



# 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

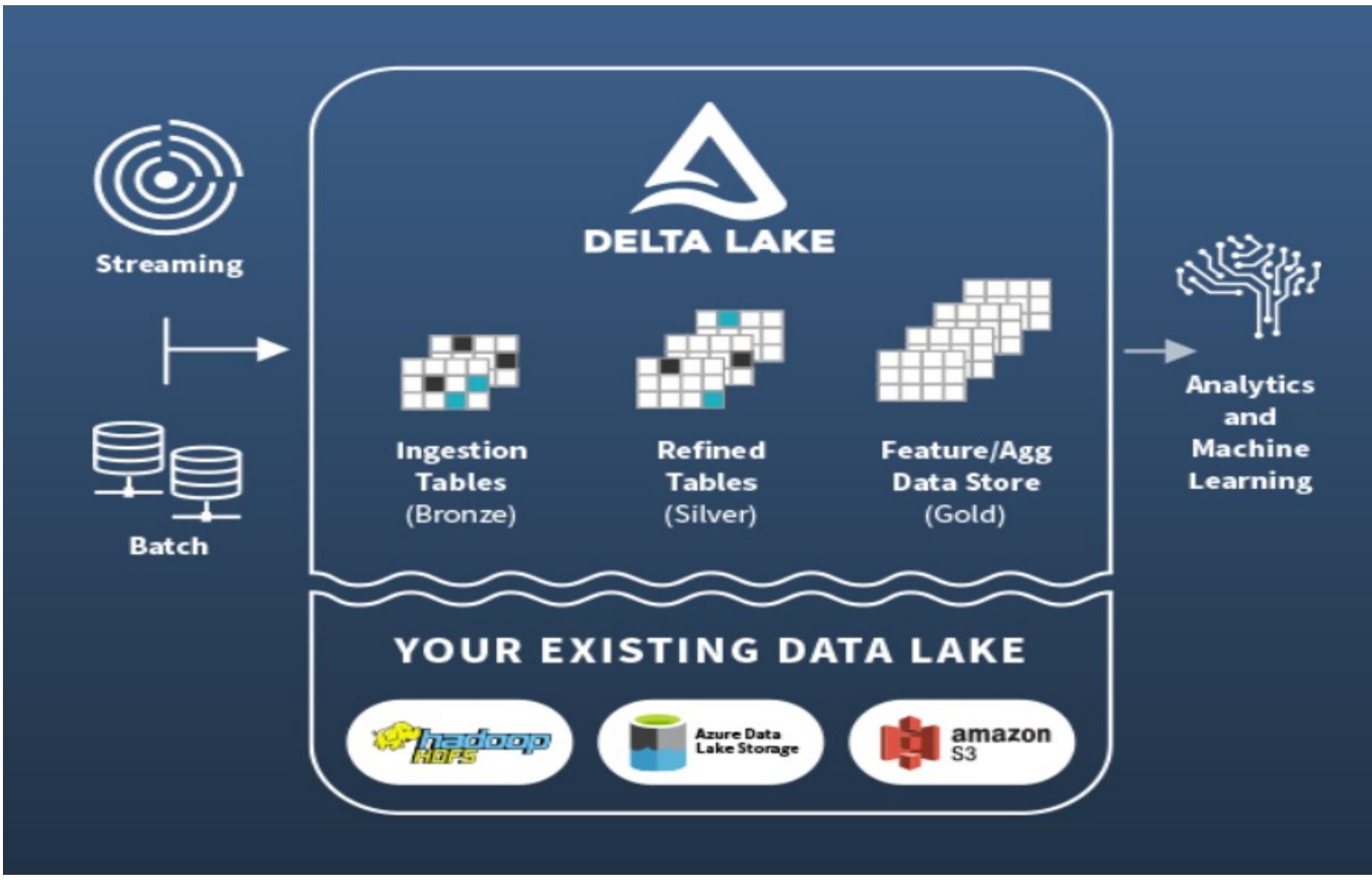
id
0
1
2

Version 1

id
0
1
2
8
9
10

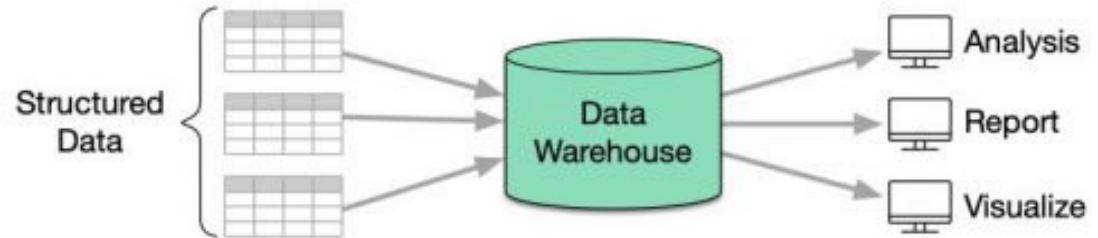
Version 2

id
55
66
77



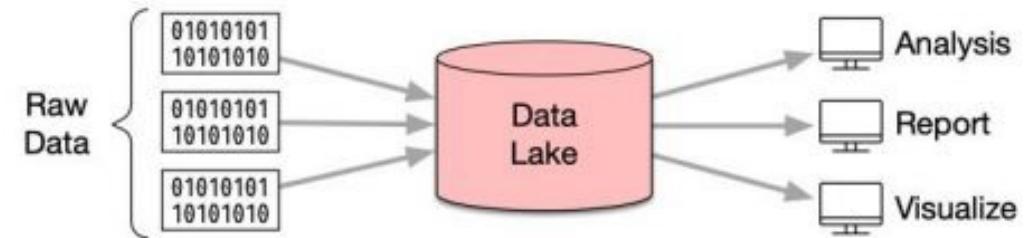
## Data Warehouse

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



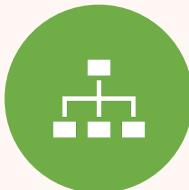
A central repository



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.

