# The big picture

databricks

# Lesson goals

- Recall foundational concepts related to the contemporary big data landscape.
- Define the core challenges with building a big data architecture.
- Explain how the core components of a Lakehouse relate to the Inmon Architecture.
- Explain how Delta Lake can be used to build a Lakehouse.
- Describe the core components of Delta Lake.
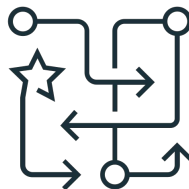
databricks

# The big data landscape

databricks

# The big data problem

**Volume**  **Velocity**  **Variety**  **Veracity**  **Value**

# At Moovio, your SLA includes:

**Data freshness**

**Query speed**

**Data reliability**

**Ease of use**

# Big data needs

High flux event data

Long-term data storage

Real-time dashboards

Periodic reporting

Artificial intelligence

# A single source of truth

# ODS and OLTP

Operational Data Store and Online Transaction Processing



**OLTP**

**ODS**

Transactions

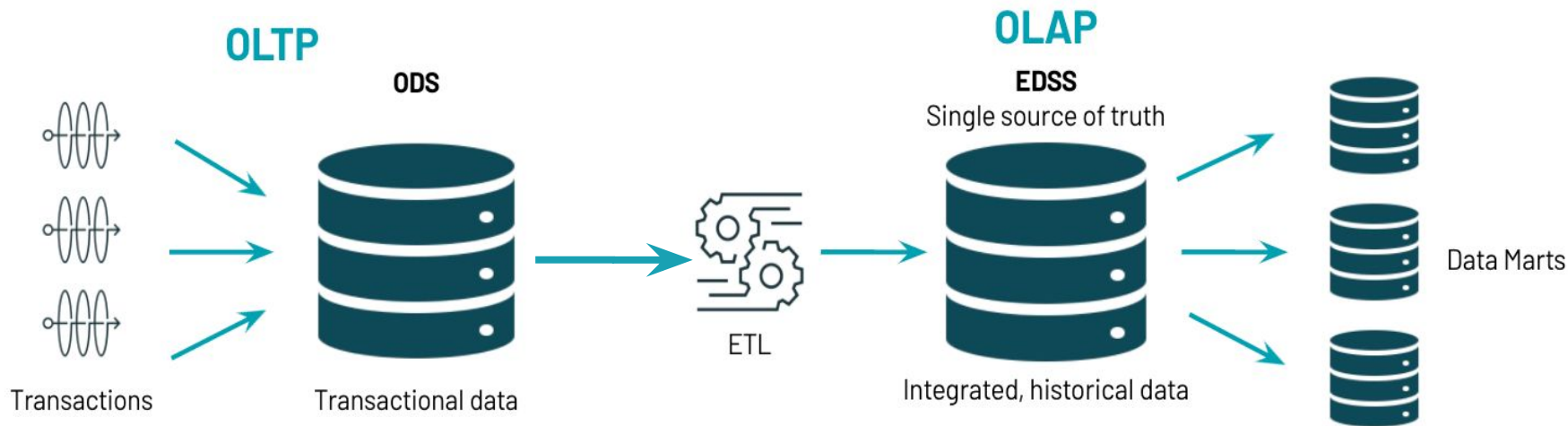Transactional data

databricks

# EDSS and OLAP
The Enterprise Decision Support System and Online Analytical Processing

# Complete data system

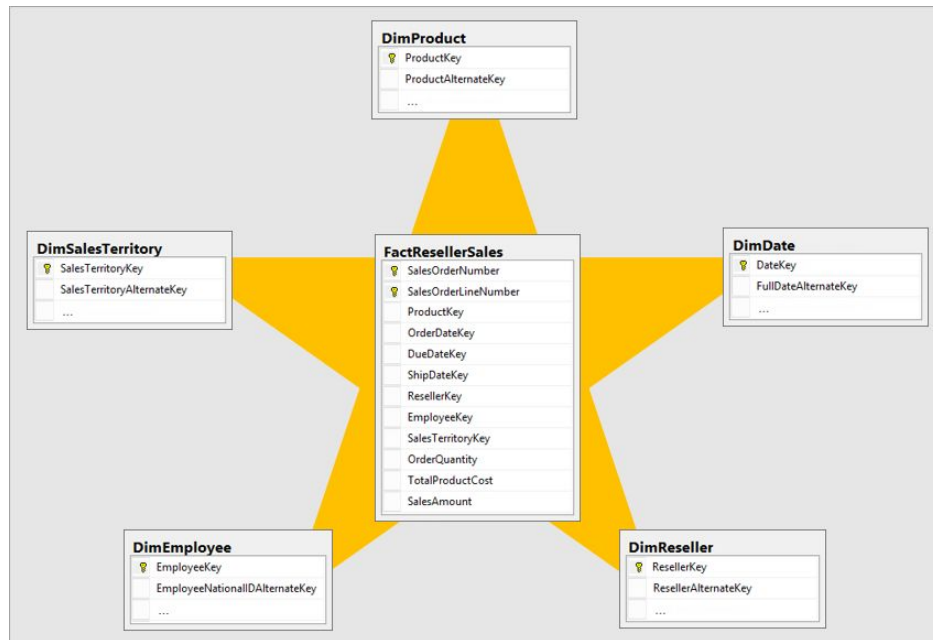- An ETL process pulls data from the ODS to be loaded into the EDSS

# Levels of data

Operational
(ODS)

Atomic
(EDSS)

Departmental
(Data Marts)

Data Marts
(Individual)

databricks

# Fact and dimension tables
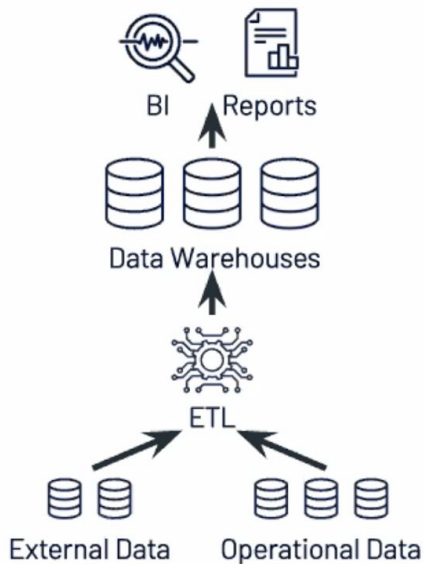
## Dimensional modeling

- Fact tables
- Dimension tables
- Aggregate fact tables

# The Lakehouse

databricks
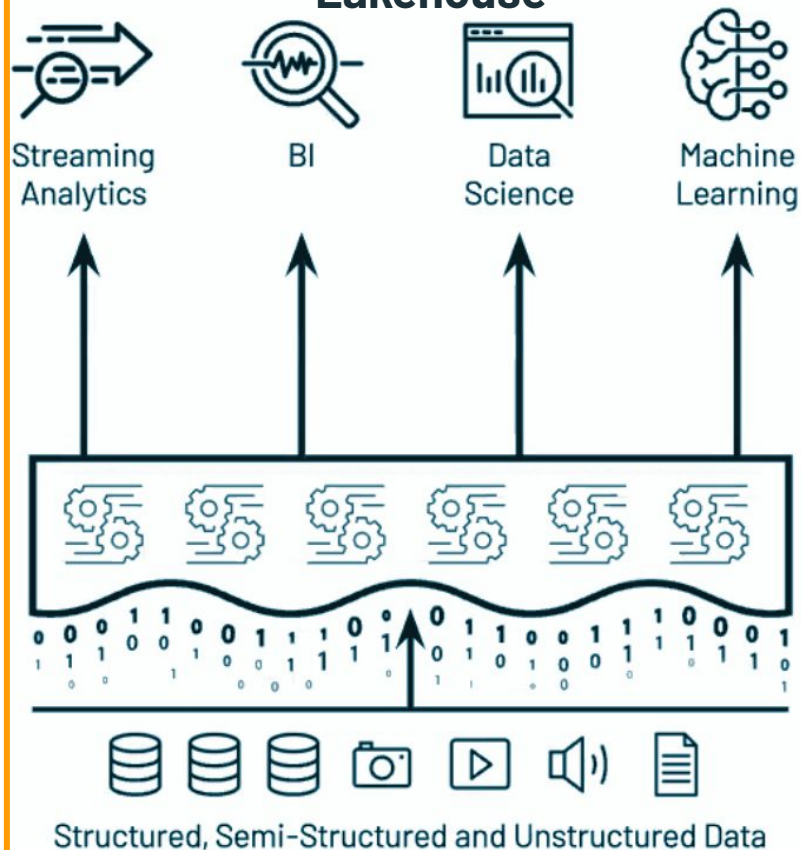
# What is a Lakehouse?



