# Building Better Data Pipelines for Apache Spark with Delta Lake

Vini Jaiswal

Customer Success Engineer

# Awesome TAs

- Emma Freeman (Curriculum Engineer)
- Spencer McGhin (Sr. Customer Success Engineer)
- Joe Widen (Lead Solutions Architect)
- Jean-Marc Spaggiari (Sr. RSA)

databricks

# Agenda

| Time | Topic |
|---|---|
| 9:00 - 9:20 | Introduction, Setup, Class Logistics |
| 9:20 - 9:50 | What is Delta Lake? |
| 9:50 - 10:00 | Break |
| 10:00 - 10:35 | Lab 1: Building a Delta Lake Data Pipeline |
| 10:35 - 10:50 | Review Lab, Q&A |
| 10:50 - 11:00 | Break |
| 11:00 - 11:15 | Delta Lake Features |
| 11:15 - 11:35 | Lab 2: Exploring Delta Lake Features |
| 11:35 - 11:55 | Summary, Q&A |

# Getting started:

- Login to the workspace
- Familiarize yourself with using chat to communicate with TAs
- Download materials from the Resources tab

# By the end of this class you will be able to:

- Explain how Delta Lake works for building data pipelines
- Build a complete Delta Lake pipeline to stream data from raw ingestion to gold tables
- Understand how Delta Lake unifies batch and streaming data in a single pipeline
- Understand how Delta Lake schema enforcement and evolution impacts data writes and updates
- Understand how to perform data deletions to comply with data privacy laws and regulations

# What is Delta Lake?

Delta Lake is a technology for **building robust data lakes**.
It is an **open source storage layer** specifically **designed
to work with Apache Spark**.

A data lake built using Delta Lake is **ACID compliant**.

# Delta Lake offers:

- ACID transactions on Spark
- Scalable metadata handling
- Streaming and batch unification
- Schema enforcement
- Time travel
- Upserts and deletes
- Fully configurable/optimizable
- Structured Streaming support

# Components of Delta Lake

# Delta Lake is comprised of:

- Delta tables
- The Delta optimization engine
- The Delta Lake storage layer

# Delta Tables

## Data files

- Parquet format
- Kept in cloud storage

## Table registered in the Metastore

- Contains the data schema and metadata

## Transaction log

- Kept in cloud storage
- Changes are stored as ordered, atomic commits
- Records every transaction that occurs
- Allows for Time Travel
- Single source of truth

# The Delta Optimization Engine

- Thanks to Apache Spark!

- File management optimizations

  - Compaction, data skipping, localized data storage

- Auto-optimized writes and file compaction

- Performance optimization via Delta caching

# The Delta Lake Storage Layer

- Highly performant and persistent
- Low-cost, easily scalable object storage
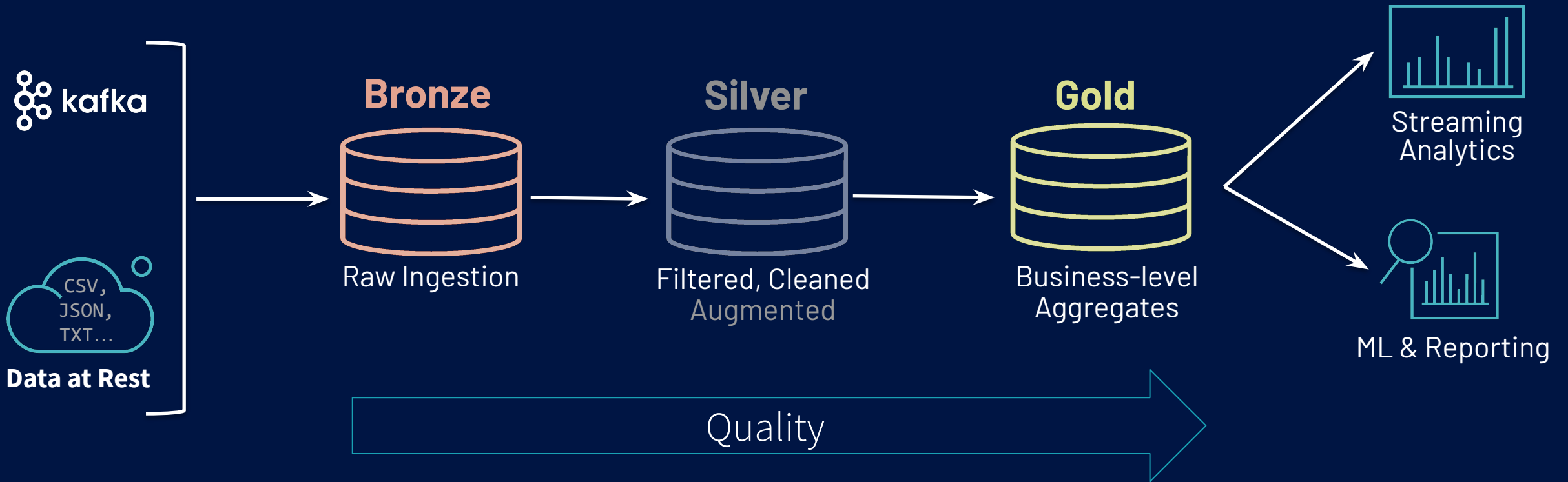- Ensures consistency
- Allows for flexibility