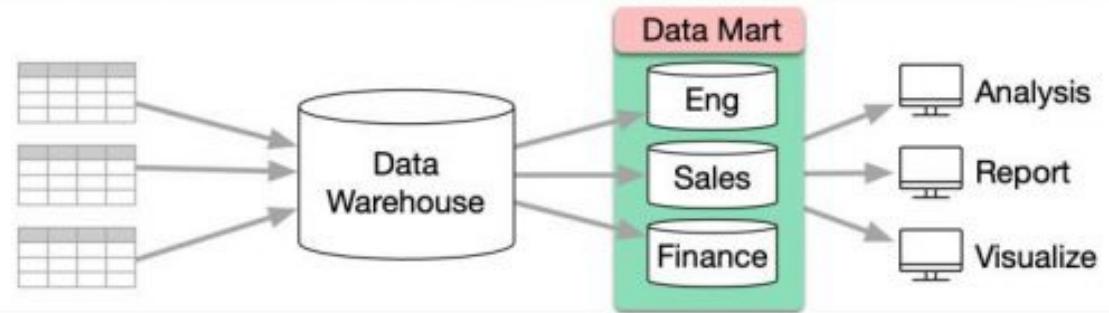


VS



## 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



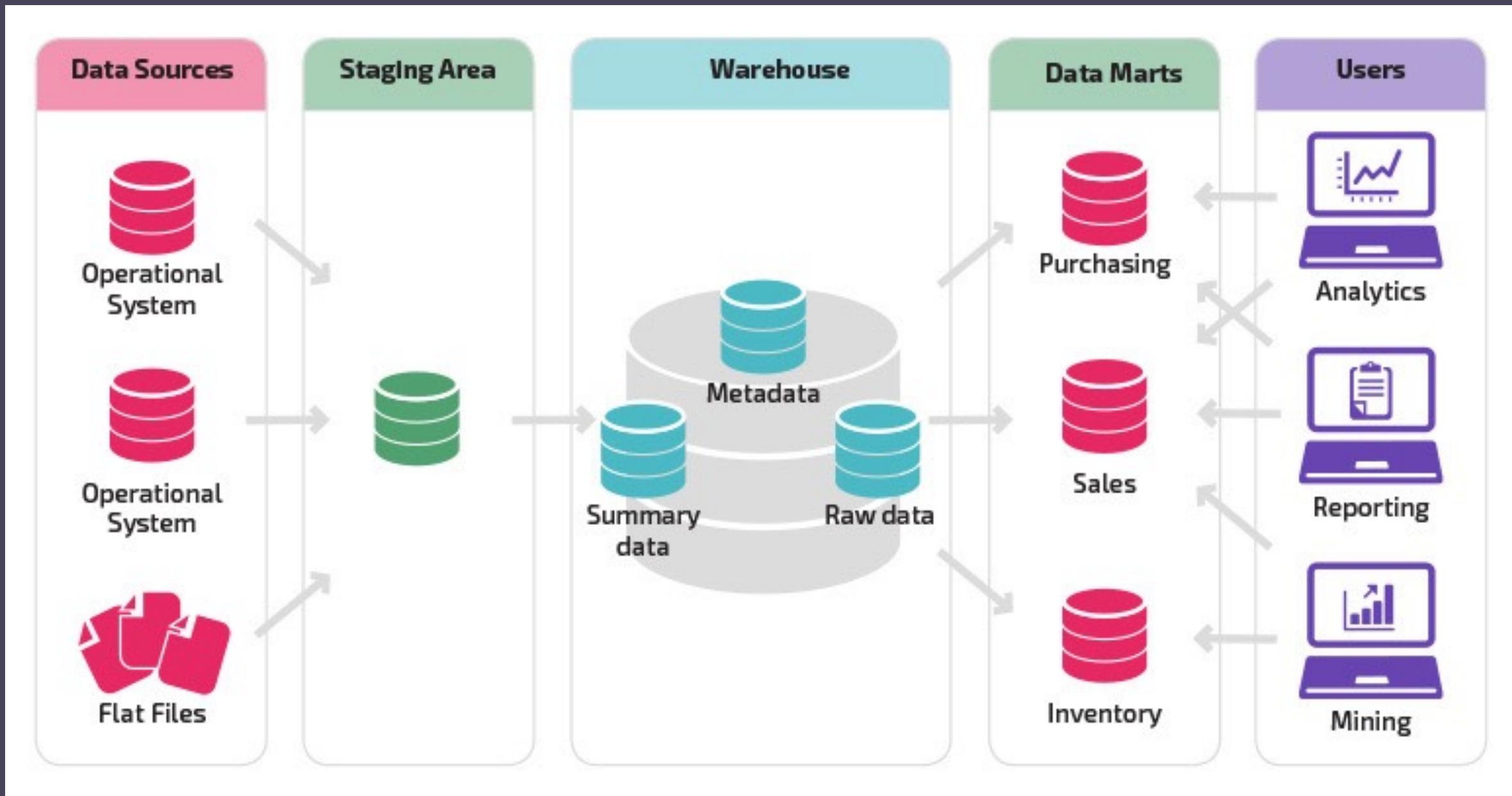
subset of a data warehouse



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

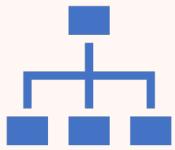


compared to the broader

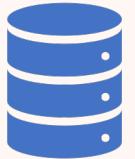


data warehouse.

# Data Mart



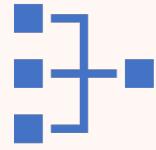
Can be  
structured



Similar to a data  
warehouse or



May use  
different



Schemas and  
models.

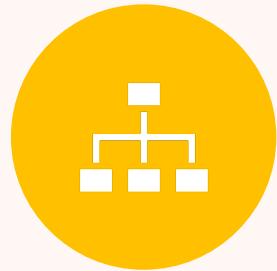
# What is Data Lake?



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?



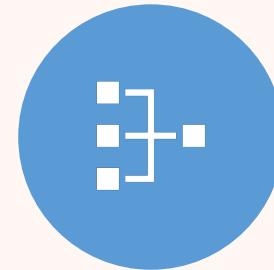
ALLOWS  
ORGANIZATIONS



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?



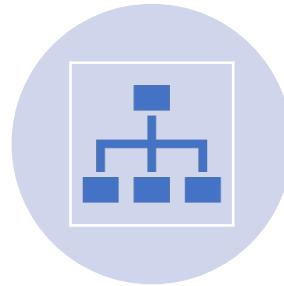
The data lake architecture typically



follows a schema-on-read approach,



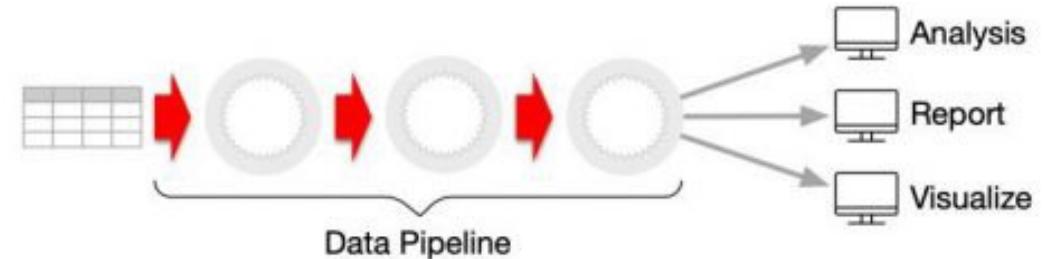
where data is transformed and



structured at the time of analysis.

