



Introduction to Delta Lake

Surendra Panpaliya

Agenda



What is Delta Lake?



Delta Lake operations on Databricks



Create a table.



Upsert to a table.



Read from a table.



Display table history

Agenda

Query an earlier version of a table.

Optimize a table

Add a Z-order index.

Liquid Clustering

Vacuum unreferenced files

What is Delta Lake?

An open-source Storage Layer

Provides ACID transactions

Atomicity, Consistency, Isolation, Durability

Scalable metadata handling

Time travel capabilities to data lakes.

Delta table versions

Version 0

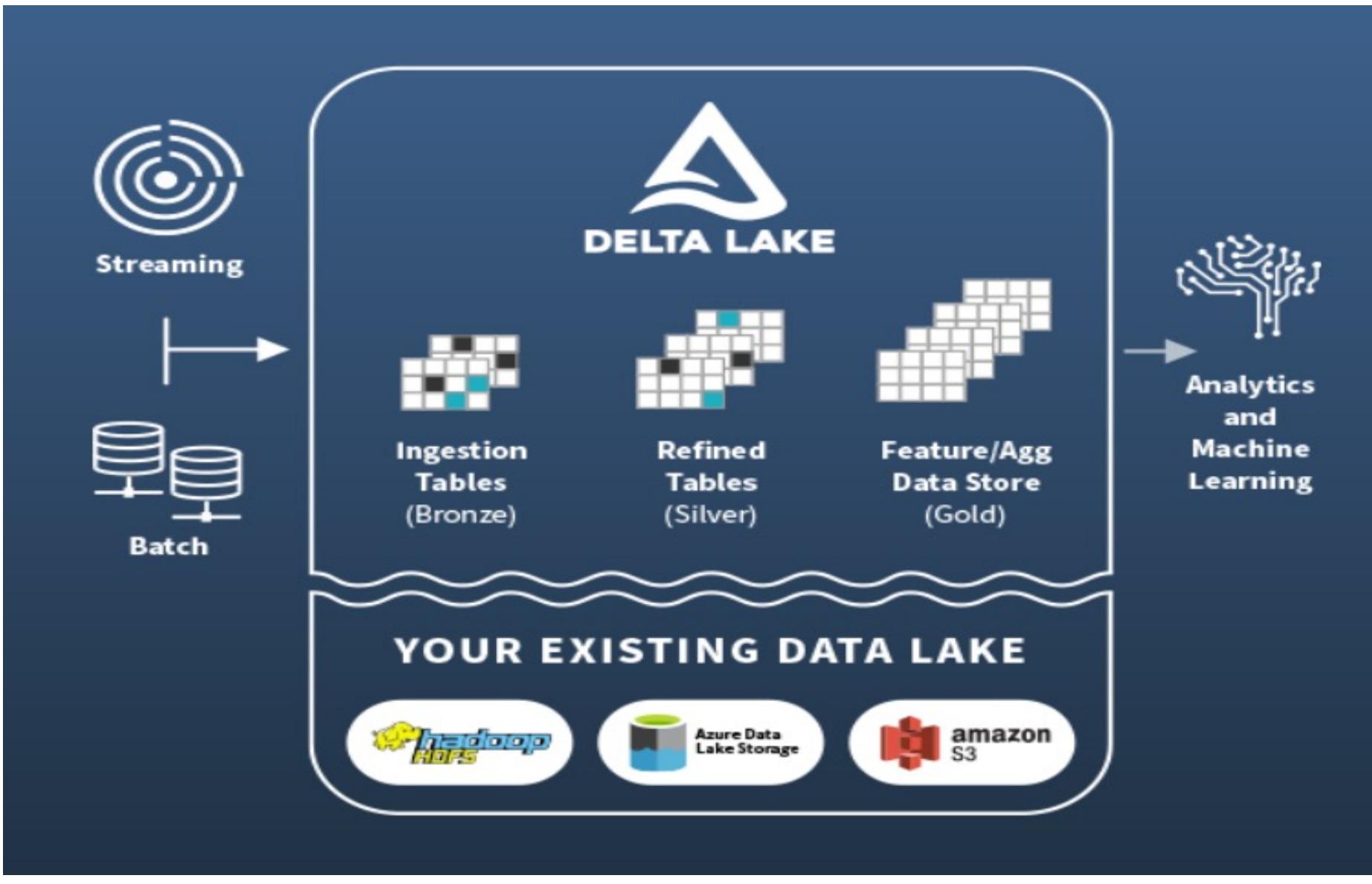
id
0
1
2

Version 1

id
0
1
2
8
9
10

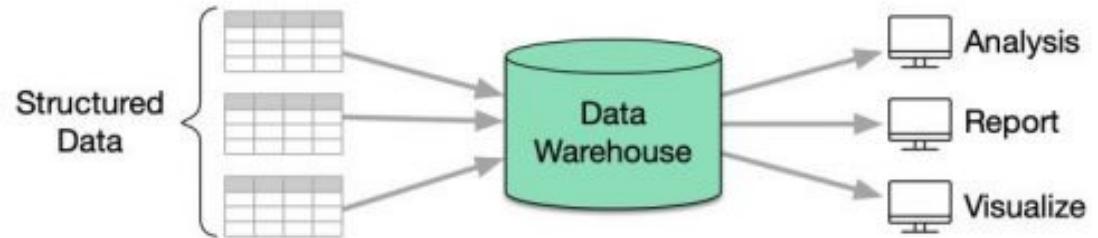
Version 2

id
55
66
77



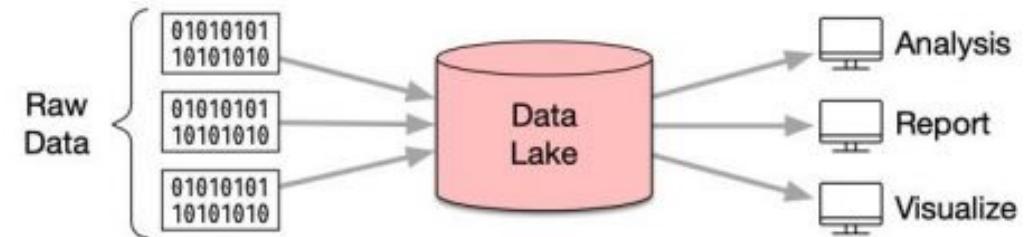
Data Warehouse

A large, structured repository of integrated data from various sources, used for complex querying and historical analysis



Data Lake

A more focused, department-specific subset of a data warehouse providing quick data retrieval and analysis



What is Data Warehouse?



A data warehouse is



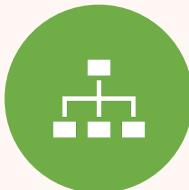
A central repository



Integrated and structured data



From various sources



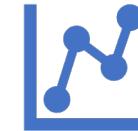
within an organization.

What is Data Warehouse?

Designed to support



Business intelligence (BI) and



Analytics activities by



Providing a unified

Consistent view of data.

What is Data Warehouse?



DW typically follow



A schema-on-write approach,



Where data is transformed and



Loaded into predefined schemas



Optimized for Reporting and Analysis.

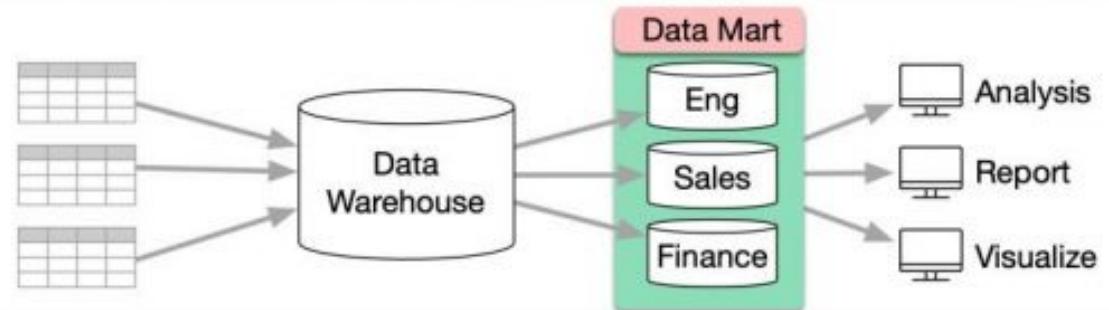


VS



Data Mart

A vast pool of raw, unstructured data stored in its native format until it's needed for use



Delta Lake

An open-source storage layer that brings reliability and ACID transactions to data lakes, unifying batch and streaming data processing