Data Warehouse

BI    Reports

Data Warehouses

ETL

External Data    Operational Data



Data Lake

Data Science    Machine Learning    Real-Time Database    Reports    BI

Data Prep and Validation    ETL    Data Warehouses

Data Lake

Structured, Semi-Structured and Unstructured Data



**Lakehouse**

Streaming Analytics    BI    Data Science    Machine Learning

Structured, Semi-Structured and Unstructured Data

databricks

# Levels of data in a Lakehouse



**Operational ODS**

**Atomic EDSS Data Lake**

ETL

Metadata, Summary data Raw data

ETL

**Departmental Data Marts Data warehouse**

**Individual End users**

databricks

# **Benefits** of a Lakehouse

- Separation of compute and storage
- Infinite storage capacity
- Leverage best aspects of a data warehouse
- Low data gravity
- High data throughput
- No limits on data structure
- Mix batch and streaming workloads

databricks

# The Delta architecture

# What is **Delta Lake**?

- Technology designed to be used with Apache Spark to build robust data lakes

# Delta Lake features

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



databricks

# Delta Lake components

**Delta Lake storage layer**

**Delta tables**

**Delta Engine**

databricks

# Delta Lake components

**Delta Lake storage layer**

**Delta tables**

**Delta Engine**

databricks

# Delta Lake storage layer

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

databricks

# Delta Lake components

**Delta Lake storage layer**

**Delta tables**

**Delta Engine**

databricks

# Delta table components

- Data in Parquet/Delta files
- Transaction log
- Registered in metastore
  (optional)

databricks

# Data – Parquet files

- File format for tabular data stored as columns
- Fast and powerful
- Delta files = Parquet + versioning + metadata

databricks

# Transaction log

- Record of all transactions on a Delta table
- Prevents read conflicts
- Commits - ordered, atomic, json files
- Created automatically in the *_delta_log* subdirectory