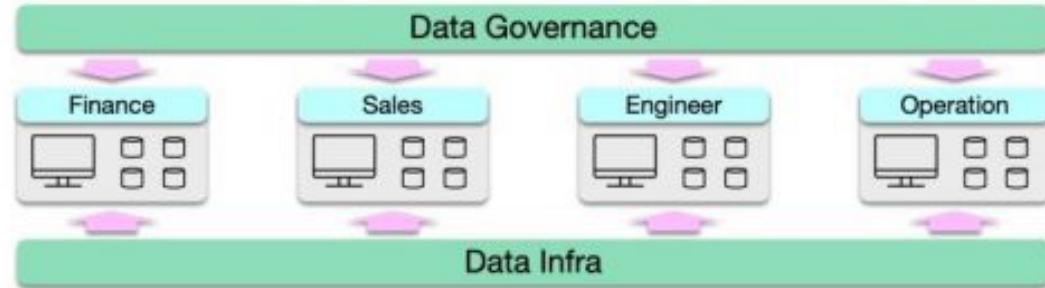
## Data Pipeline

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

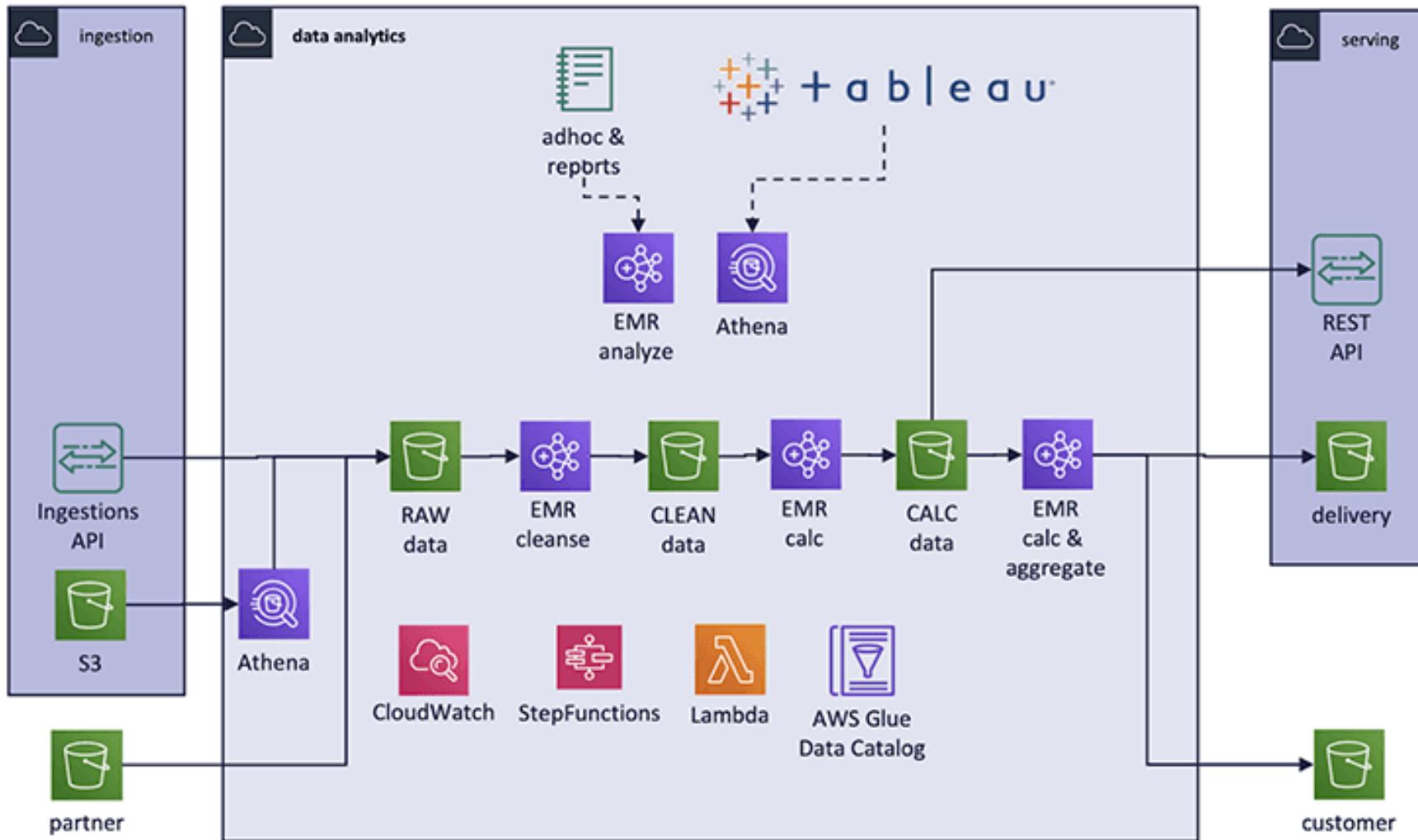


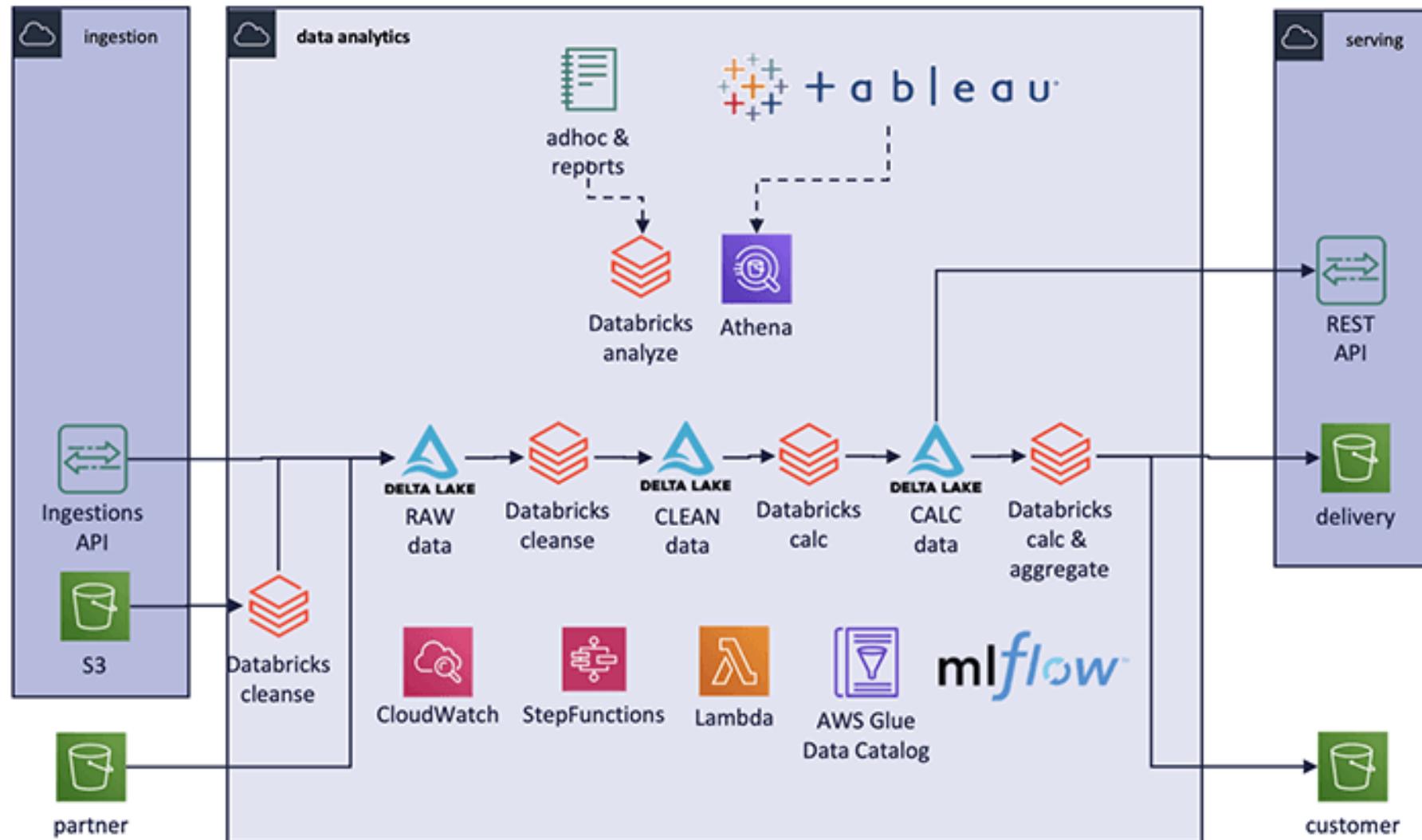
Target storage or processing system.



# High level architecture & tech stack

How our world looked before Databricks





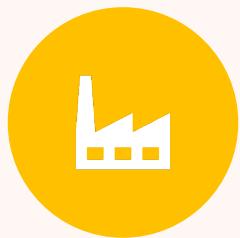
# Data Pipeline



ENABLES THE  
SMOOTH FLOW OF  
DATA



FROM SOURCE  
SYSTEMS



TO DATA  
WAREHOUSES,



DATA LAKES, OR



OTHER  
DESTINATIONS.

# Data Pipeline



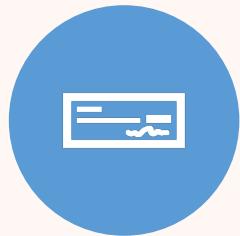
DATA PIPELINES  
HANDLE TASKS



DATA INGESTION,



DATA  
TRANSFORMATION,



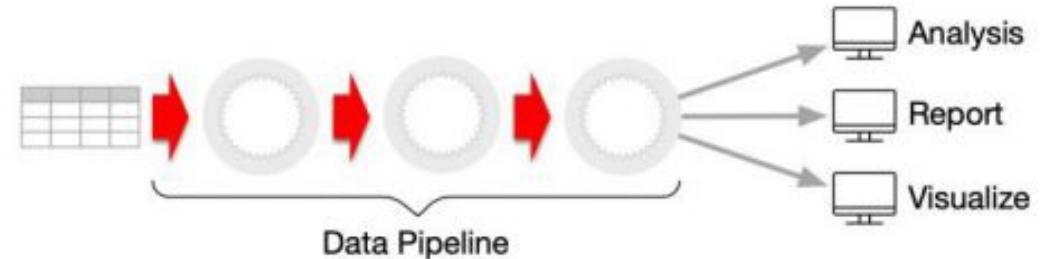
DATA QUALITY  
CHECKS,



DATA DELIVERY.

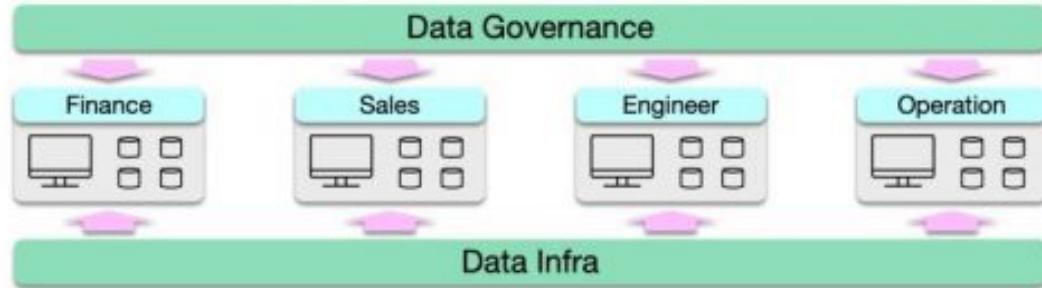
## Data Pipeline

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.