# The Delta Architecture

kafka

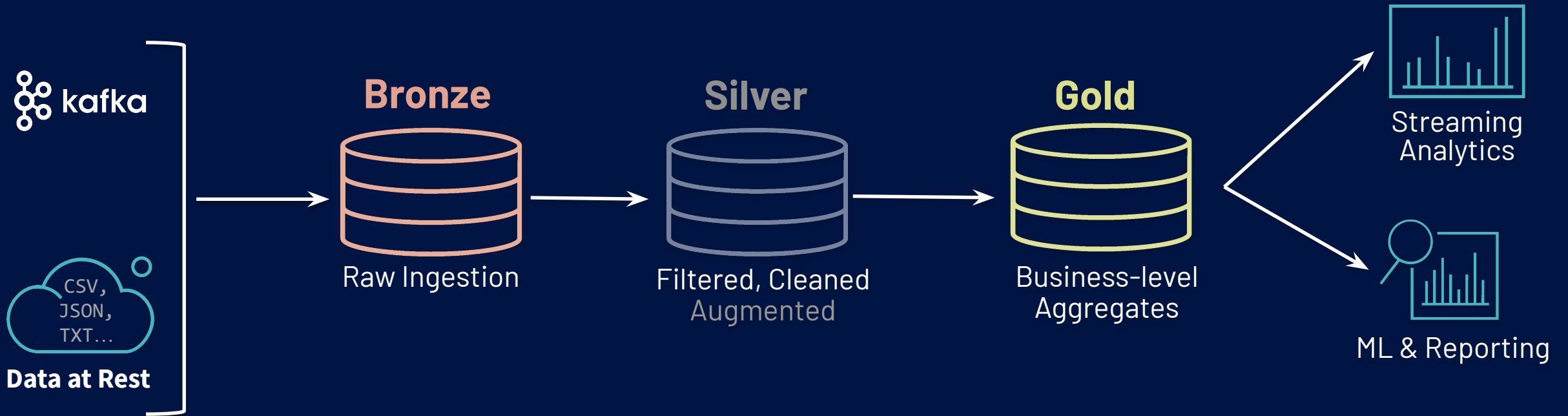CSV, JSON, TXT...

**Data at Rest**

**Bronze**

Raw Ingestion

**Silver**

Filtered, Cleaned
Augmented

**Gold**

Business-level
Aggregates
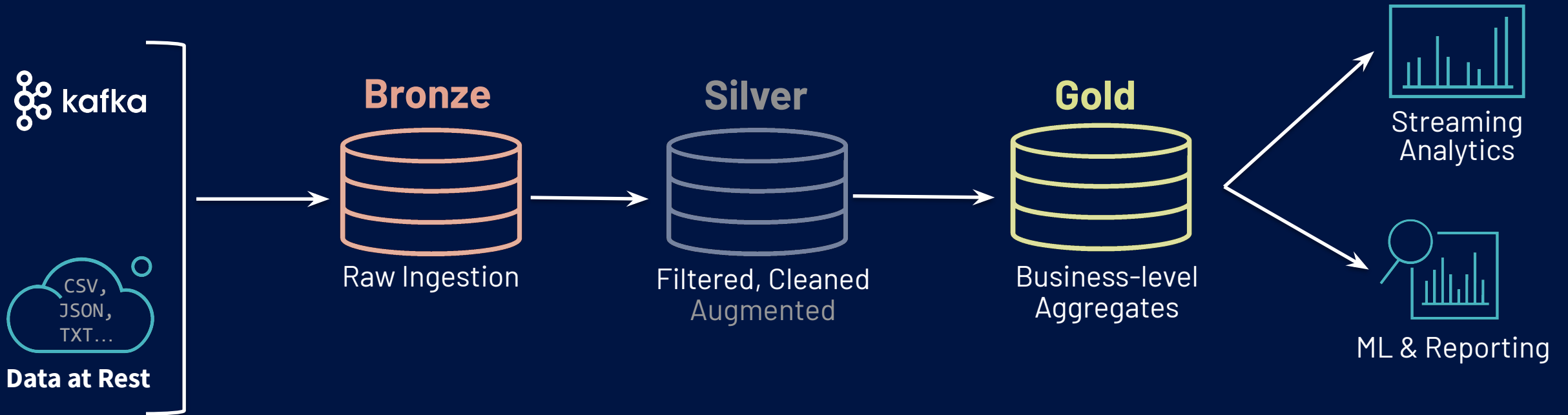
Streaming
Analytics