Data Mart



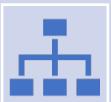
A data mart is a



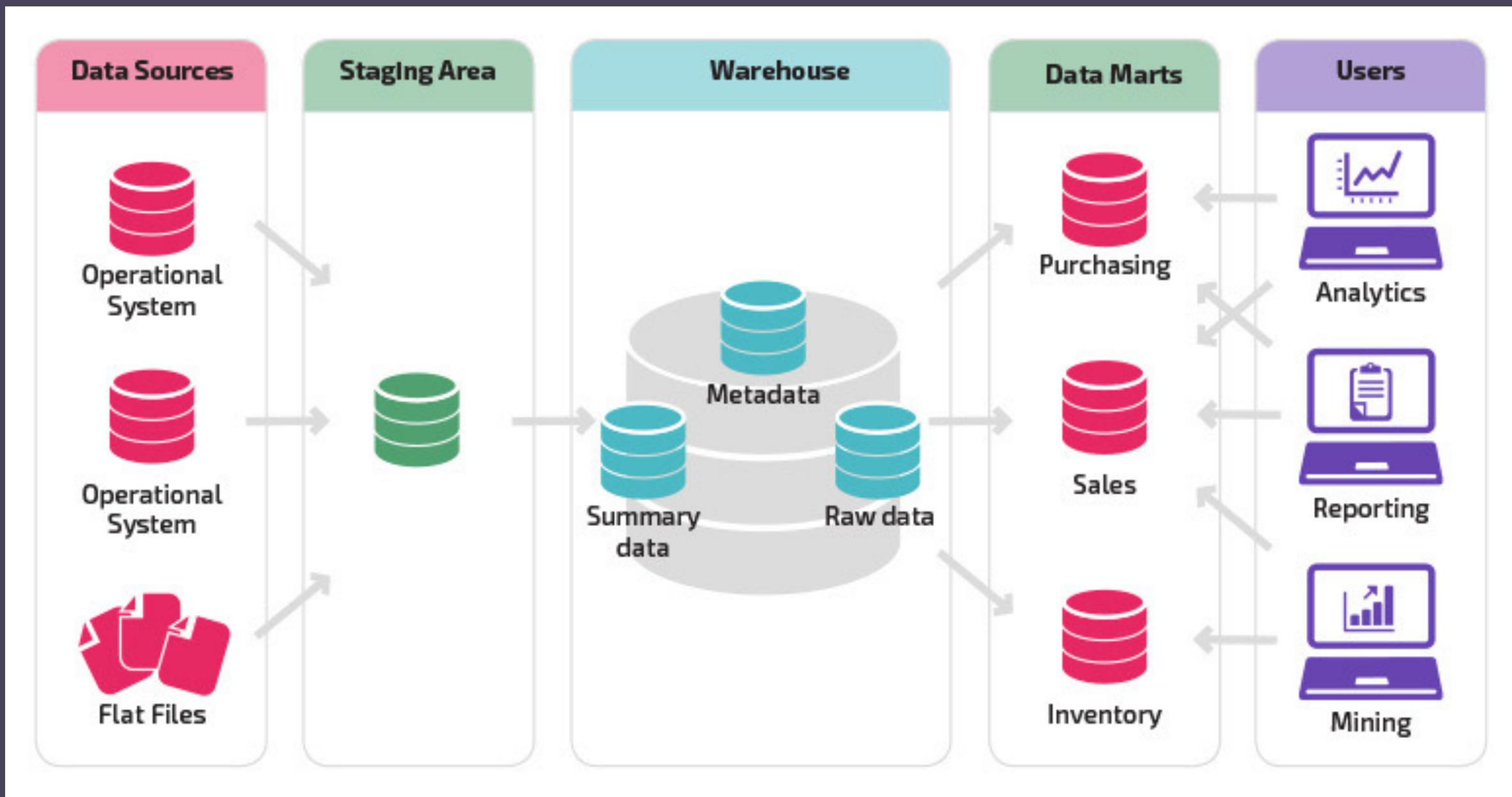
subset of a data warehouse



focused on a



specific business function or department.



Data Mart



Contains a subset of data



Relevant to a particular group of users,



making it more



targeted and specialized.

Data Mart



Data marts are often designed



to provide faster and



more specific insights

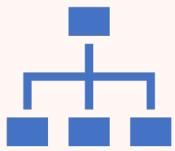


compared to the broader

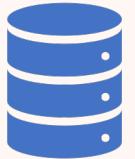


data warehouse.

Data Mart



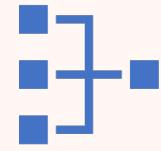
Can be
structured



Similar to a data
warehouse or



May use
different



Schemas and
models.

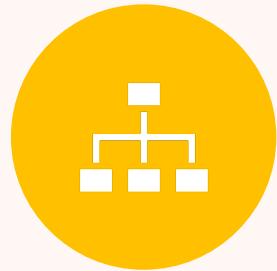
What is Data Lake?



A DATA LAKE IS A
LARGE,



CENTRALIZED
REPOSITORY



THAT STORES
STRUCTURED, SEMI-
STRUCTURED, AND



UNSTRUCTURED DATA
IN ITS RAW AND
UNPROCESSED FORM.

What is Data Lake?



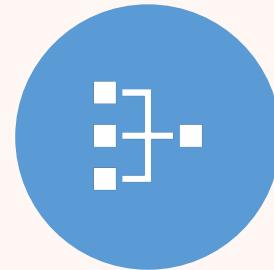
ALLOWS
ORGANIZATIONS



TO STORE VAST
AMOUNTS OF DATA



FROM DIFFERENT
SOURCES



WITHOUT PREDEFINED
SCHEMAS.

What is Data Lake?



DATA LAKES ENABLE
DATA EXPLORATION,



ADVANCED ANALYTICS,
AND



MACHINE LEARNING
BY PROVIDING
FLEXIBILITY



IN DATA PROCESSING
AND ANALYSIS.

What is Data Lake?



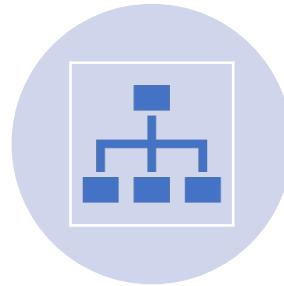
The data lake architecture typically



follows a schema-on-read approach,



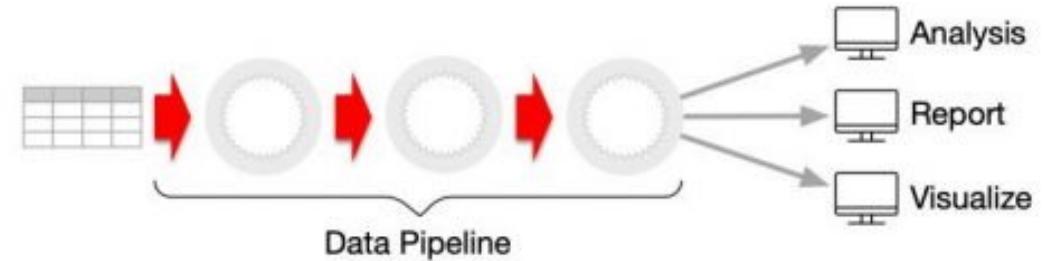
where data is transformed and



structured at the time of analysis.

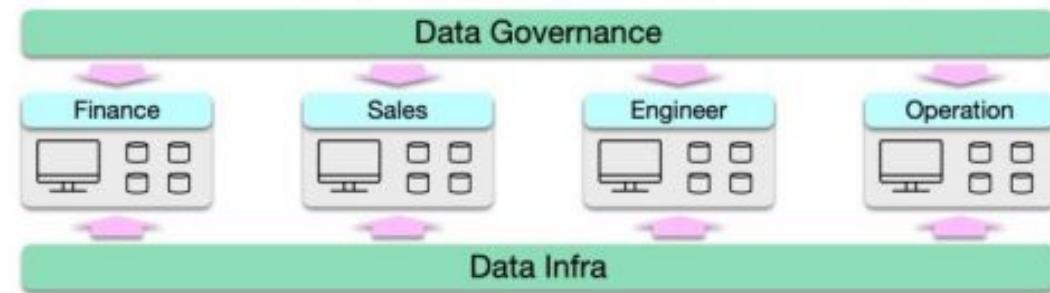
Data Pipeline

A process that moves and transforms data from one system to another, often used to populate data warehouses and data lakes



Data Mesh

An architectural and organizational approach where data ownership and delivery are decentralized across domain-specific, cross-functional teams