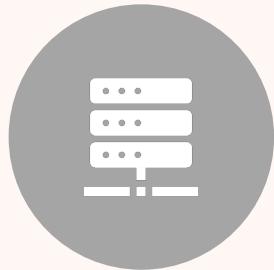
# Data Mesh Architecture



# Data Mesh



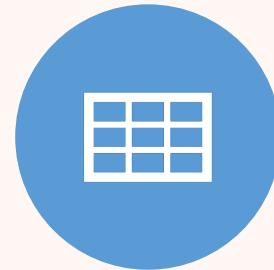
IT EMPHASIZES  
DOMAIN-ORIENTED,



SELF-SERVE DATA  
PRODUCTS,



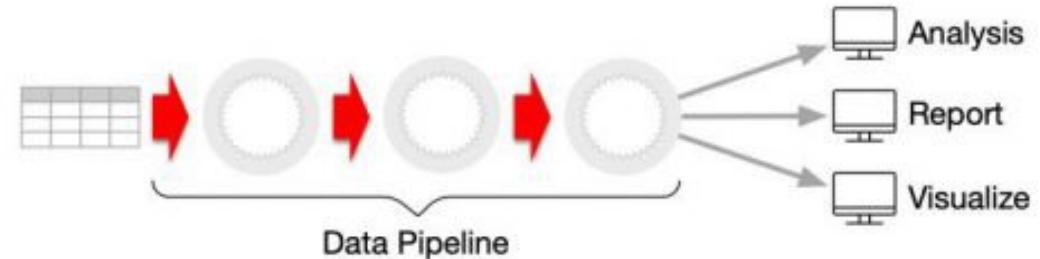
WHERE CROSS-  
FUNCTIONAL TEAMS



TAKE OWNERSHIP OF  
THEIR DATA DOMAINS.

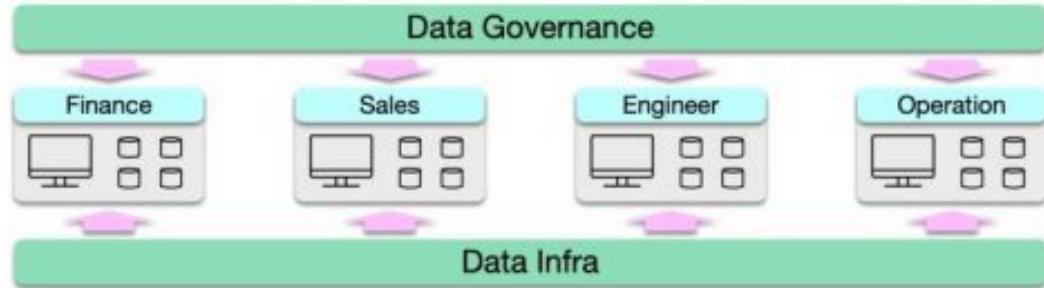
## Data Pipeline

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



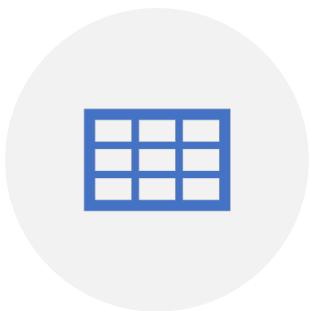
# Data Mesh



DATA MESH AIMS TO  
EMPOWER TEAMS



WITH DATA AUTONOMY,  
ALLOWING THEM



TO DEFINE THEIR DATA  
MODELS,



DATA PIPELINES, AND  
DATA PRODUCTS.