ML & Reporting

Quality

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

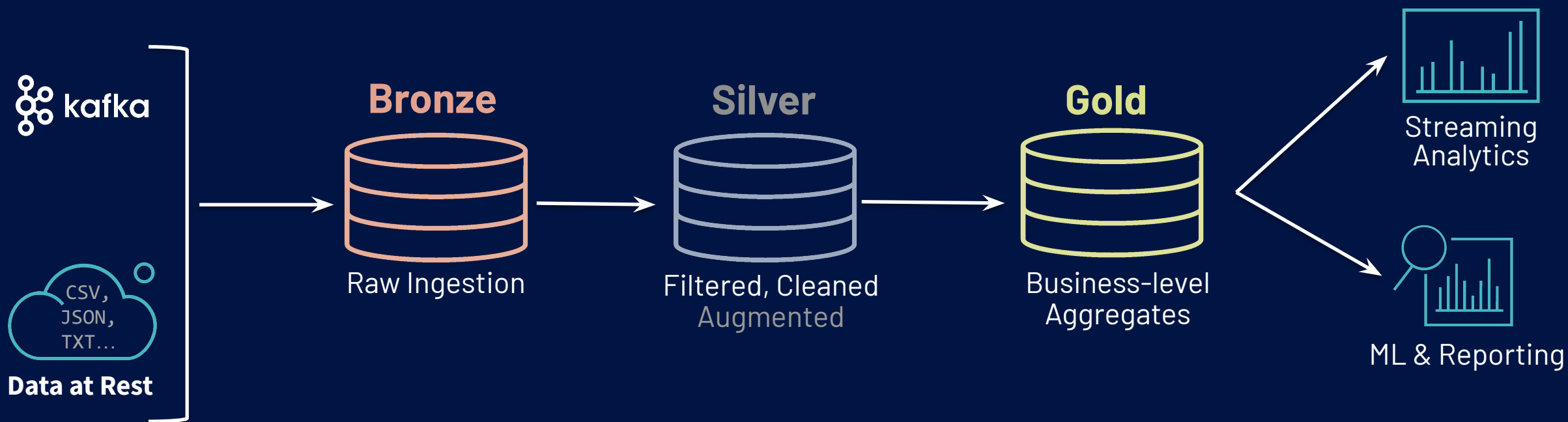databricks

# The Delta Architecture

**Bronze** — Raw Ingestion

**Silver** — Filtered, Cleaned Augmented

**Gold** — Business-level Aggregates

kafka

CSV, JSON, TXT...
**Data at Rest**

Streaming Analytics

ML & Reporting

- Dumping ground for raw data
- Often with long retention (years)
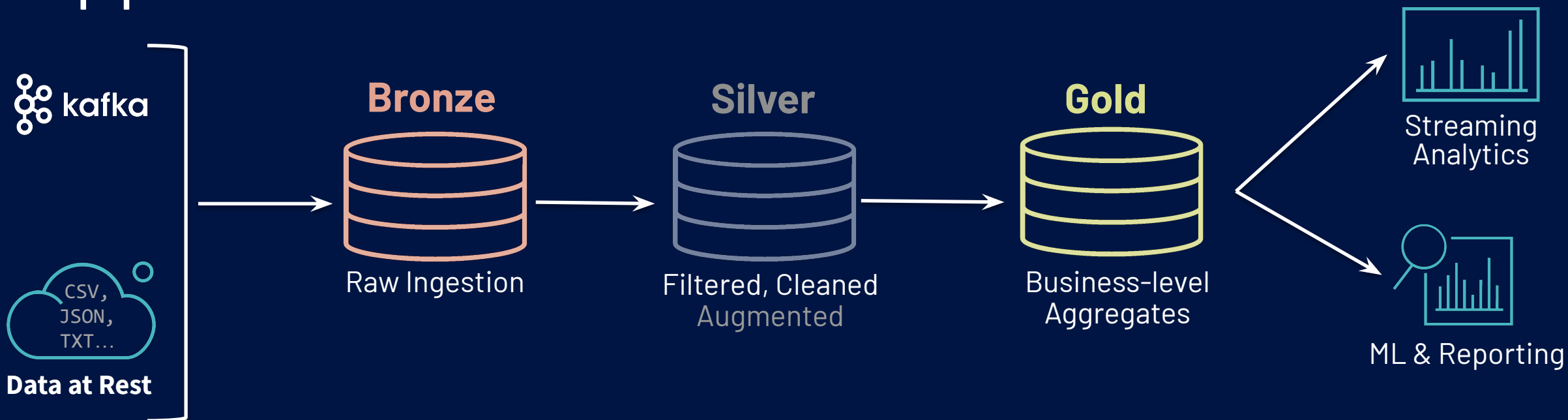- Raw data with minimal parsing\

databricks

# The Delta Architecture



kafka

CSV, JSON, TXT...

**Data at Rest**

**Bronze**

Raw Ingestion

**Silver**

Filtered, Cleaned
Augmented

**Gold**

Business-level
Aggregates

Streaming
Analytics

ML & Reporting

Intermediate data with some cleanup applied.

Queryable for easy debugging!

databricks

# The Delta Architecture

kafka

CSV,
JSON,
TXT...

**Data at Rest**

**Bronze**
Raw Ingestion

**Silver**
Filtered, Cleaned
Augmented

**Gold**
Business-level
Aggregates

Streaming
Analytics

ML & Reporting

Clean data, ready for consumption

databricks

# Simplify real-time / streaming data applications

kafka

CSV, JSON, TXT...

**Data at Rest**

**Bronze**

Raw Ingestion

**Silver**

Filtered, Cleaned
Augmented

**Gold**

Business-level
Aggregates

Streaming
Analytics

ML & Reporting

Streams move data through the Delta Lake
- Low-latency or manually triggered
- Eliminates management of schedules and jobs

databricks

# 10 minute break

# Lab 1: Building Delta Pipelines

# 10 minute break

SPARK+AI SUMMIT

# Features of Delta Lake

# Unification of batch and streaming

- Continuous data flow model

- Same engine. Same APIs. Same user code. No need to reason about system complexities separately.

- A table in Delta Lake is a batch table as well as a streaming source and sink.

- Can stream late arriving data and backfilled historical data through the same pipeline without having to delay downstream processing

- Incrementally load the new data efficiently. No need to do state management on what are the new files added.