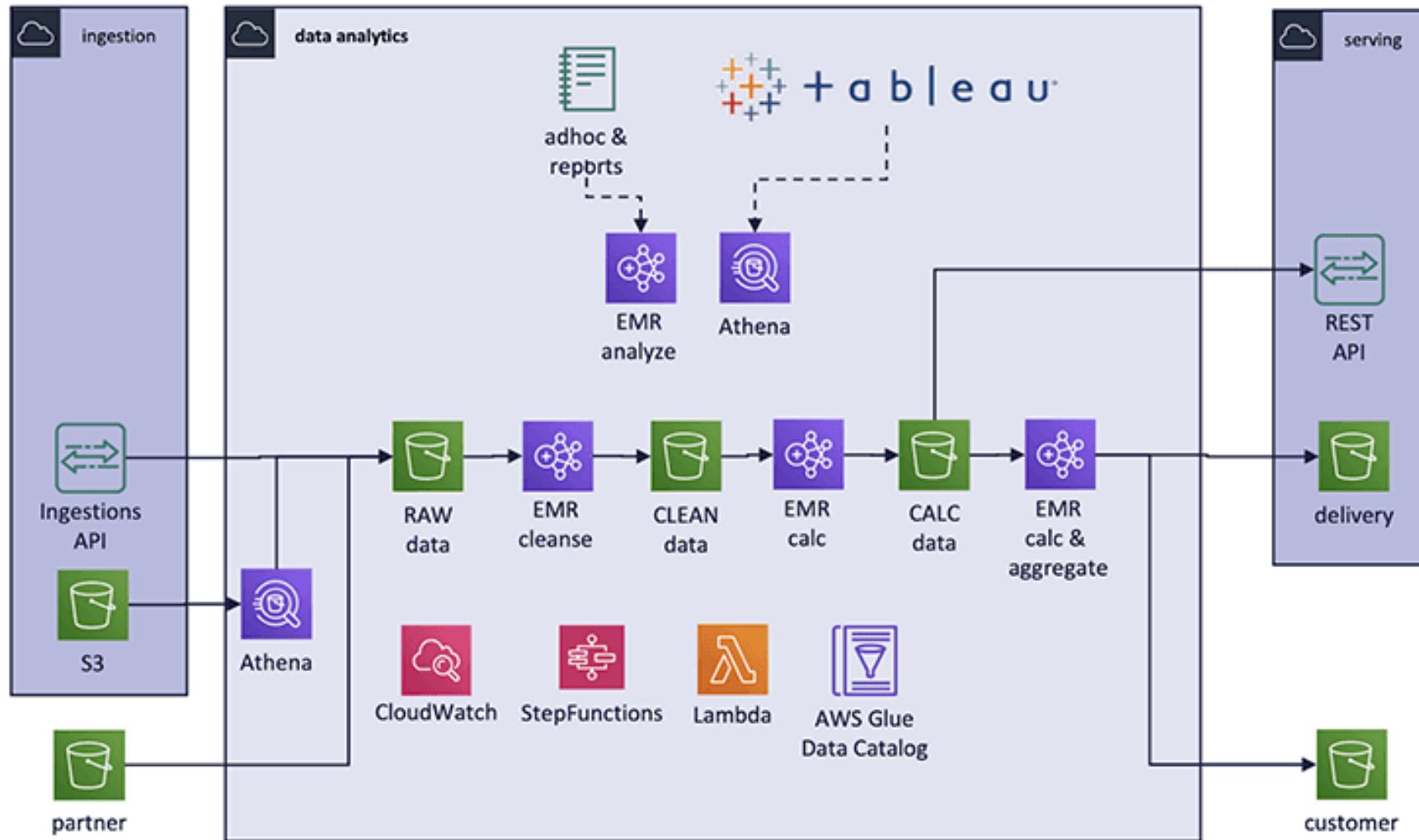
Data Pipeline

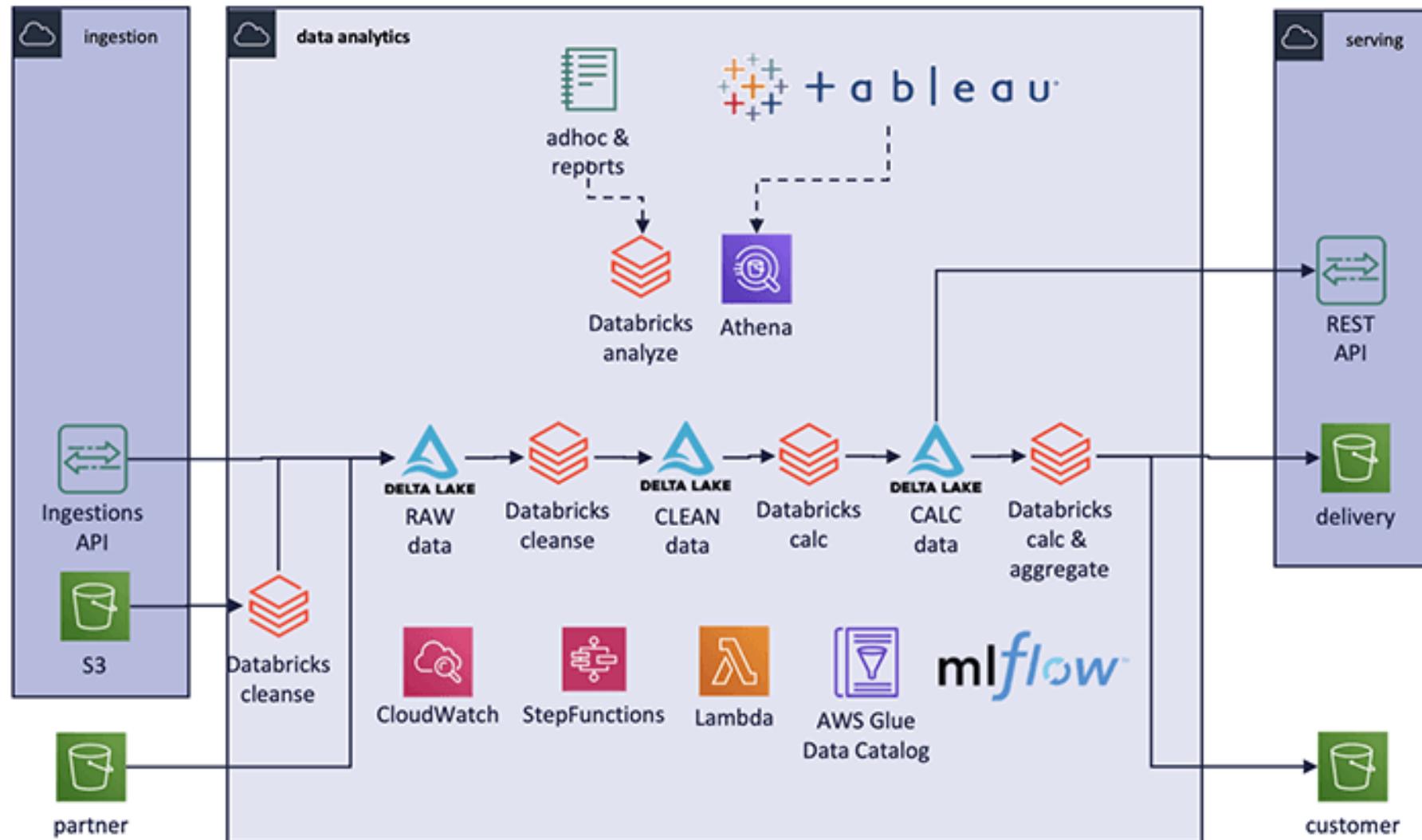
- A system or framework
- Facilitates the automated
- Extraction, Transformation, Loading (ETL)
- From various sources into
- Target storage or processing system.



High level architecture & tech stack

How our world looked before Databricks





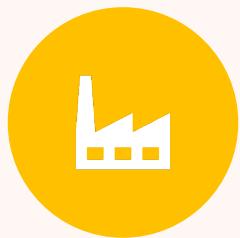
Data Pipeline



ENABLES THE
SMOOTH FLOW OF
DATA



FROM SOURCE
SYSTEMS



TO DATA
WAREHOUSES,



DATA LAKES, OR



OTHER
DESTINATIONS.