# Data Mesh



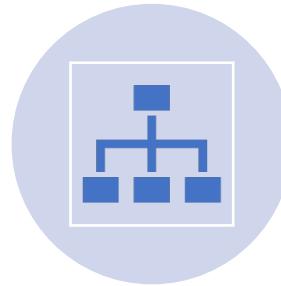
It promotes the idea of



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



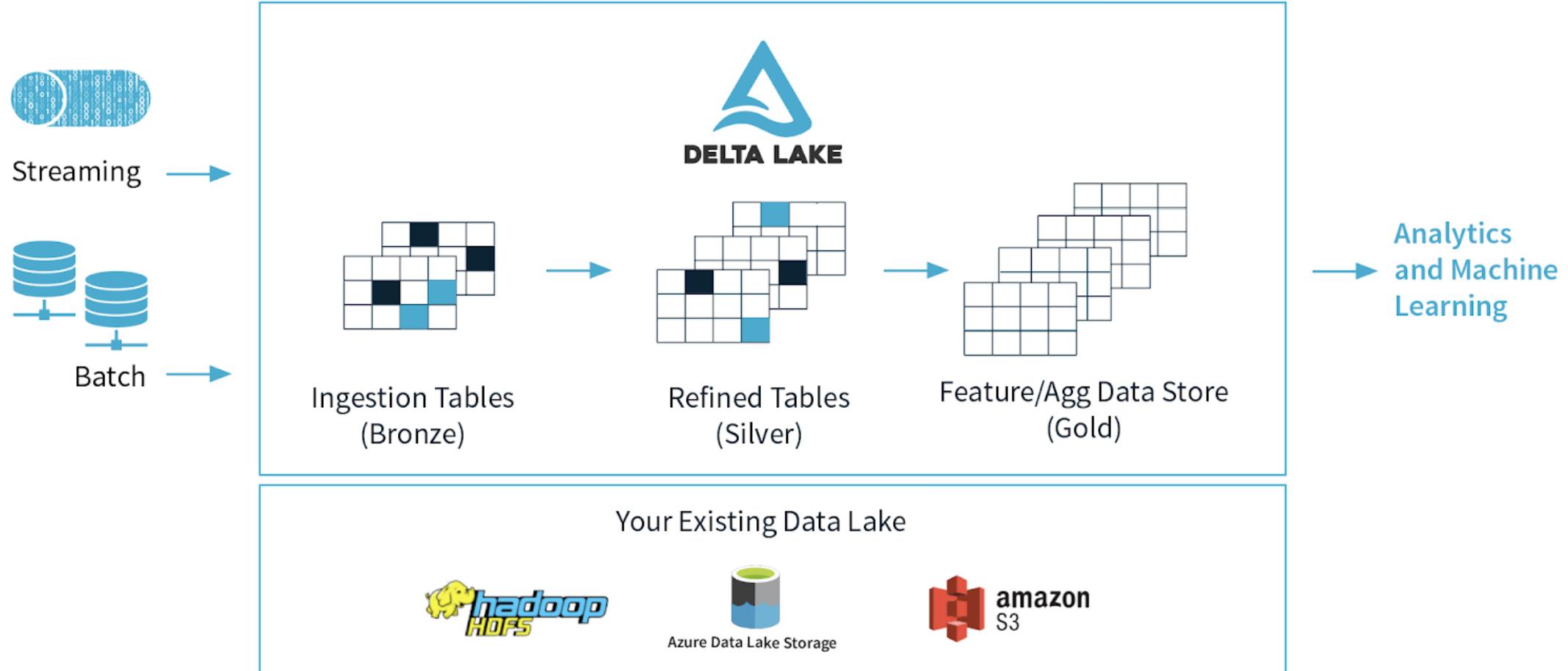
With the  
Reliability



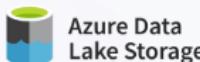
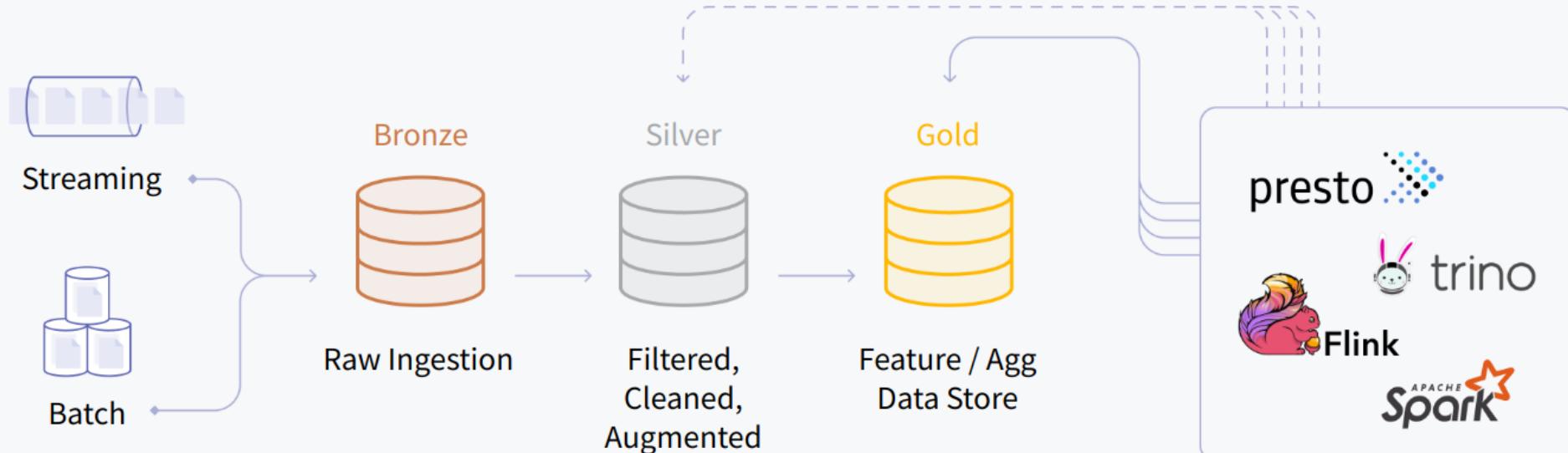
Transactional  
capabilities



of a Data  
Warehouse.