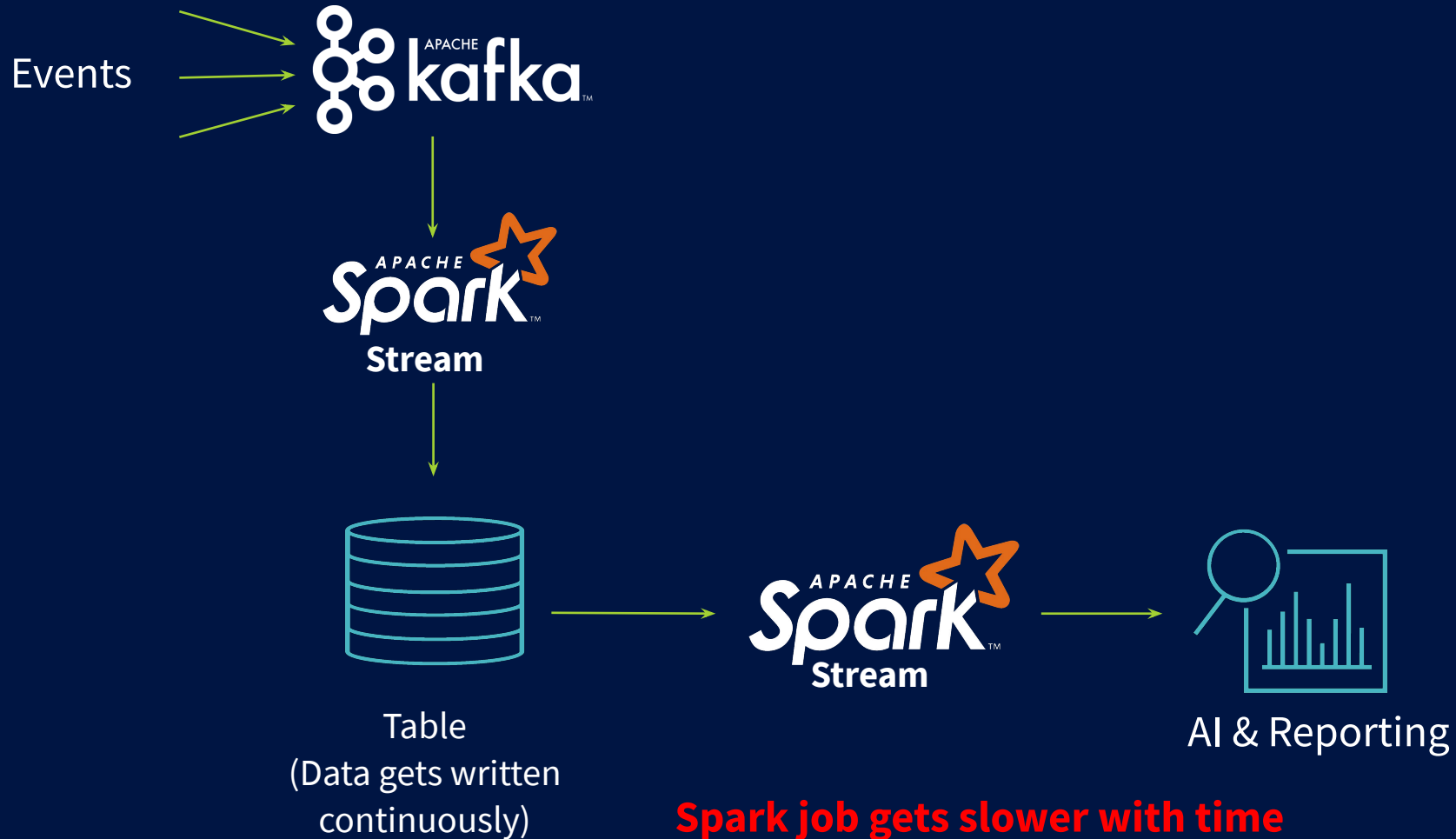- Process the data quickly as it arrives without any delays.

