Real-time and Ad Hoc Analytics at Massive Scale

Databricks Delta Under the Hood

MASSIVE SCALE

- Decouple Compute & Storage

RELIABILITY

- ACID Transactions & Data Validation

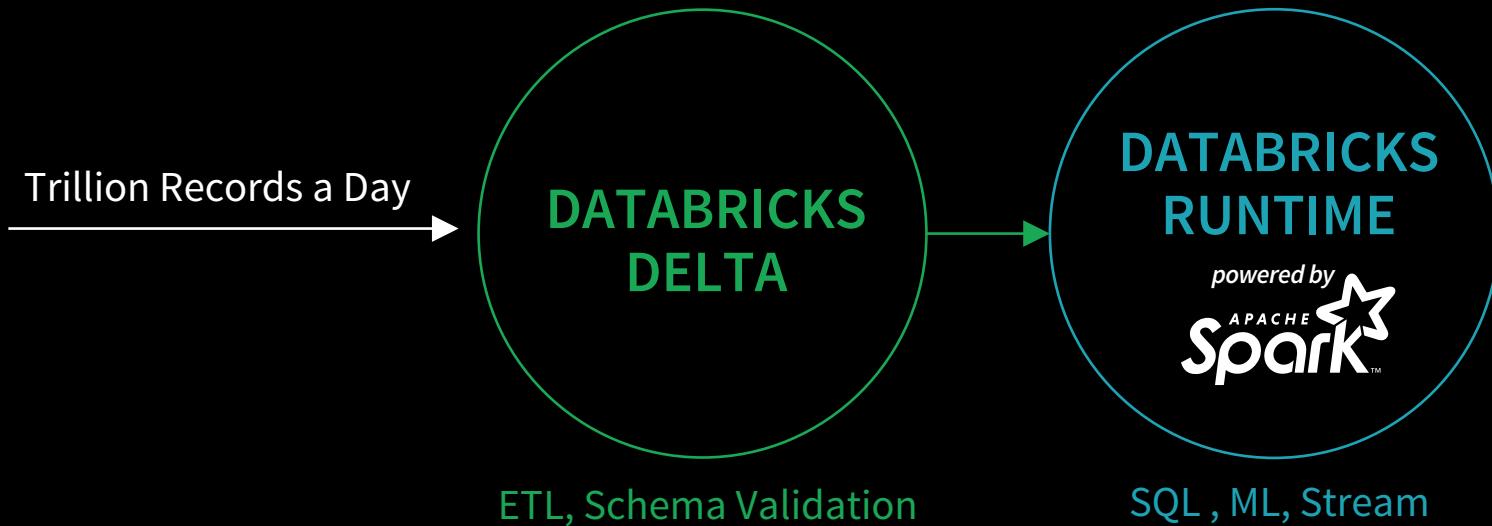
PERFORMANCE

- Data Indexing & Caching (10-100x)

LOW-LATENCY

- Real-Time Streaming Ingest

Info Sec with Databricks Delta



ADVANTAGES

- AI capable data warehouse at the scale of a data lake
- Interactive analysis on 2 years of data
- 2 Weeks to build with a 5 person data platform team