Data Pipeline



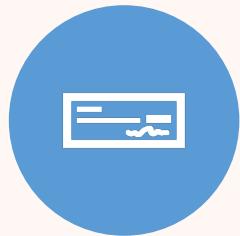
DATA PIPELINES
HANDLE TASKS



DATA INGESTION,



DATA
TRANSFORMATION,



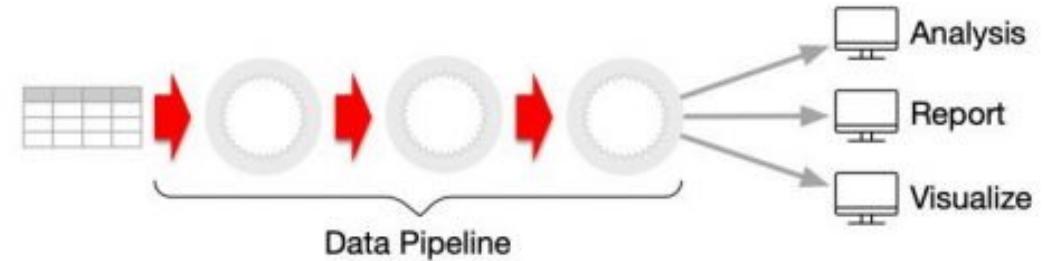
DATA QUALITY
CHECKS,



DATA DELIVERY.

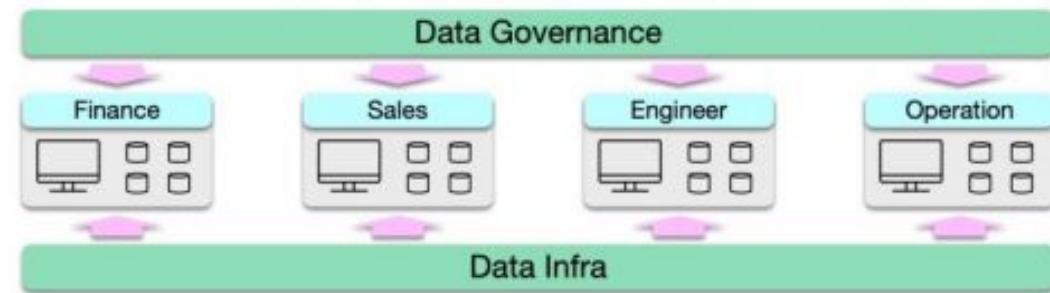
Data Pipeline

A process that moves and transforms data from one system to another, often used to populate data warehouses and data lakes



Data Mesh

An architectural and organizational approach where data ownership and delivery are decentralized across domain-specific, cross-functional teams



Data Mesh

Data Mesh is a

decentralized approach

to data architecture and

management.

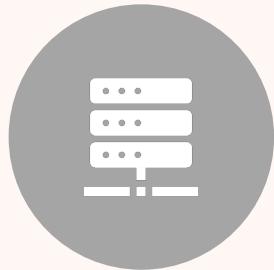
Data Mesh Architecture



Data Mesh



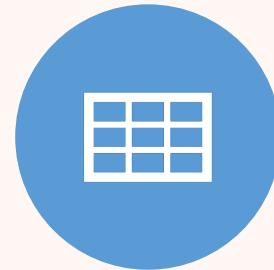
IT EMPHASIZES
DOMAIN-ORIENTED,



SELF-SERVE DATA
PRODUCTS,



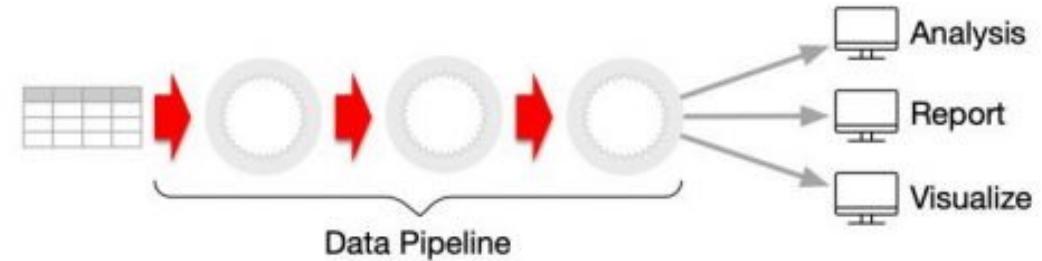
WHERE CROSS-
FUNCTIONAL TEAMS



TAKE OWNERSHIP OF
THEIR DATA DOMAINS.

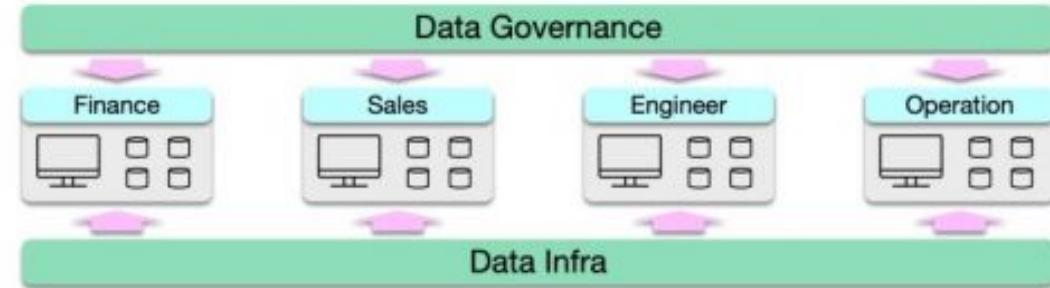
Data Pipeline

A process that moves and transforms data from one system to another, often used to populate data warehouses and data lakes



Data Mesh

An architectural and organizational approach where data ownership and delivery are decentralized across domain-specific, cross-functional teams



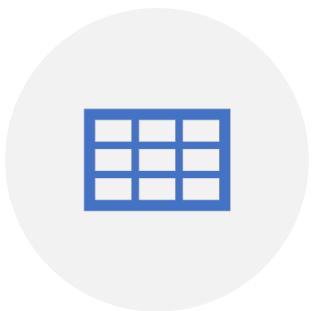
Data Mesh



DATA MESH AIMS TO
EMPOWER TEAMS



WITH DATA AUTONOMY,
ALLOWING THEM



TO DEFINE THEIR DATA
MODELS,



DATA PIPELINES, AND
DATA PRODUCTS.

Data Mesh



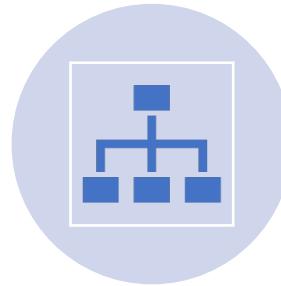
It promotes the idea of



treating data as a product and



advocates for data democratization and



data collaboration across an organization.



databricks Lakehouse Platform

SIMPLE ◦ OPEN ◦ COLLABORATIVE

Data Engineering

BI & SQL
Analytics

Real-time Data
Applications

Data Science
& Machine Learning

Data Management & Governance



DELTA LAKE

Open Data Lake



Structured



Semi-structured



Unstructured



Streaming



Microsoft
Azure



Google Cloud

What is Delta Lake?



DELTA LAKE IS AN OPEN-SOURCE STORAGE LAYER



BUILT ON TOP OF A DATA LAKE,



OFTEN USED IN CONJUNCTION



WITH APACHE SPARK.

What is Delta Lake?

Provide Atomicity,
Consistency,
Isolation,
Durability

Transactions,

Data versioning

Data reliability

Improve the
quality

Reliability of data
in a data lake

What is Delta Lake?