# A Data Engineer's Dream...

Process data **continuously** and **incrementally** as new data arrive in a **cost efficient way** without having to *choose* between batch or streaming
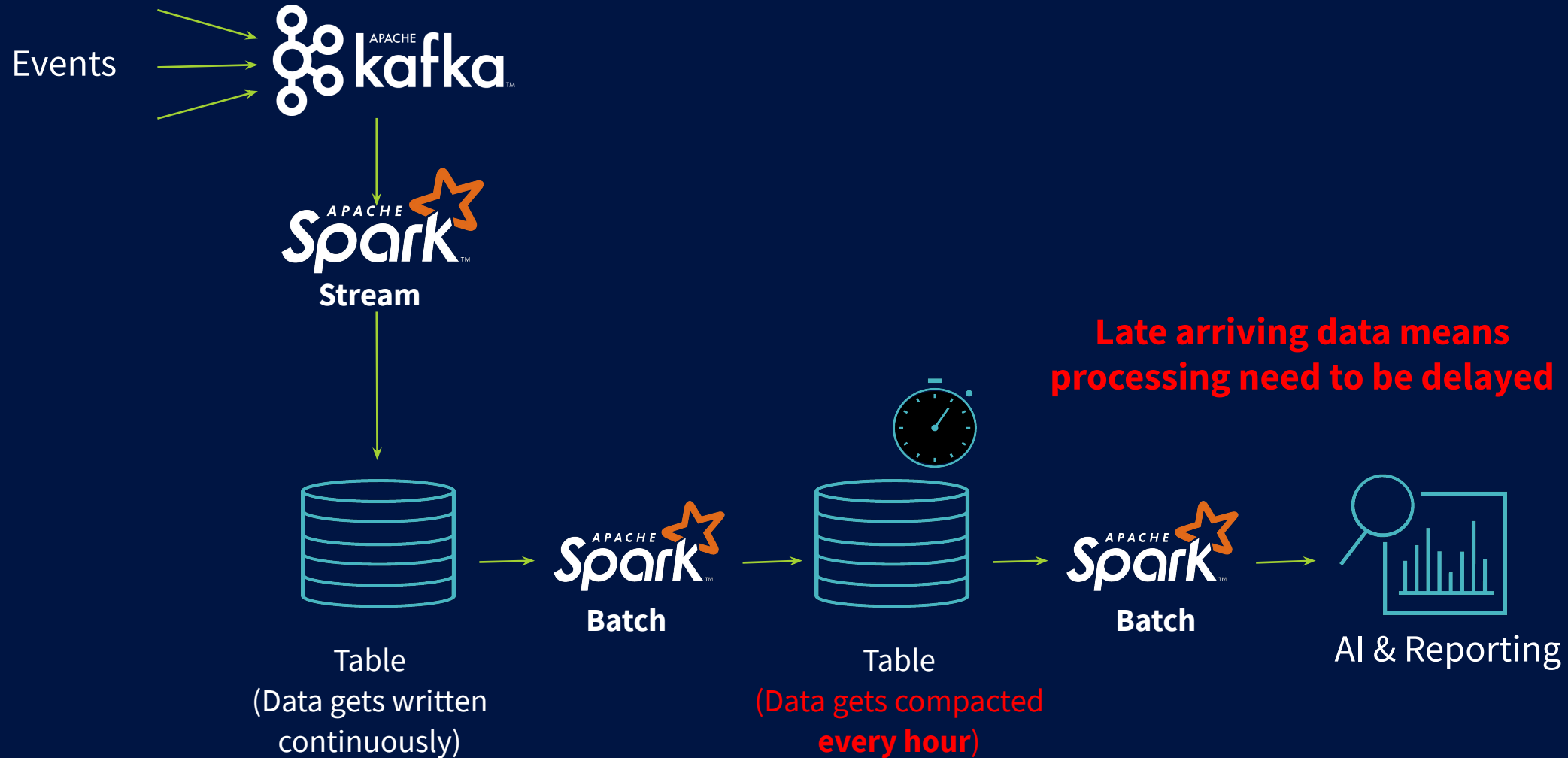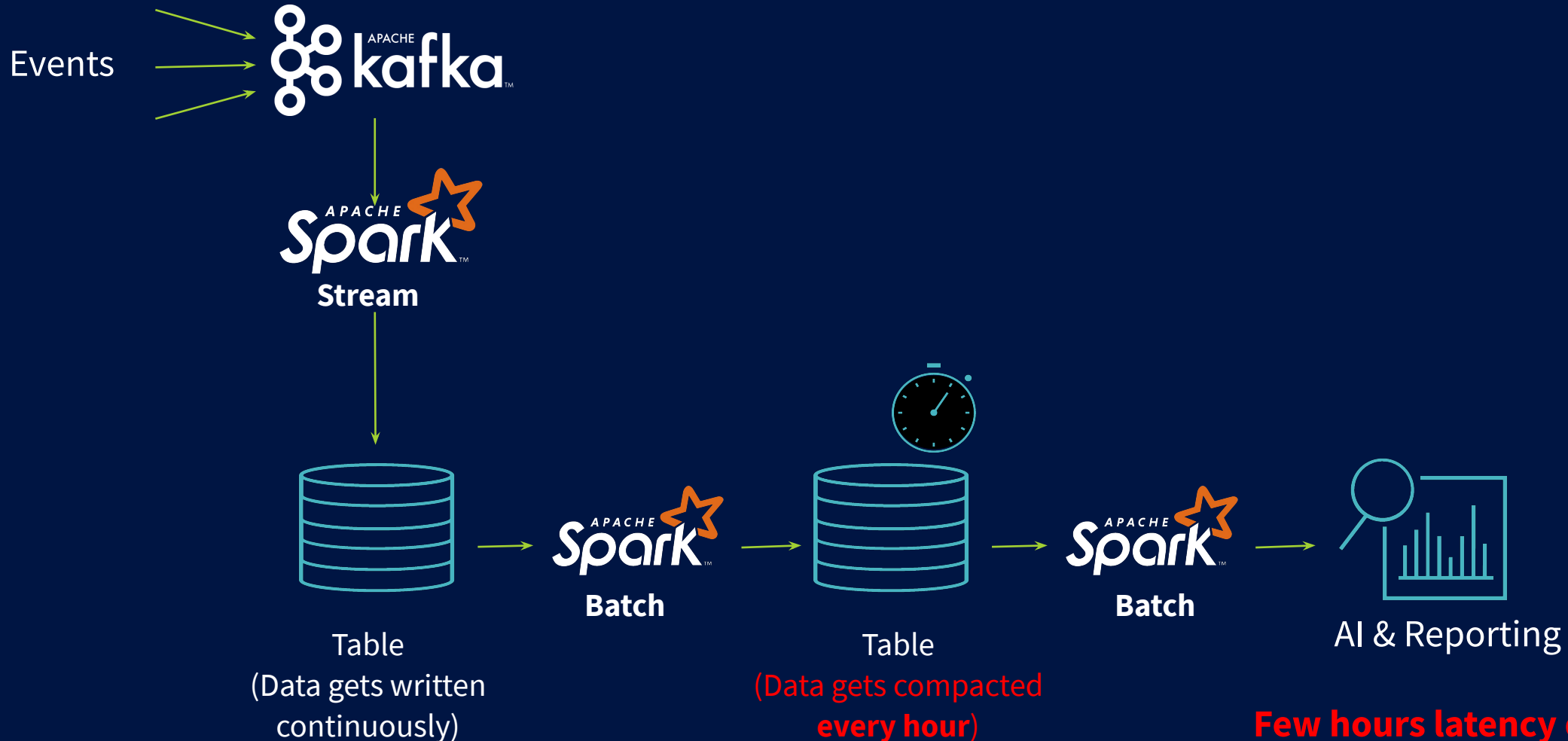
# The Data Engineer's Journey...