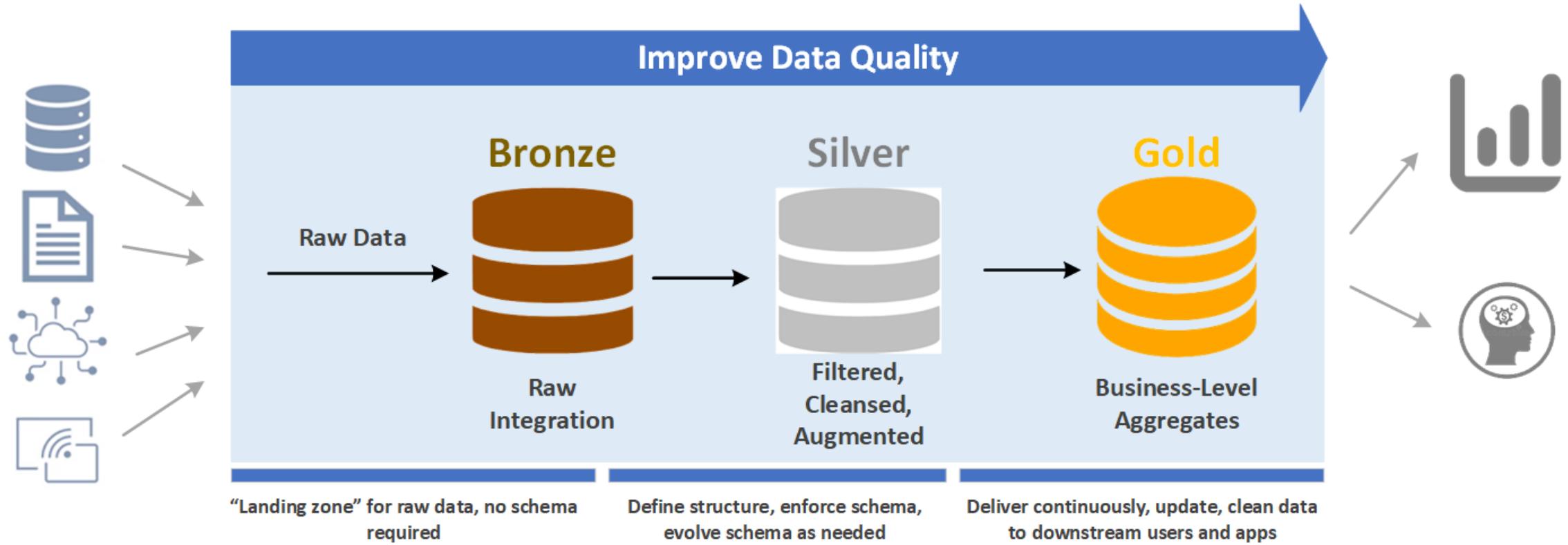


# DELTA LAKE

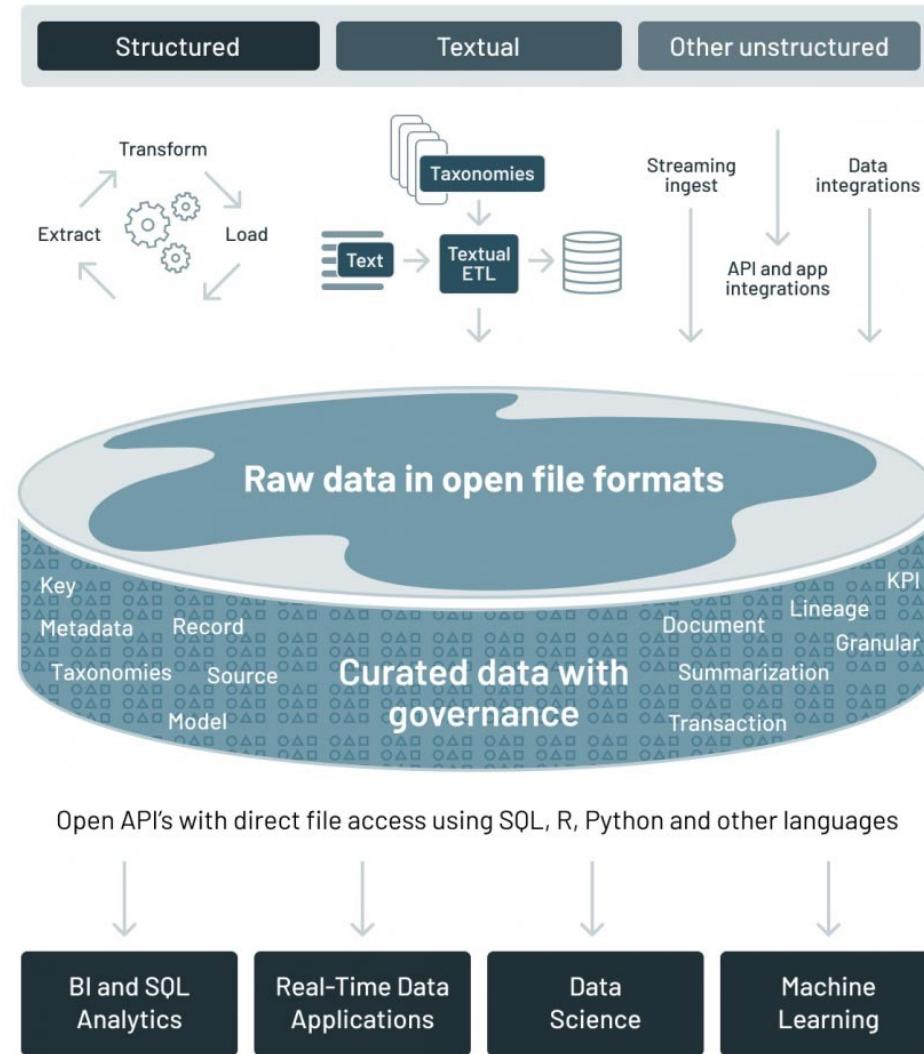


Your Existing Data Lake

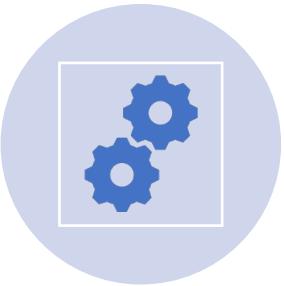
BIGDATA



# Data Lakehouse



# What is Delta Lake?



Designed to enhance the reliability,



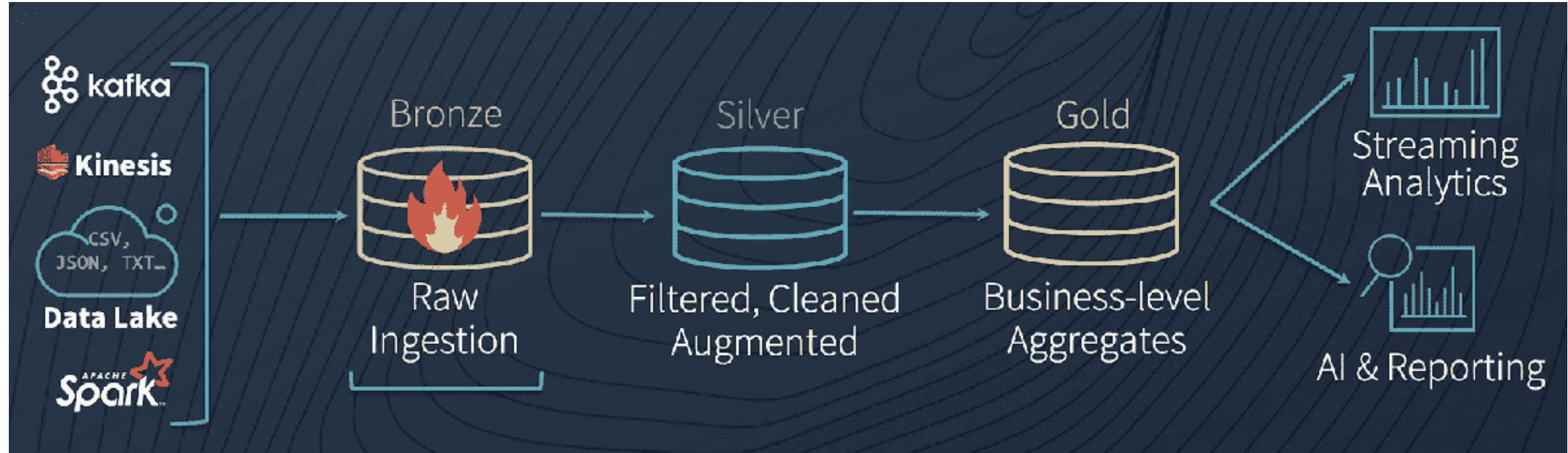
Operational ease of data processing



Performance, 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.

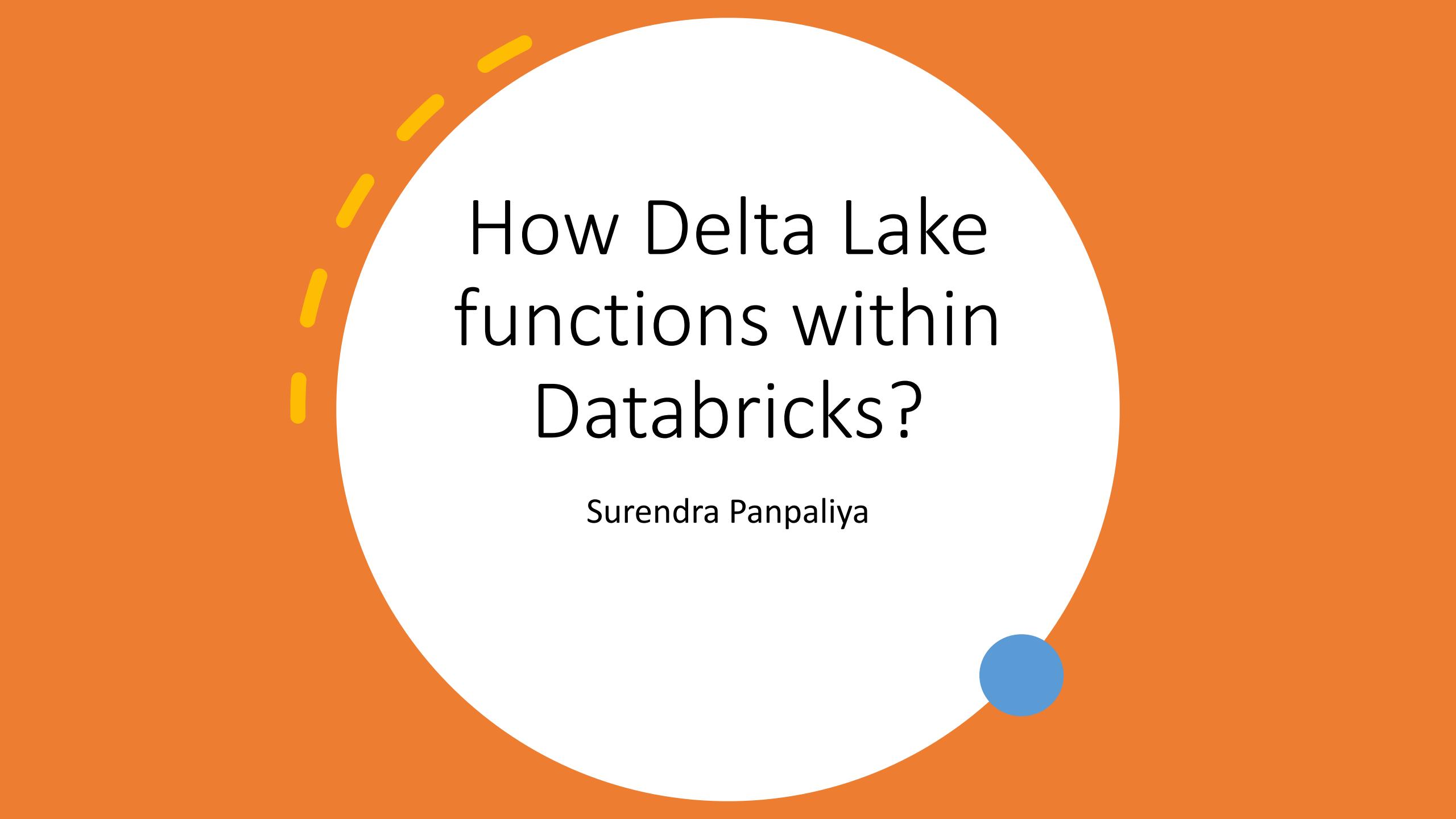
# What is Delta Lake?



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



Utilizing a transaction log



Captures all changes to the data.