Combines the
Scalability



Flexibility of a
data lake



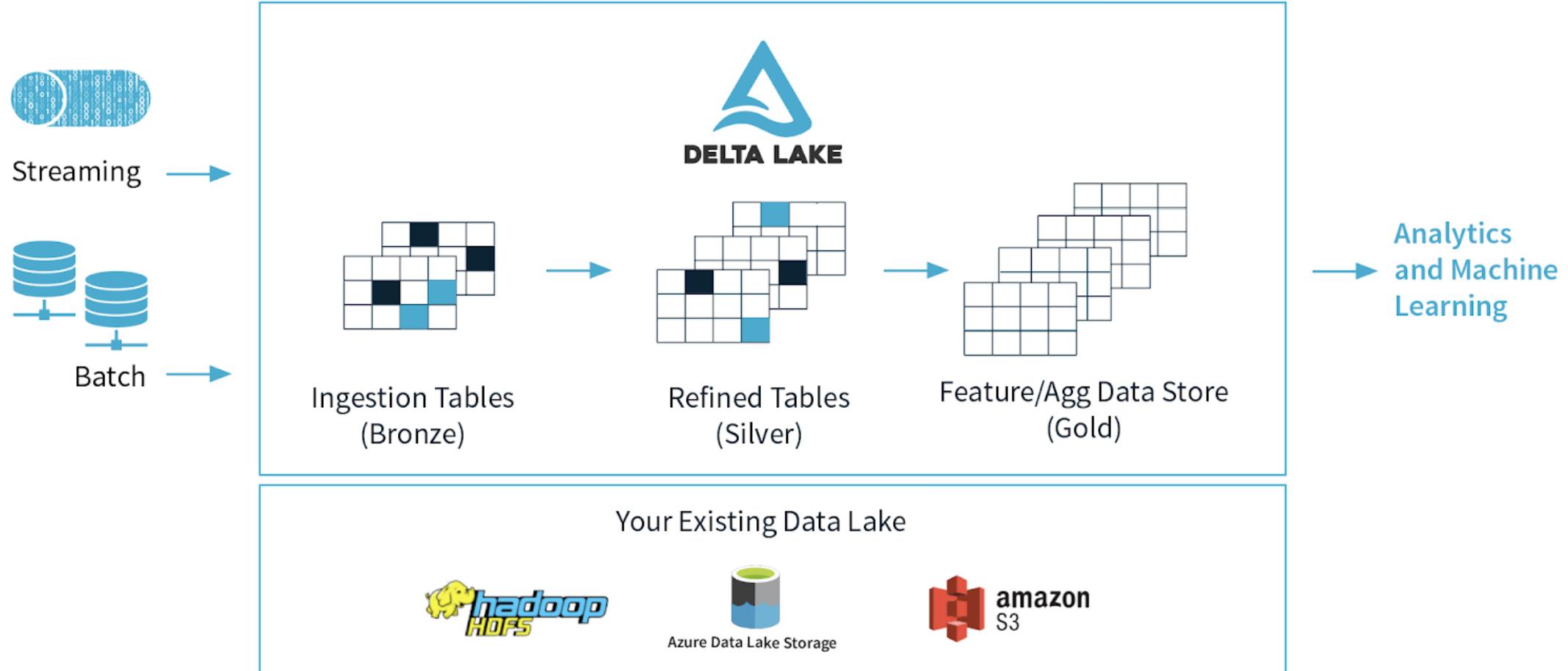
With the
Reliability



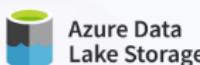
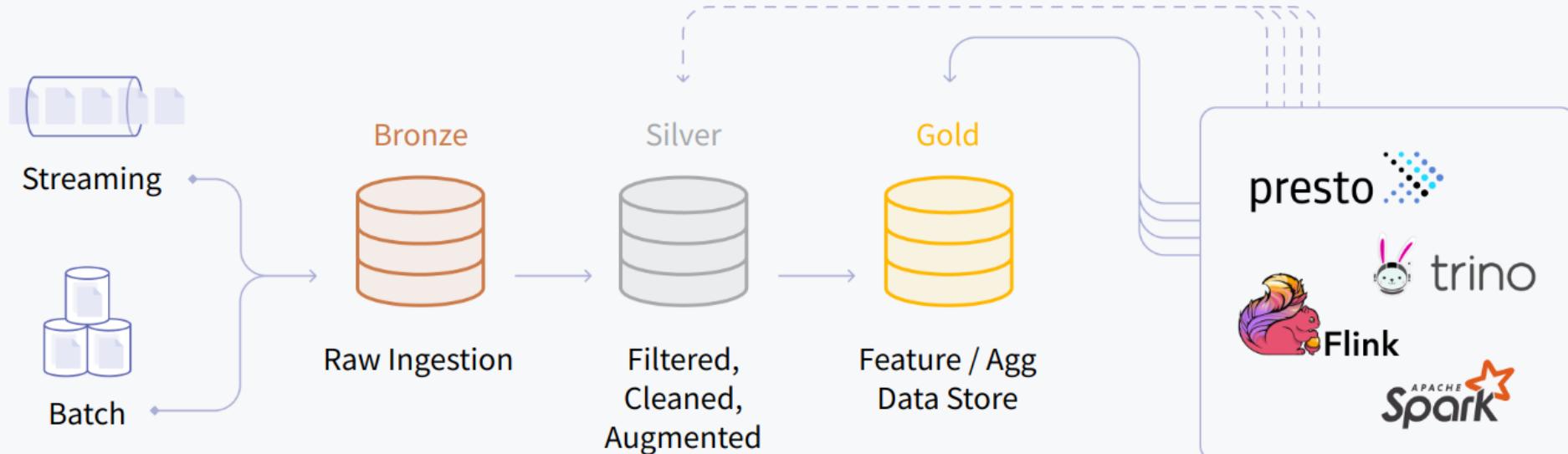
Transactional
capabilities



of a Data
Warehouse.