Unified Analytics Platform

UNIFIED
EXPERIENCE
ACROSS TEAMS

Notebooks, Dashboards, Reports



Unified Analytics Platform

UNIFIED

DATA

MANAGEMENT

Reliable Transactions, Performance

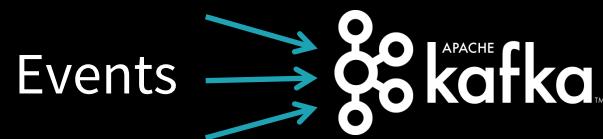


UNIFIED
EXPERIENCE
ACROSS TEAMS

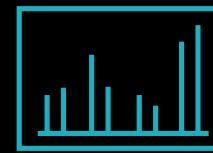
Notebooks, Dashboards, Reports

Demo by Michael Armbrust

Evolution of a Cutting-Edge Data Pipeline



Data Lake



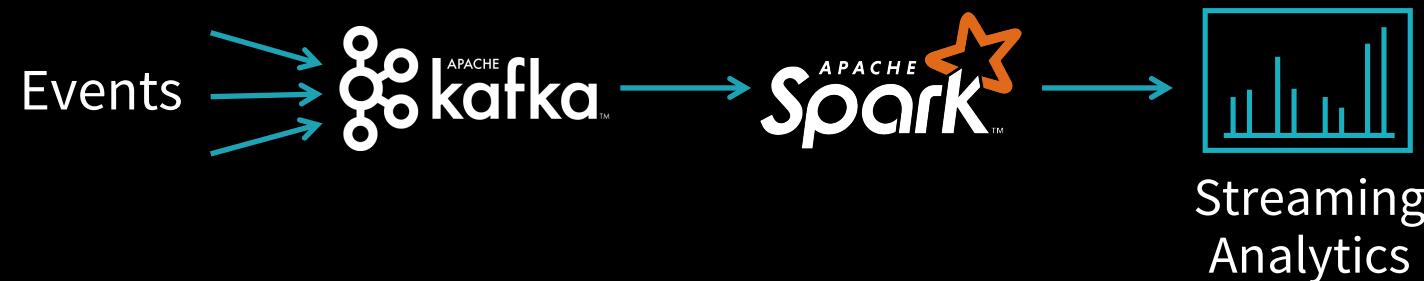
Streaming
Analytics

A graphic of a bar chart with several vertical bars of different heights, enclosed in a light blue square frame.

Reporting

A graphic of a bar chart with several vertical bars of different heights, enclosed in a light blue square frame, with a magnifying glass icon positioned over the top-left corner of the frame.