# Scalable Metadata Handling

Ensures Quick access

to Metadata,

Regardless of

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,



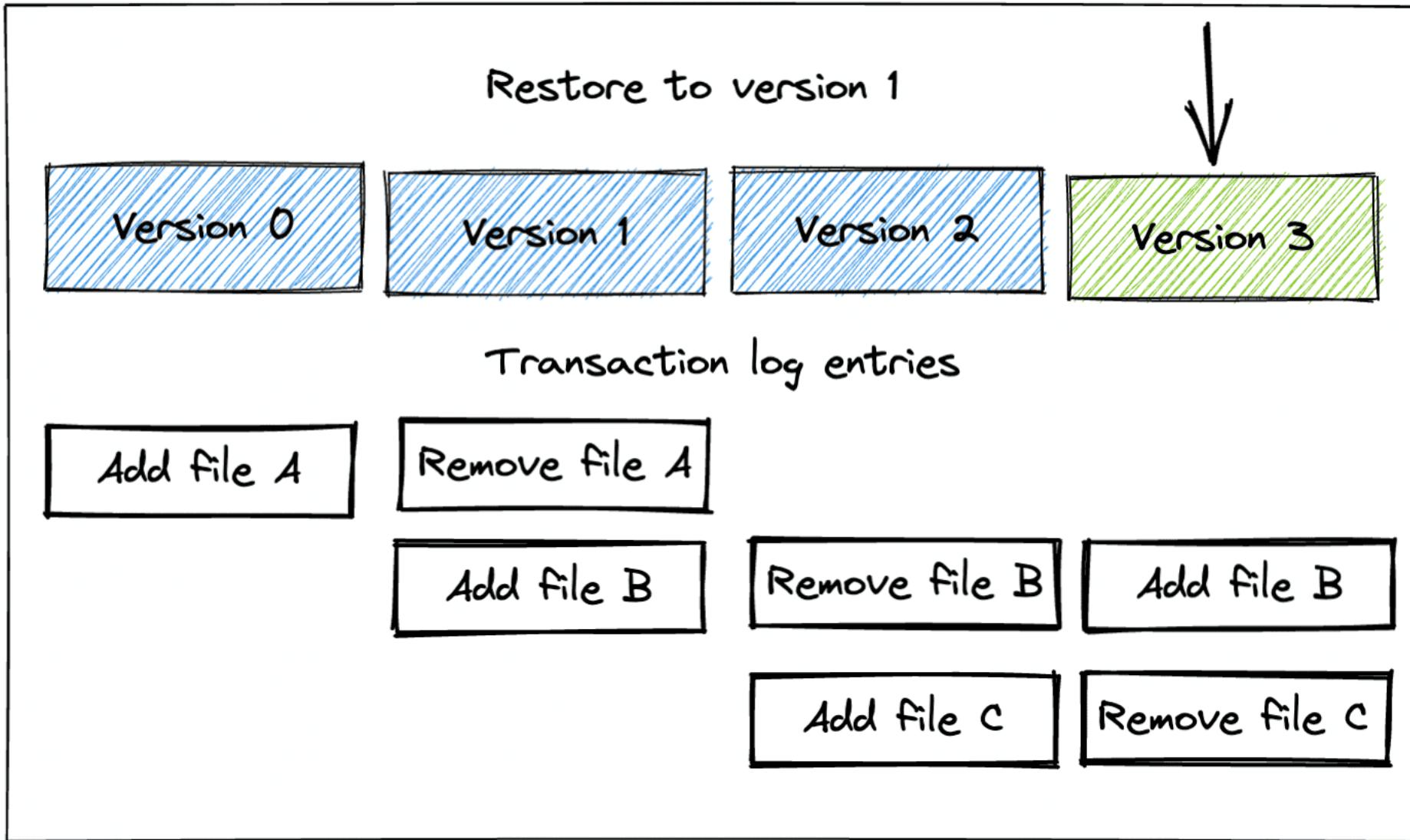
Aiding in Debugging,



Auditing



Historical analysis.



# Schema Evolution



Supports schema evolution



Enabling you to modify the data schema



Without requiring complex ETL processes.



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

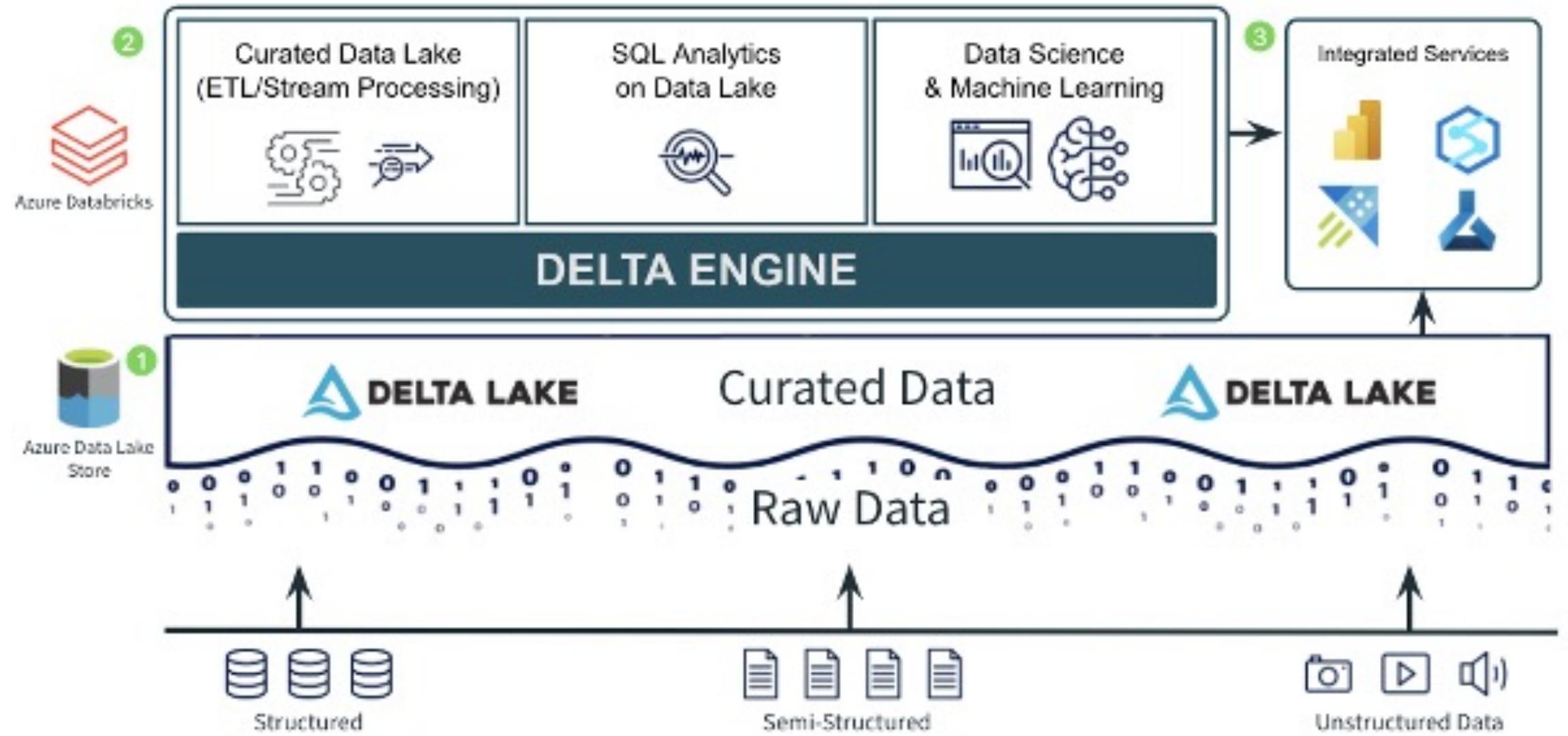
Automatically performs

Data consistency checks,

Including Data Integrity,

Schema validation, and Statistics validation.

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



Using Advanced Indexing



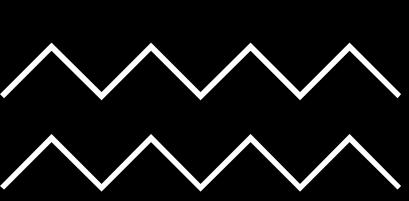
Data Skipping Techniques,



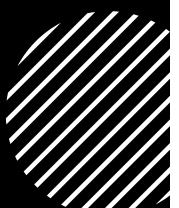
Leading to improved query performance



Reduced I/O overhead.



# Unified Batch and Streaming



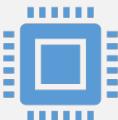
Seamlessly integrate



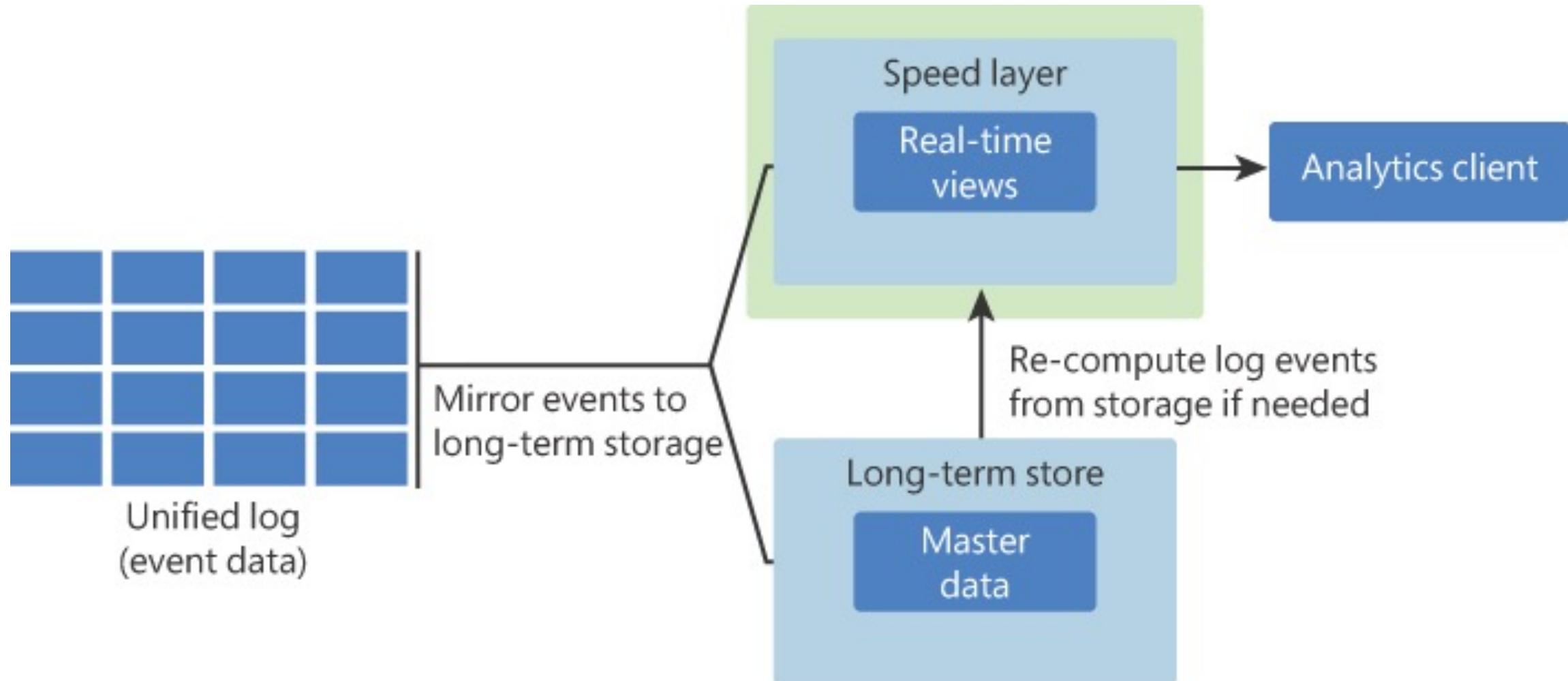
batch and structured streaming workloads using Delta Lake.