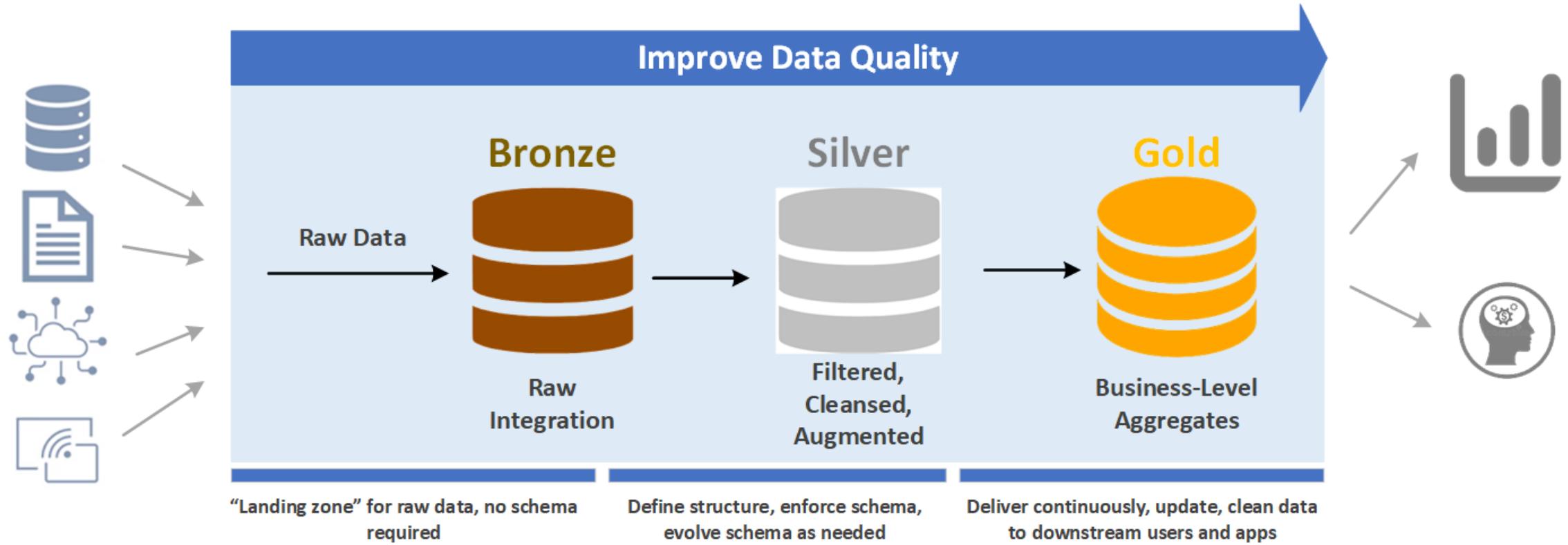


DELTA LAKE

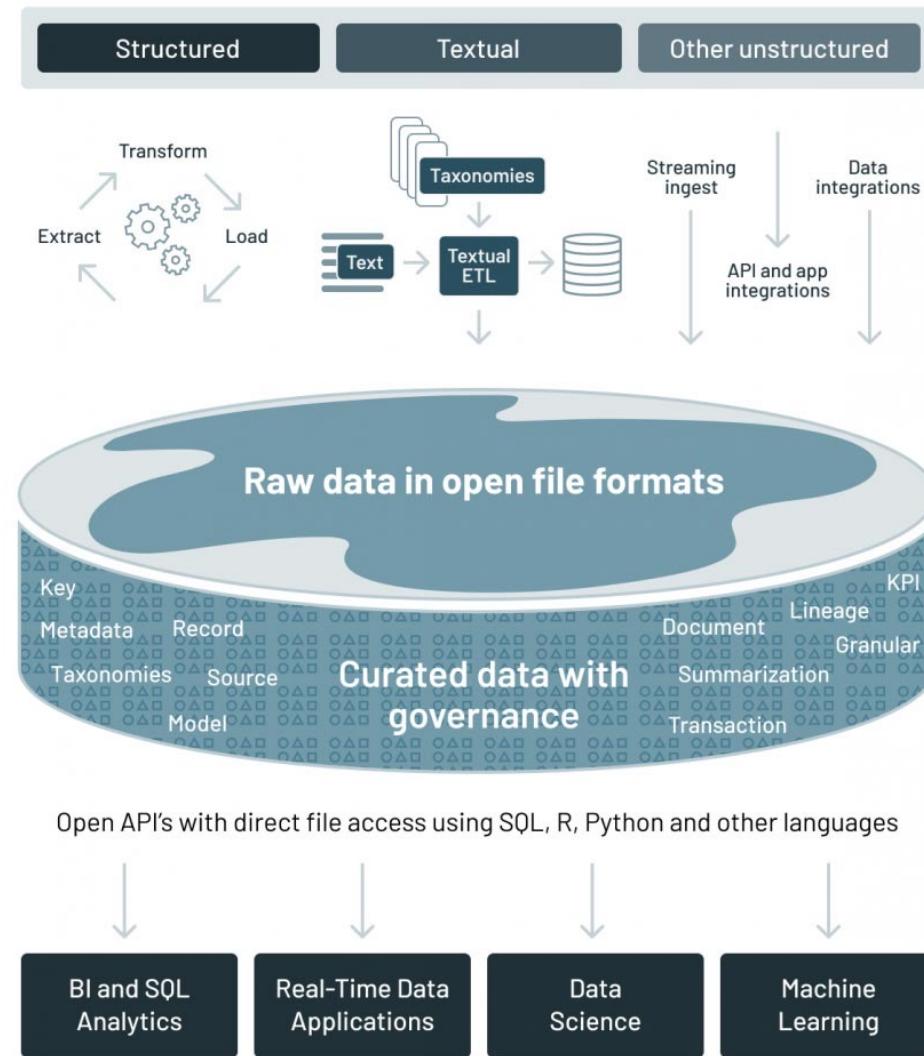


Your Existing Data Lake

BIGDATA

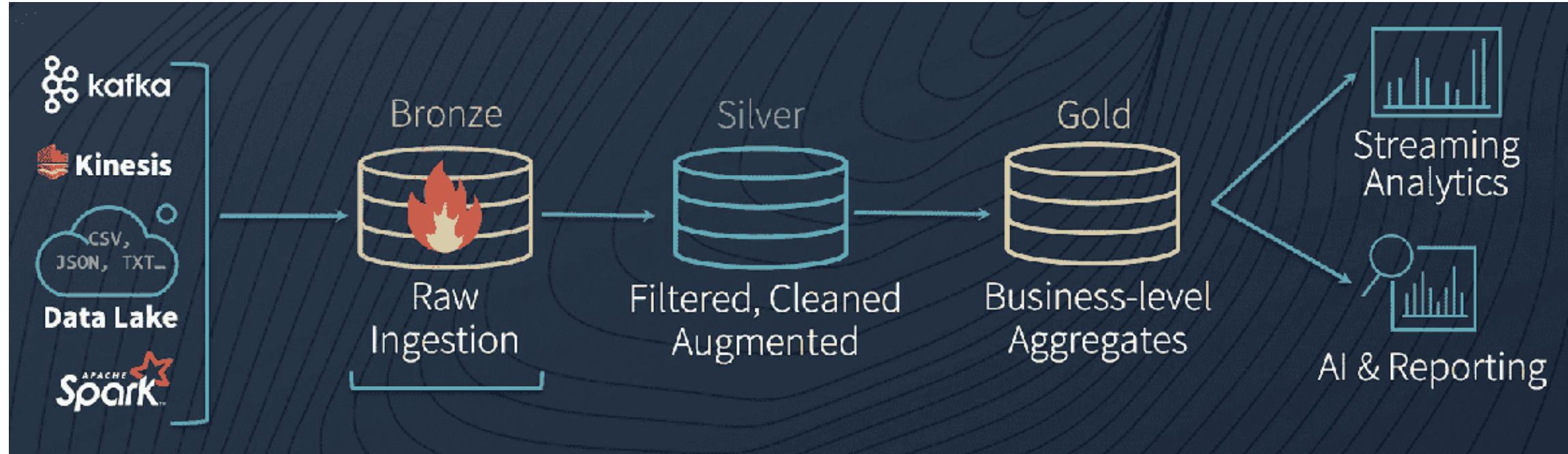


Data Lakehouse



What is Delta Lake?

- Designed to enhance the reliability,
- performance, and
- operational ease of data processing and
- analytics on large-scale datasets.



What is Delta Lake?

Delta Lake is tightly integrated

with the Databricks platform and

offers significant advantages

for managing and analyzing data

within Databricks.

What is Delta Lake?



Serves as a powerful addition



to your data processing toolkit.

How Delta Lake functions within Databricks?

Surendra Panpaliya

ACID Transactions

- DL provides Atomicity, Consistency,
- Isolation, and Durability (ACID) guarantees for data operations.
- Perform updates, inserts, and deletes on your data
- while ensuring data integrity and consistency.

a

Atomicity:
Transactions
are all or
nothing

c

Consistency:
Only valid data
is saved

i

Isolation:
Transactions
do not affect
each other

d

Durability:
Written data
will not be lost



Scalable Metadata Handling

- Delta Lake efficiently manages metadata
- by utilizing a transaction log
- that captures all changes to the data.

Scalable Metadata Handling

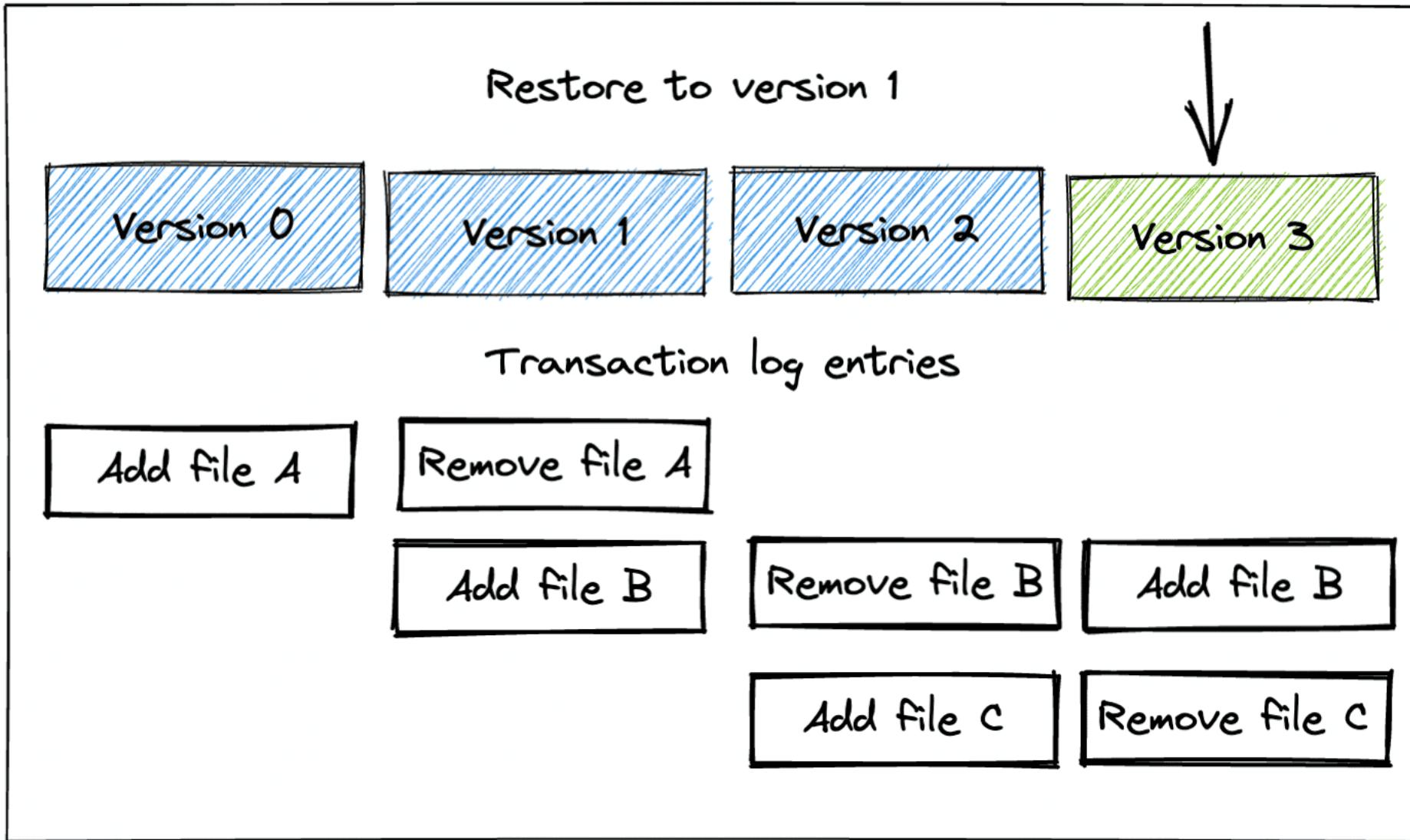
- This scalable approach ensures
- quick access to metadata,
- regardless of the data volume.