Evolution of a Cutting-Edge Data Pipeline

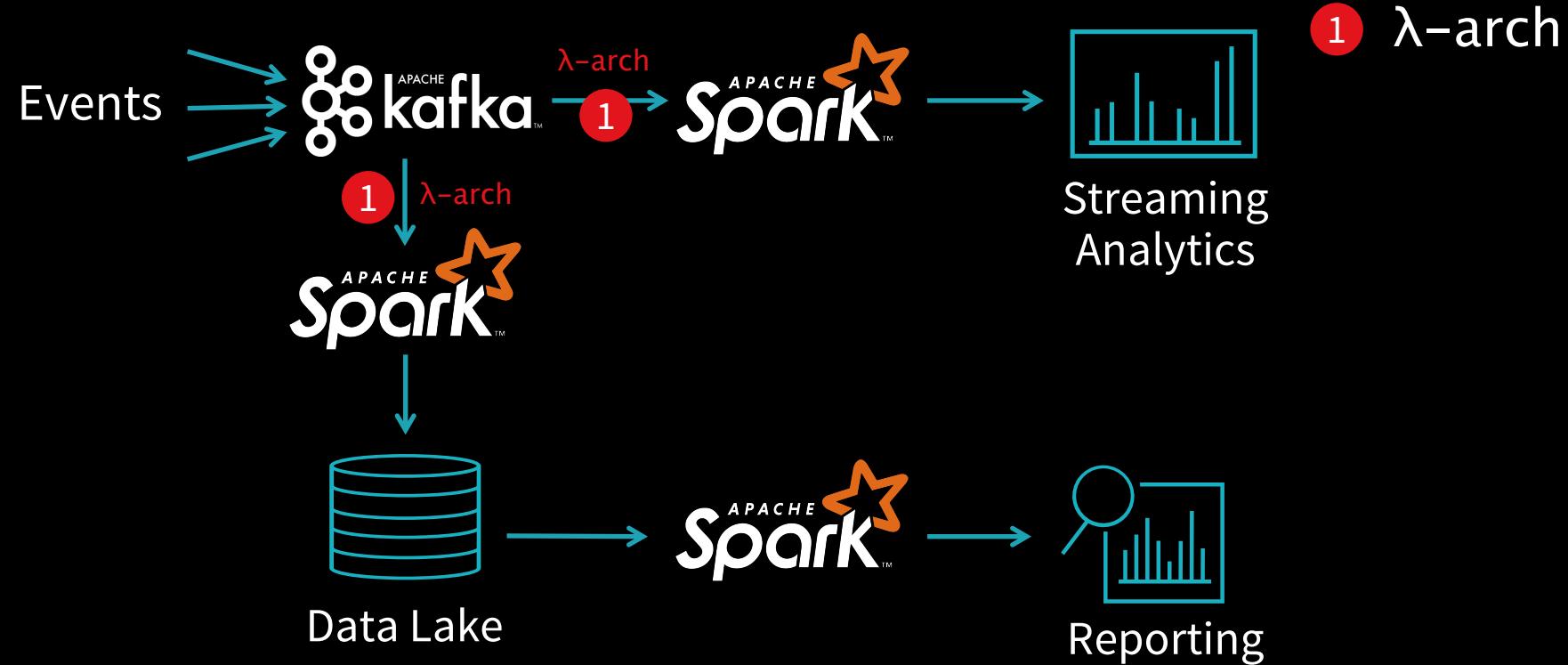


Data Lake

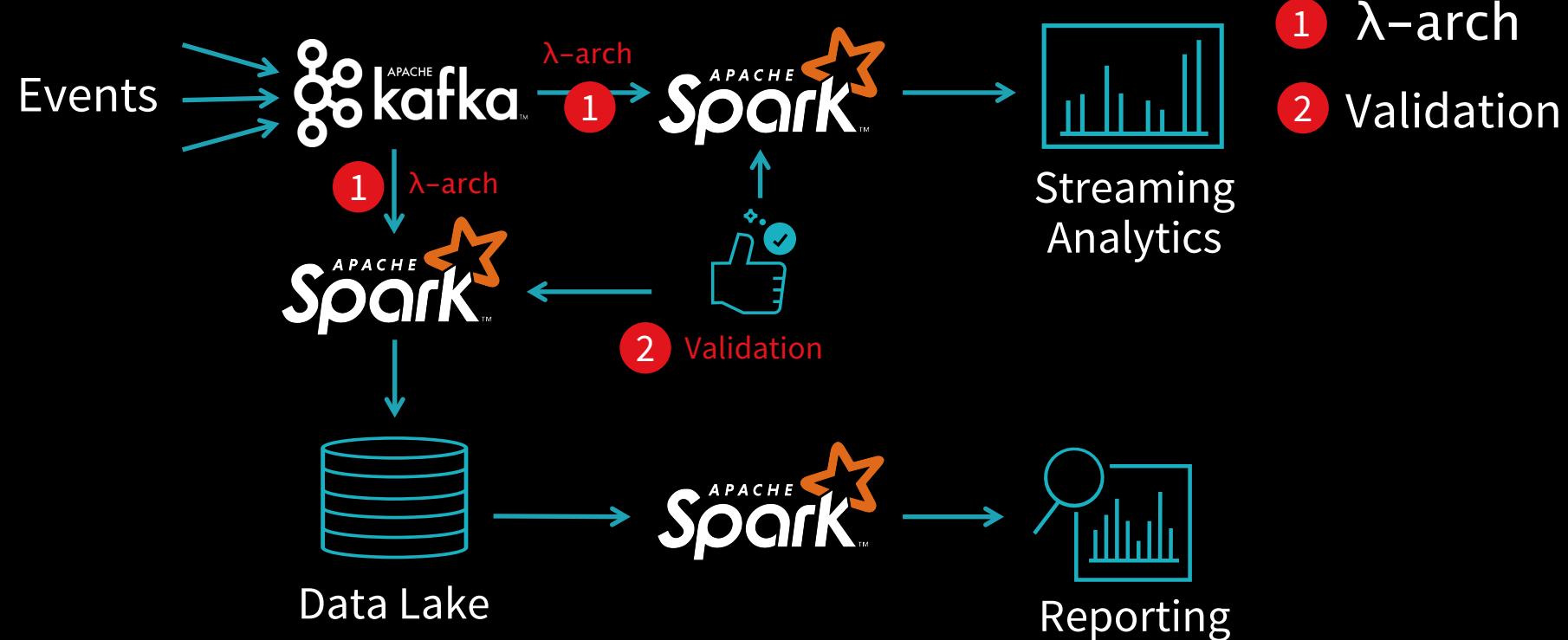


Reporting