Events

APACHE kafka™

APACHE Spark™
**Stream**

Table
(Data gets written
continuously)

APACHE Spark™
**Stream**

AI & Reporting

**Spark job gets slower with time
due to small files.**

# The Data Engineer's Journey...

Events

Apache Kafka

Apache Spark
**Stream**

**Late arriving data means
processing need to be delayed**

Table
(Data gets written
continuously)

Apache Spark
**Batch**

Table
(Data gets compacted
**every hour**)

Apache Spark
**Batch**

AI & Reporting

SPARK+AI SUMMIT

databricks

# The Data Engineer's Journey…



Events

Apache Kafka

Apache Spark
**Stream**

Apache Spark
**Batch**

Apache Spark
**Batch**

AI & Reporting

Table
(Data gets written continuously)

Table
(Data gets compacted **every hour**)

**Few hours latency doesn't satisfy business needs**

SPARK+AI SUMMIT

# The Data Engineer's Journey...



Events → **Apache kafka**

**Apache Spark Stream**

**Apache Spark Stream**

Table
(Data gets written continuously)

**Apache Spark Batch**

Table
(Data gets compacted **every hour**)

**Apache Spark**

Unified View → AI & Reporting

**Lambda arch increases operational burden**

#Datateams  #SparkAISummit

databricks

# The Data Engineer's Journey...



Events → **kafka** → **Spark** Stream → Unified View → AI & Reporting

**Spark** Stream ← Validation

Table
(Data gets written continuously) → **Spark** Batch → Table
(Data gets compacted **every hour**) → **Spark** Batch → Unified View

**Validations and other cleanup actions need to be done twice**

# The Data Engineer's Journey...



Events

APACHE kafka

APACHE Spark — Stream

APACHE Spark — Stream

Validation

APACHE Spark — Batch

Reprocessing

Table
(Data gets written continuously)

Table
(Data gets compacted **every hour**)

APACHE Spark — Batch

Unified View

AI & Reporting

**Fixing mistakes means blowing up partitions and doing atomic re-publish**

# The Data Engineer's Journey...

Events

APACHE kafka

APACHE Spark Stream

Validation

APACHE Spark Stream

APACHE Spark Stream

Table
(Data gets written continuously)

APACHE Spark Batch

Reprocessing

Table
(Data gets compacted **every hour**)

APACHE Spark Batch

Update & Merge

Unified View

AI & Reporting

**Updates & Merge get complex with data lake**

# Schema Enforcement and Evolution

- Schema enforcement (aka validation): prevents bad writes

- Schema evolution: allows for updates to data schema

- Delta saves table schema in JSON format in the transaction log

- Data to be written cannot have:

  - Additional columns

  - Data type mismatches

  - Column names that only differ by case

- Schema evolution commonly used for appends or overwrites

# A brief look at time travel

- Utilizes the transaction log
- Can recreate a table's state at any point in time
- Automatic versioning
- Rollback to previous versions in case of bad writes
- Easily perform deletes

# A brief look at time travel

# The Delta Architecture



Easy to recompute when business logic changes:
- Clear tables
- Restart streams

Lab 2: Exploring Delta Lake features

# Wrap Up

# Summary

In this class, we learned:

- The benefits of the Delta architecture
- How to build a data pipeline using Delta Lake
- How to incorporate batch and streaming data in a Delta Lake pipeline
- How Delta Lake schema enforcement and evolution impacts data writes and updates
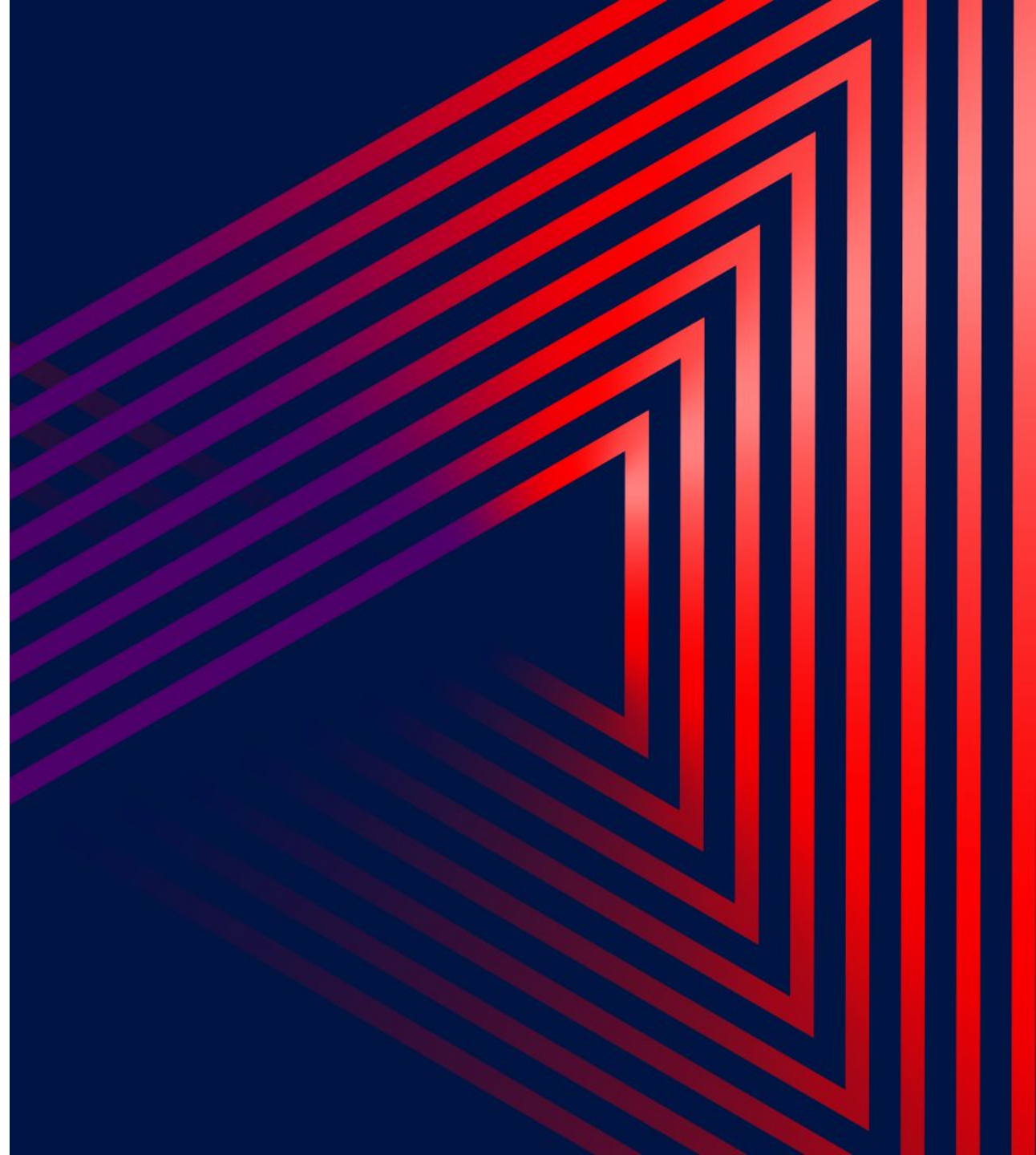- How to perform data deletions and version rollbacks

Q & A

# Course evaluation

Thank you for taking this class!

Please take a moment to fill out
the feedback form.

# Feedback

Your feedback is important to us.

Don't forget to rate and review the sessions.

**SPARK+AI** SUMMIT

#Datateams #SparkAISummit