Time Travel

- Delta Lake's time travel capabilities
- Allow you to access and query
- historical versions of your data.

Time Travel

- Can query the data
 - as it appeared at different points in time,
 - aiding in debugging, auditing,
 - historical analysis.



Schema Evolution

- Supports schema evolution
- Enabling you to modify the data schema
- without requiring complex ETL processes.
- This flexibility makes it easier
- to adapt to changing business requirements.



Delta Lake schema evolution

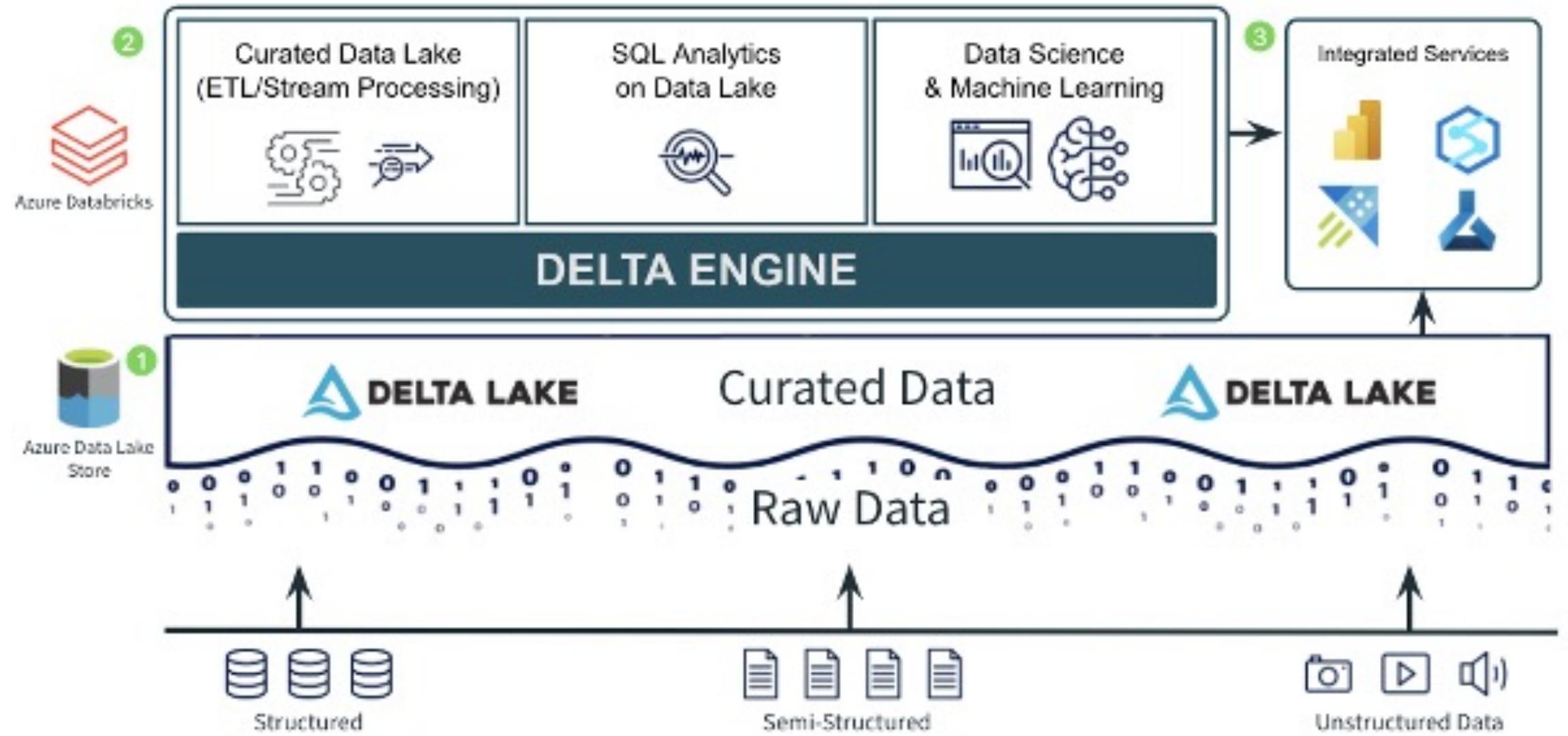
name	age
bob	3
sue	5



name	age	country
bob	3	usa
sue	5	uk

Data Consistency Checks

- Delta Lake automatically performs
- data consistency checks,
- including data integrity,
- schema validation, and statistics validation.
- This ensures that your data is accurate and reliable.



Bronze



Raw Ingestion
and History



DELTA LAKE

Silver



Filtered, Cleaned,
Augmented

Gold



Business-level
Aggregates

Curated Data

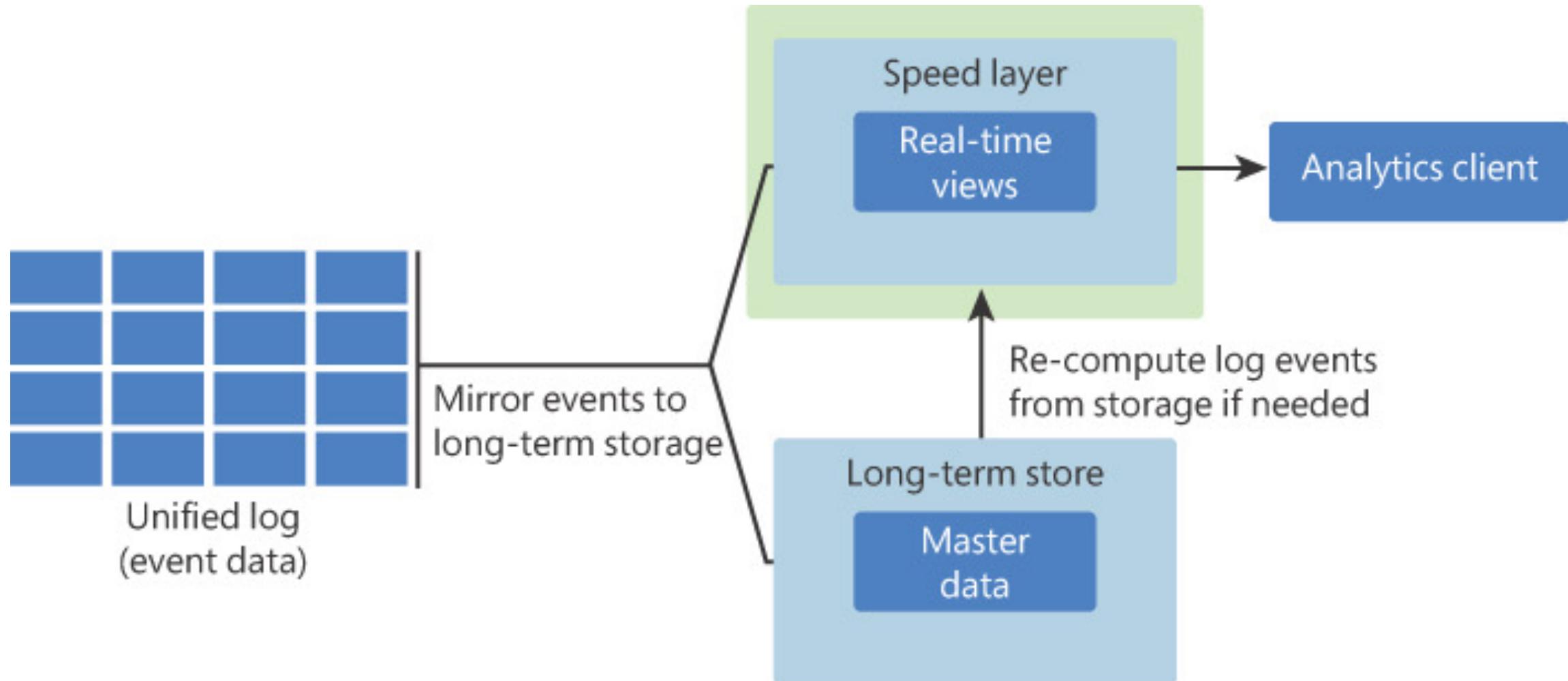


Improved Performance

- Delta Lake optimizes data storage
- by using advanced indexing and
- data skipping techniques,
- leading to improved query performance and
- reduced I/O overhead.