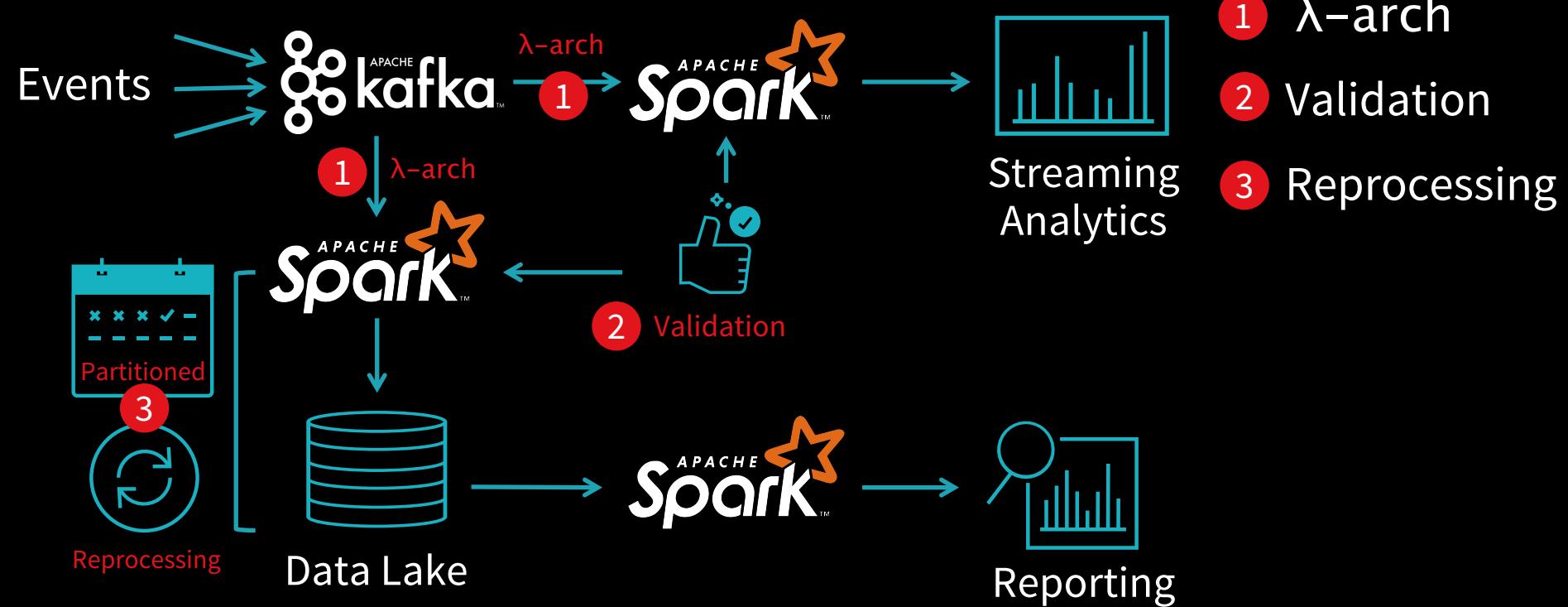
Challenge #1: Historical Queries?



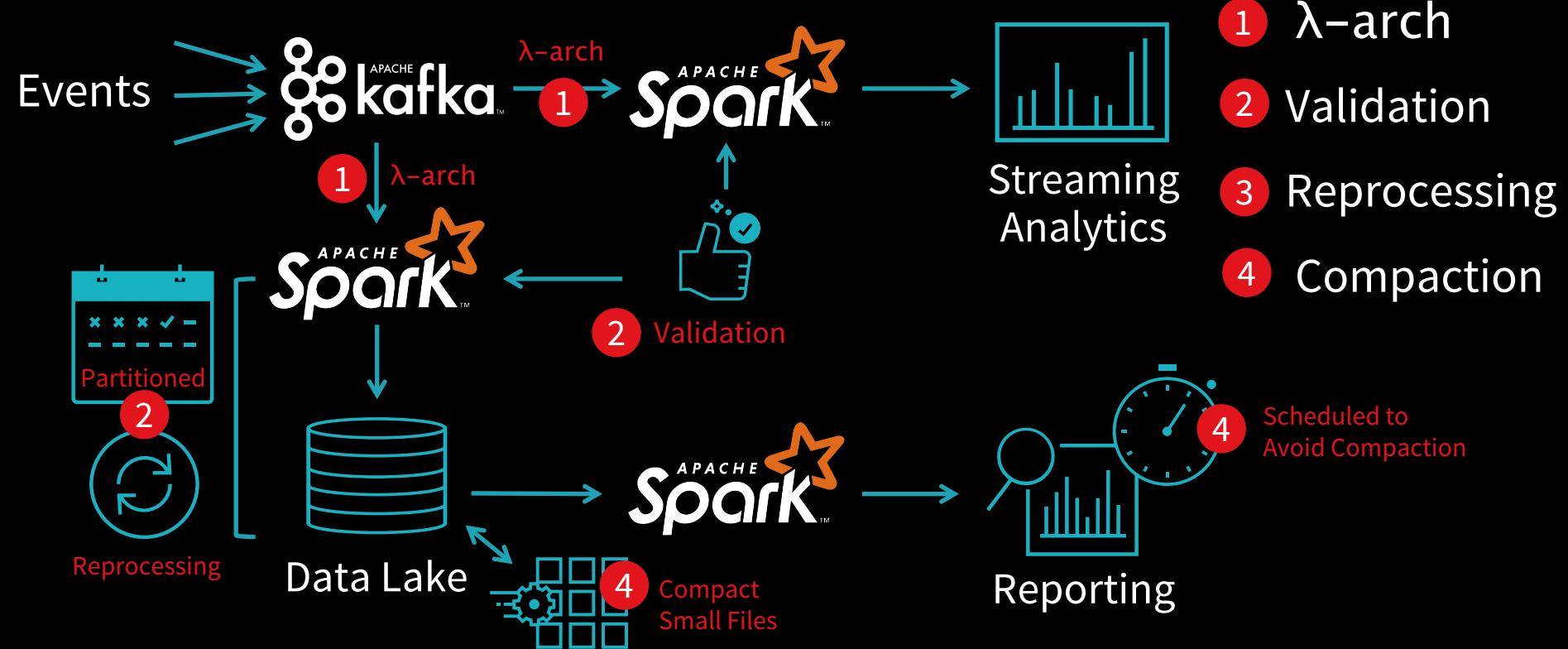
Challenge #2: Messy Data?