databricks

# Delta Lake components



**Delta Lake storage layer**

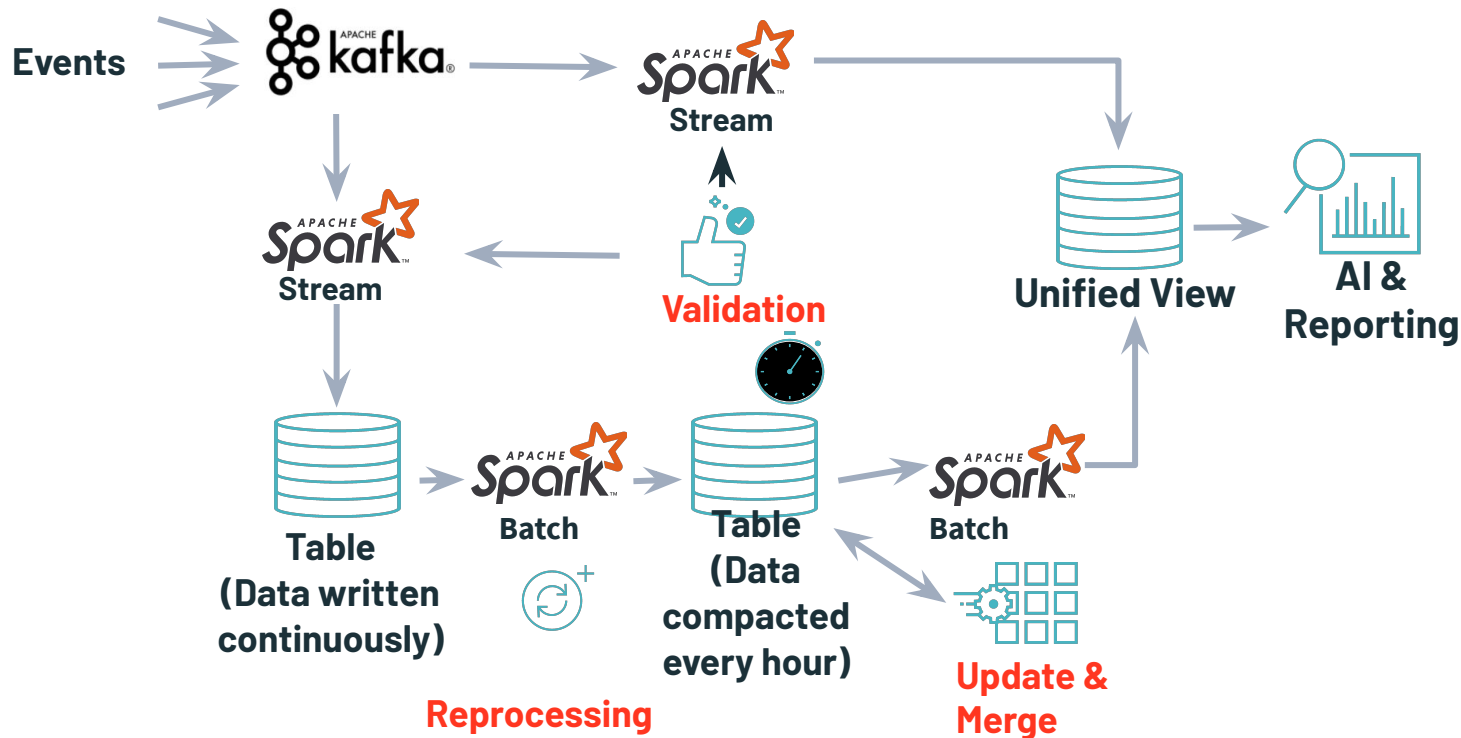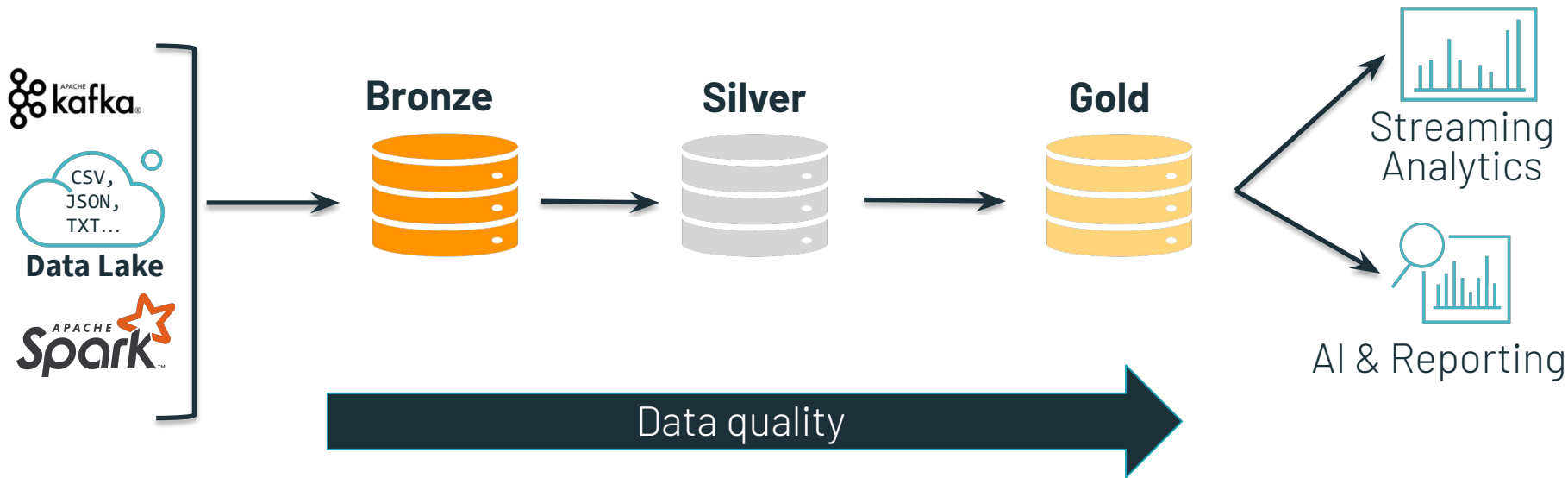**Delta tables**