Unified Batch and Streaming

- Seamlessly integrate
- batch and structured streaming workloads using Delta Lake.
- Write streaming data into Delta Lake tables
- query them using both batch and streaming processing.



Integration with Databricks

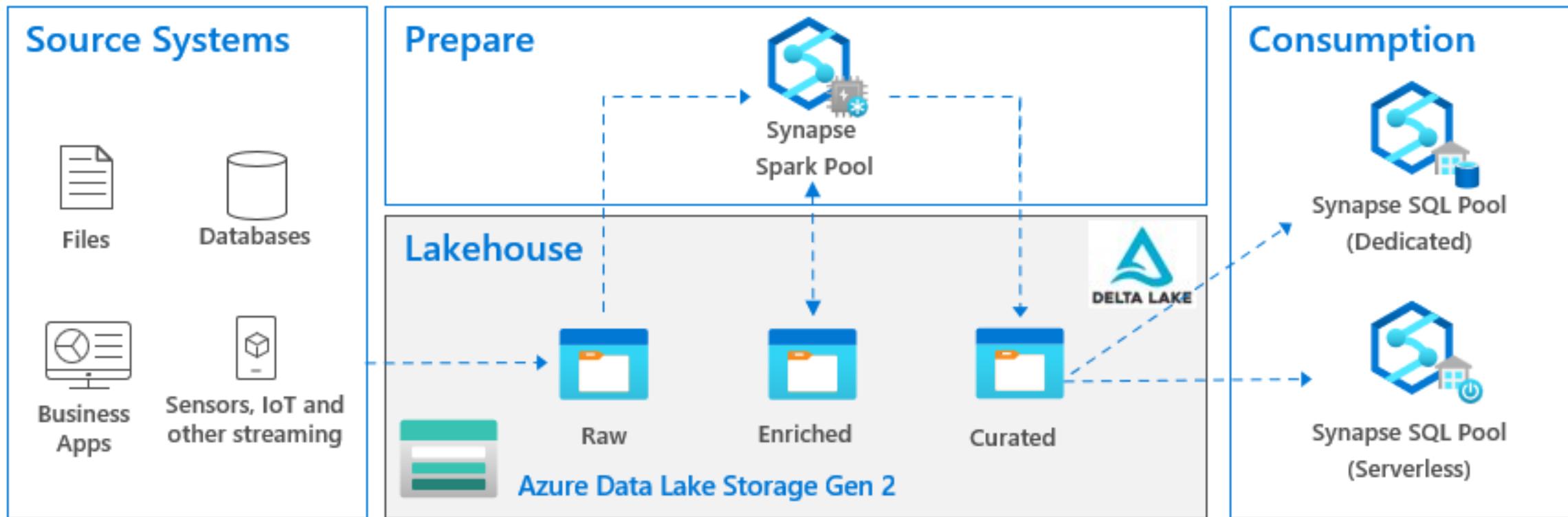
- Delta Lake is closely integrated with Databricks,
- making it straightforward
- to create, manage, and
- analyze Delta Lake tables
- directly within the Databricks environment.

Integration with Databricks

- Databricks provides built-in commands and
- features to work with
- Delta Lake efficiently.

	Delta Lake	Lake ETLs
Lock-in	High. Need to change ingestion and query interfaces to Delta, no support for reliable concurrent writes.	Low. No change in interfaces and no proprietary metadata Lake ETL vendor can be replaced with home-grown ETL.
Ingestion Performance	Low. ACID transactions and indexes.	High. Append-only writes.
Entry barrier	High. Requires non-trivial expertise in Spark coding	Low Visual interface and SQL
Ease-of-use	Medium. ACID operations replace ETLs but all ingestion and query interfaces need to be migrated to Delta. Delta requires a DBA for operations like Vacuum and Optimize	Depends on ETL platform Upsolver offers a turn-key solution, automating ETLs.

DATA WAREHOUSE	vs.	DATA LAKE
structured, processed	DATA	structured / semi-structured / unstructured, raw
schema-on-write	PROCESSING	schema-on-read
expensive for large data volumes	STORAGE	designed for low-cost storage
less agile, fixed configuration	AGILITY	highly agile, configure and reconfigure as needed
mature	SECURITY	maturing
business professionals	USERS	data scientists et. al.



Summary

- Delta Lake in Databricks enhances your ability
- to build reliable and efficient data pipelines,
- enables historical analysis, and
- simplifies the process of managing large-scale data lakes.

Summary

- It's a valuable tool
- for organizations seeking
- to leverage the full potential of their data
- within the Databricks platform.

Data Warehouse	A large, structured repository of integrated data from various sources, used for complex querying and historical analysis	<pre> graph LR SD[Structured Data] --> DW[Data Warehouse] DW --> A[Analysis] DW --> R[Report] DW --> V[Visualize] </pre>
Data Lake	A more focused, department-specific subset of a data warehouse providing quick data retrieval and analysis	<pre> graph LR RD[Raw Data] --> DL[Data Lake] DL --> A DL --> R DL --> V </pre>
Data Mart	A vast pool of raw, unstructured data stored in its native format until it's needed for use	<pre> graph LR RD --> DW[Data Warehouse] DW --> DM[Data Mart] DM --> Eng[Eng] DM --> Sales[Sales] DM --> Finance[Finance] Eng --> A Eng --> R Eng --> V Sales --> A Sales --> R Sales --> V Finance --> A Finance --> R Finance --> V </pre>
Delta Lake	An open-source storage layer that brings reliability and ACID transactions to data lakes, unifying batch and streaming data processing	<pre> graph LR RD --> DL[Delta Lake] DL --> EDSL[Existing Data Lake Solution] EDSL --> A EDSL --> R EDSL --> V </pre>
Data Pipeline	A process that moves and transforms data from one system to another, often used to populate data warehouses and data lakes	<pre> graph LR D[Data] --> P1(()) P1 --> P2(()) P2 --> P3(()) P3 --> A[Analysis] P3 --> R[Report] P3 --> V[Visualize] </pre>
Data Mesh	An architectural and organizational approach where data ownership and delivery are decentralized across domain-specific, cross-functional teams	<pre> graph TD subgraph DG [Data Governance] subgraph F [Finance] FInfra[Data Infra] FGov[Data Governance] FGov --> FInfra end subgraph S [Sales] SInfra[Data Infra] SGov[Data Governance] SGov --> SInfra end subgraph E [Engineer] EInfra[Data Infra] EGov[Data Governance] EGov --> EInfra end subgraph O [Operation] OInfra[Data Infra] OGov[Data Governance] OGov --> OInfra end end </pre>

Data Warehouse vs Data Mart

Data Warehouse:

- A large, structured repository of integrated data
- from various sources, used for complex querying and historical analysis.

Data Mart:

- A more focused, department-specific subset of a data warehouse
- providing quick data retrieval and analysis.

Data Lake vs Delta Lake

- Data Lake:
 - A vast pool of raw, unstructured data
 - stored in its native format until it's needed for use.
- Delta Lake: An open-source storage layer
 - that brings reliability and ACID transactions to data lakes,
 - unifying batch, and streaming data processing.

Data Pipeline vs Data Mesh

- Data Pipeline
 - A process that moves and transforms data from one system to another,
 - often used to populate data warehouses and data lakes.
- Data Mesh:
 - An architectural and organizational approach
 - where data ownership and delivery are decentralized
 - across domain-specific, cross-functional teams.

Delta Lake Hands On Session

Surendra Panpaliya

Delta Lake operations

- Perform a variety of operations on Delta Lake tables
- To manage, query, and transform your data.
- Delta Lake offers capabilities for data manipulation,
- schema evolution, time travel

Create a Delta Lake Table

- create a new Delta Lake table from existing data
- Such as Parquet files or structured data.

```
from pyspark.sql import SparkSession
```

```
# Create a Spark session
spark = SparkSession.builder.appName("DeltaLakeExample").getOrCreate()
```

Create a Delta Lake Table

```
# Read data from a source (e.g., Parquet files)
data_df = spark.read.parquet("/path/to/source/data")

# Write data to a Delta Lake table
data_df.write.format("delta").save("/path/to/delta-table")
```