Challenge #3: Mistakes and Failures?

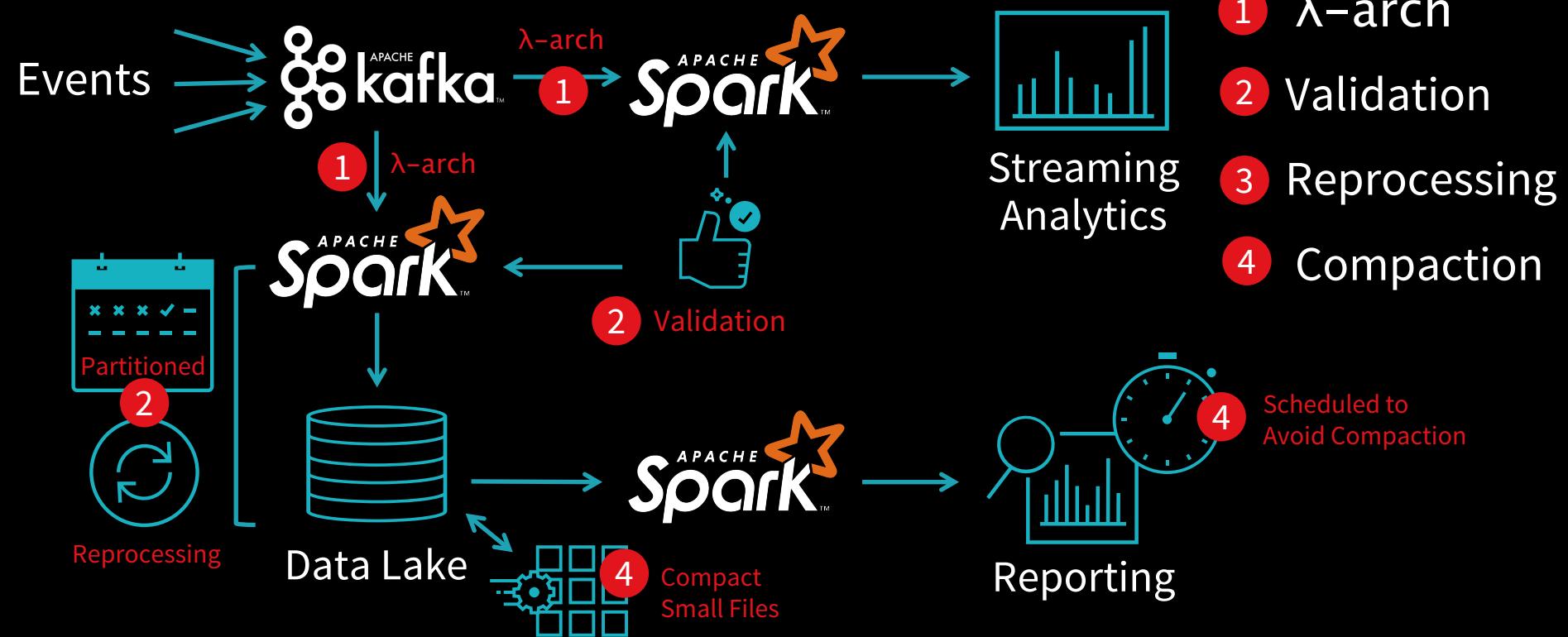


Challenge #4: Query Performance?