**Delta Engine**

databricks

# Delta Engine

- File management optimizations
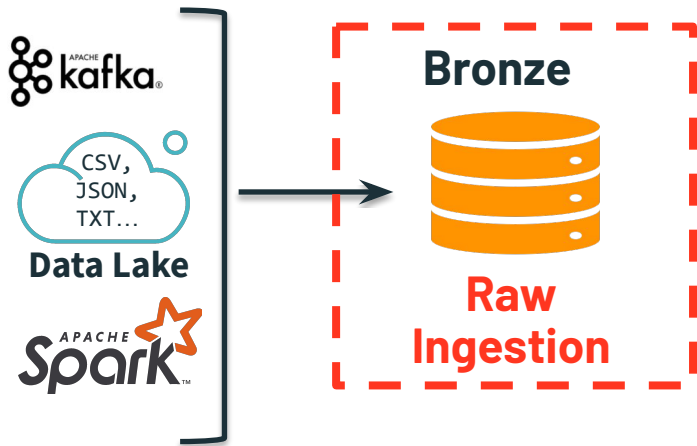- Auto-optimized writes
- Performance optimization via Delta caching

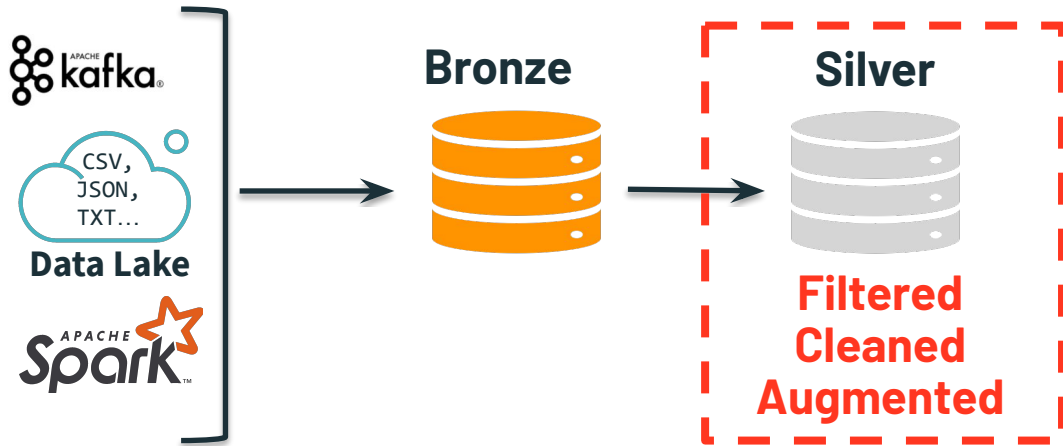databricks

# The goal of a data engineer

# The Delta architecture design pattern

# Delta architecture – Bronze



Data Lake
Bronze
Raw Ingestion
databricks

# Delta architecture - Silver

# Delta architecture - Gold