Write streaming data into Delta Lake tables

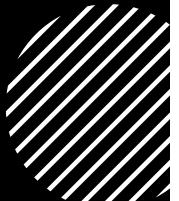


query them using both batch and streaming processing.





# Integration with Databricks



Delta Lake is closely integrated with Databricks,



making it straightforward



to create, manage, and

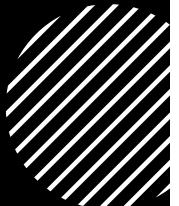


analyze Delta Lake tables



directly within the Databricks environment.

# Integration with Databricks



Databricks provides built-in  
commands and



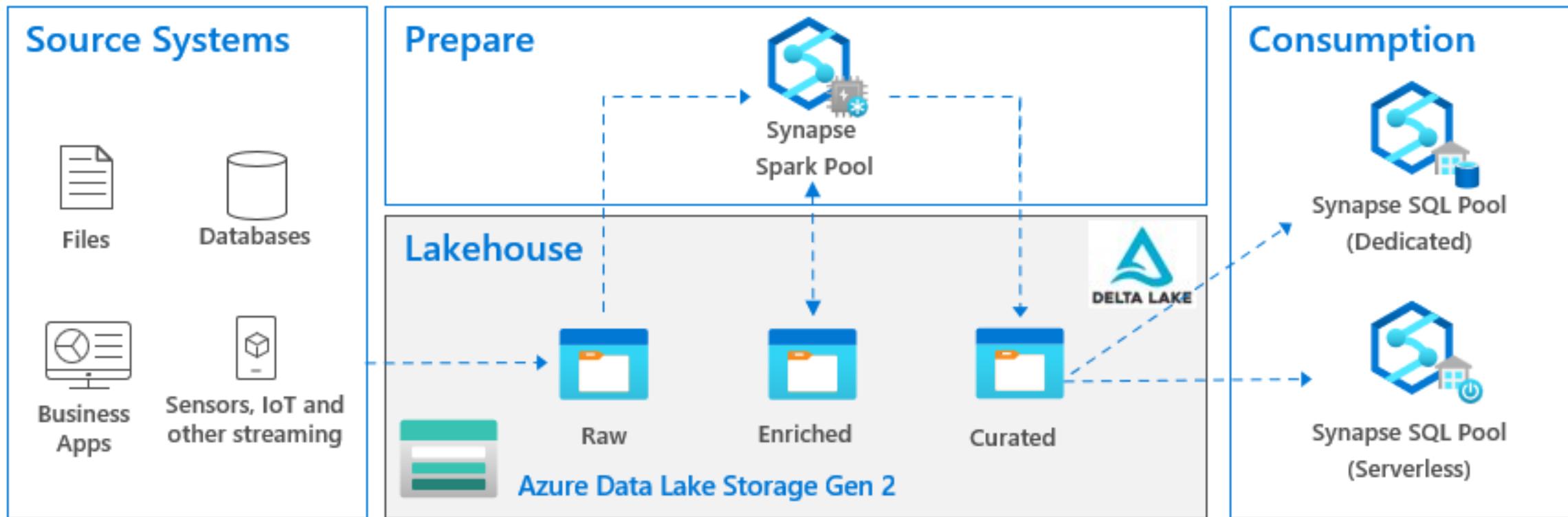
features to work with



Delta Lake efficiently.

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

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



# Summary

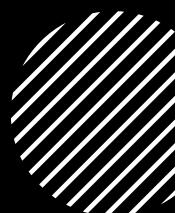
Delta Lake in Databricks enhances your ability

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.

<b>Data Warehouse</b>	A large, structured repository of integrated data from various sources, used for complex querying and historical analysis	<pre> graph LR     SD[Structured Data] --&gt; DW[Data Warehouse]     DW --&gt; A[Analysis]     DW --&gt; R[Report]     DW --&gt; V[Visualize]   </pre>
<b>Data Lake</b>	A more focused, department-specific subset of a data warehouse providing quick data retrieval and analysis	<pre> graph LR     RD[Raw Data] --&gt; DL[Data Lake]     DL --&gt; A     DL --&gt; R     DL --&gt; V   </pre>
<b>Data Mart</b>	A vast pool of raw, unstructured data stored in its native format until it's needed for use	<pre> graph LR     RD --&gt; DW[Data Warehouse]     DW --&gt; DM[Data Mart]     DM --&gt; Eng[Eng]     DM --&gt; Sales[Sales]     DM --&gt; Finance[Finance]     Eng --&gt; A     Eng --&gt; R     Eng --&gt; V     Sales --&gt; A     Sales --&gt; R     Sales --&gt; V     Finance --&gt; A     Finance --&gt; R     Finance --&gt; V   </pre>
<b>Delta Lake</b>	An open-source storage layer that brings reliability and ACID transactions to data lakes, unifying batch and streaming data processing	<pre> graph LR     RD --&gt; DL[Delta Lake]     DL --&gt; A     DL --&gt; R     DL --&gt; V   </pre> <p style="text-align: center;">Existing Data Lake Solution</p>
<b>Data Pipeline</b>	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] --&gt; P1(( ))     P1 --&gt; P2(( ))     P2 --&gt; P3(( ))     P3 --&gt; A[Analysis]     P3 --&gt; R[Report]     P3 --&gt; V[Visualize]   </pre> <p style="text-align: center;">Data Pipeline</p>
<b>Data Mesh</b>	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 --&gt; FInfra         end         subgraph S [Sales]             SInfra[Data Infra]             SGov[Data Governance]             SGov --&gt; SInfra         end         subgraph E [Engineer]             EInfra[Data Infra]             EGov[Data Governance]             EGov --&gt; EInfra         end         subgraph O [Operation]             OInfra[Data Infra]             OGov[Data Governance]             OGov --&gt; 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")
```

# Z-Ordering



Z-Ordering is a technique



to optimize query performance



by co-locating related data



in the same files.

# Z-Ordering



You can use the OPTIMIZE command



to optimize the table's layout



based on specified columns



`delta_table.optimize("column_name")`

# Cache Metadata



Caching the metadata of a



Delta Lake table can



improve query performance.



`delta_table.cache()`

# Update Statistics



Delta Lake maintains statistics



Help improve query optimization.



Periodically updating these statistics



can lead to better query plans



`delta_table.generate("STATISTICS")`

# Liquid Clustering in Delta Lake Databricks



Available in Public Preview in Databricks Runtime 13.2 and above.



Delta Lake liquid clustering



Replaces table partitioning &



ZORDER to simplify data layout decisions



Optimize query performance.

# Liquid Clustering in Delta Lake Databricks



Liquid clustering provides flexibility



to redefine clustering keys



without rewriting existing data,



allowing data layout



to evolve alongside analytic needs over time.

# Liquid Clustering in Delta Lake Databricks

- <https://docs.databricks.com/en/delta/clustering.html>