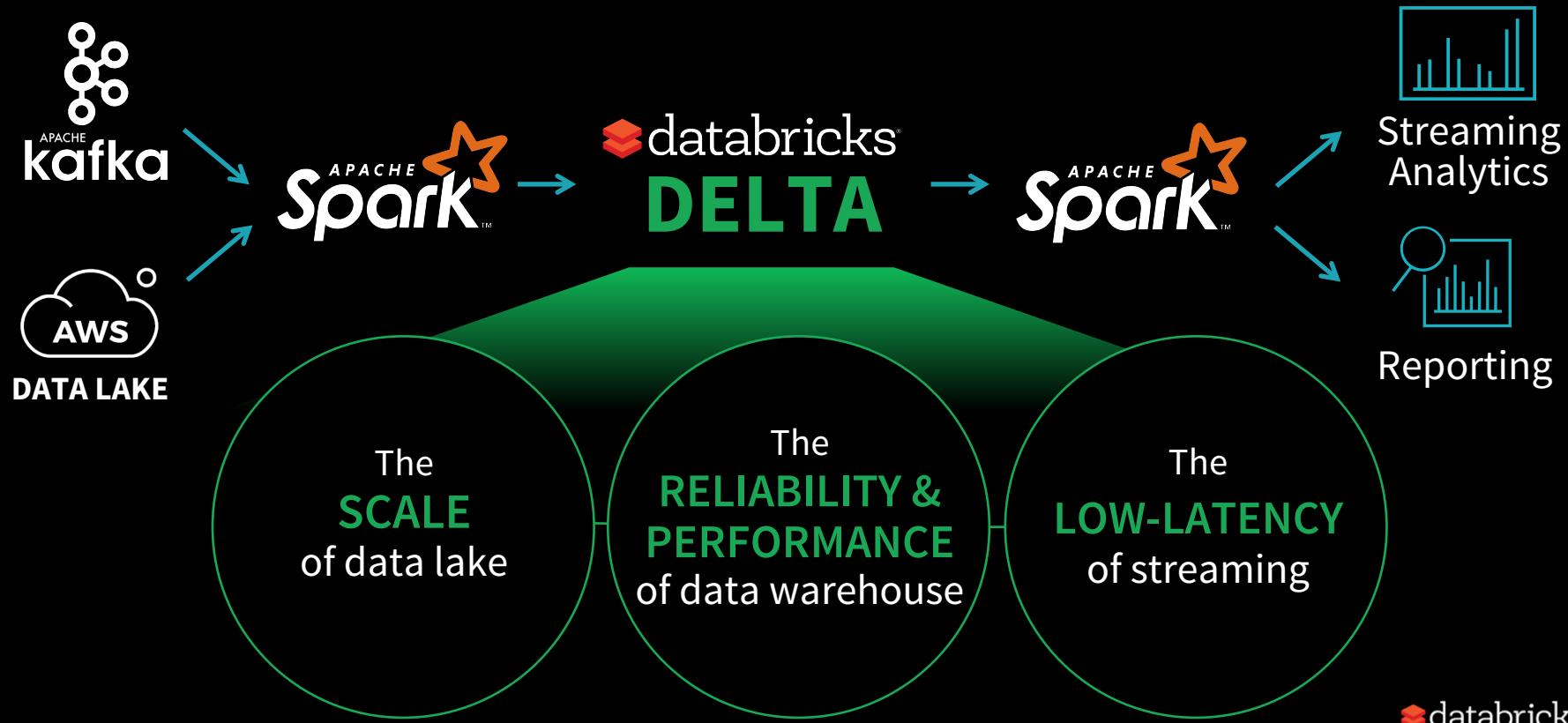


Let's try it instead with
DELTA

The Canonical Data Pipeline



The Delta Architecture



Sign up for the Private Beta
visit **databricks.**