

iconsulting



Data Lakehouse Platforms

Andrea Carmè

Senior Manager

a.carme@iconulting.biz

Iconsulting



Andrea Carmè

Senior Manager

a.carme@iconulting.biz



Areas of interest:

Tech:

- Data Lakehouse
- Data Architectures
- Data Modeling
- GenAI Engineering

Iconsulting

Business:

- Fashion & Luxury Industry



ABOUT US

23 years on the cutting edge of innovation.

Our journey began in 2001. We ventured beyond the borders of university research so we could imagine methodologies, algorithms and technologies capable of exploiting the greatest asset of today's market: data.

Today we operate from 4 offices in Bologna, Milan, Rome and Naples, helping big companies and international organizations to make better decisions supported by the passion and talents of our professionals, advisors and data scientists.

We believe that data is an invaluable resource to support the evolution of humanity and we stand alongside our partners to help them have a positive impact on society.

20⁰¹
24

OUR OFFICES

Bologna

Rome

iconsulting

Milan

Naples

Romagna (coming soon)

**Ignite the right decision
through human potential.**

For us, data are firestones to spark knowledge and bring the light in the darkness of arbitrary informations.
We help our clients evolve their relationship with data and gain awareness in decision making.

ABOUT US

Smart people

450 +

Happy customers

200 +

iconsulting

Challenging projects

1500 +

95%

blessed NPS *

* Net Promoter Score

ABOUT US

20⁺

Great
technology
partners

iconsulting



CLOUDERA
Connect



Google Cloud
Partner



The sparkling interaction between data intelligence and human genius

In a context in which data and algorithms are increasingly driving companies' decisions, the ability to look at the existing corporate information asset from a fresh perspective is becoming the real competitive advantage.

Advisory

Application Maintenance

Big Data Platform

Blockchain

Business Analytics

Customer Data Platform

Data Governance

Location Analytics

Machine Learning & AI

Performance Management

Agenda

First Session:

1. Data Lakehouse: Concepts & Formats
2. Delta Lake format Deep-Dive (with Practical Session)

Second Session:

1. Data Lakehouse: Vendors
2. Govern the Data Lakehouse implementation in real-world scenarios
3. Takeaways
4. Q&A Session



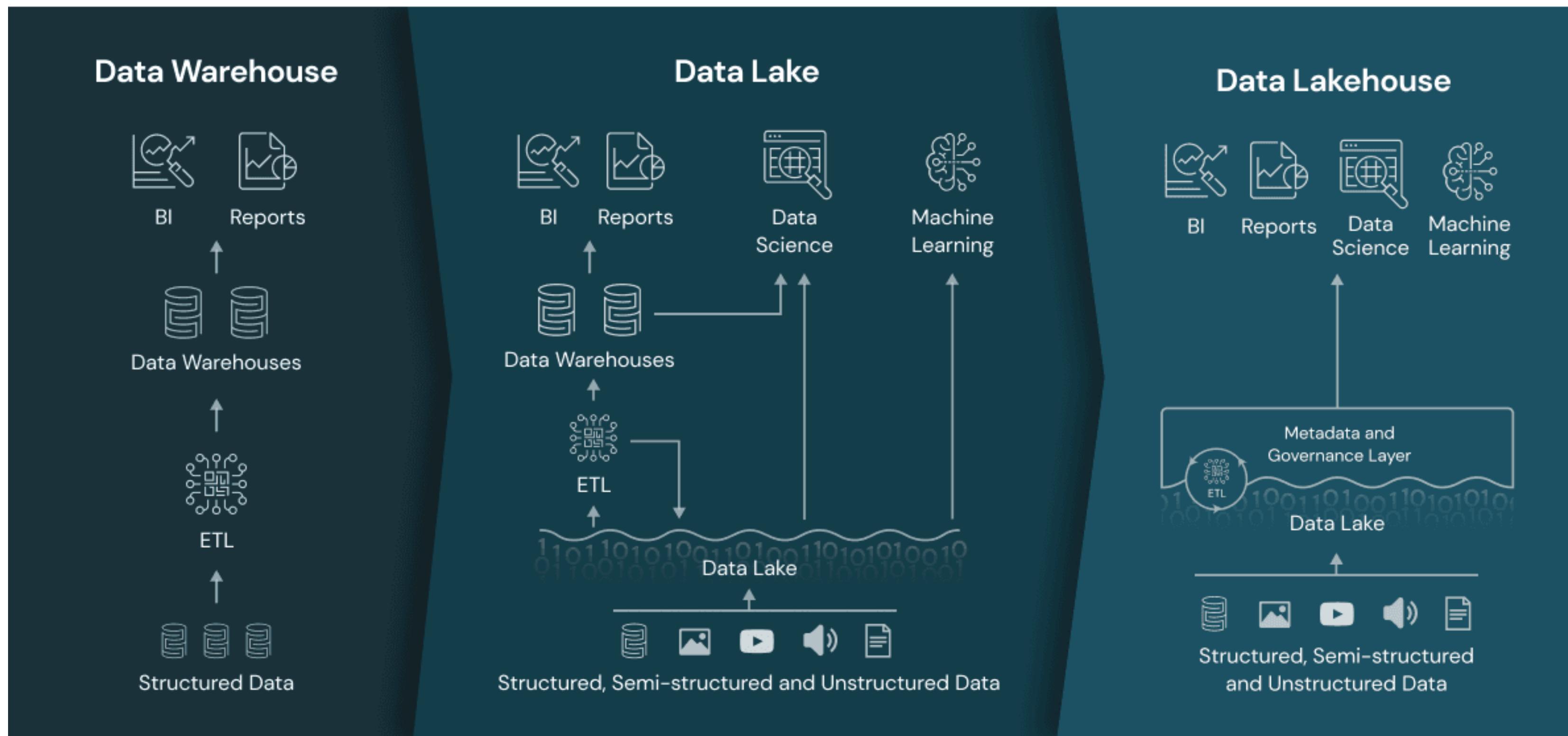
Data Lakehouse: Concepts and Formats

From Data Warehouse & Data Lake to Data Lakehouse: bring best of both to the data market

Lakehouse: A New Generation of Open Platforms that Unify Data Warehousing and Advanced Analytics

Michael Armbrust¹, Ali Ghodsi^{1,2}, Reynold Xin¹, Matei Zaharia^{1,3}

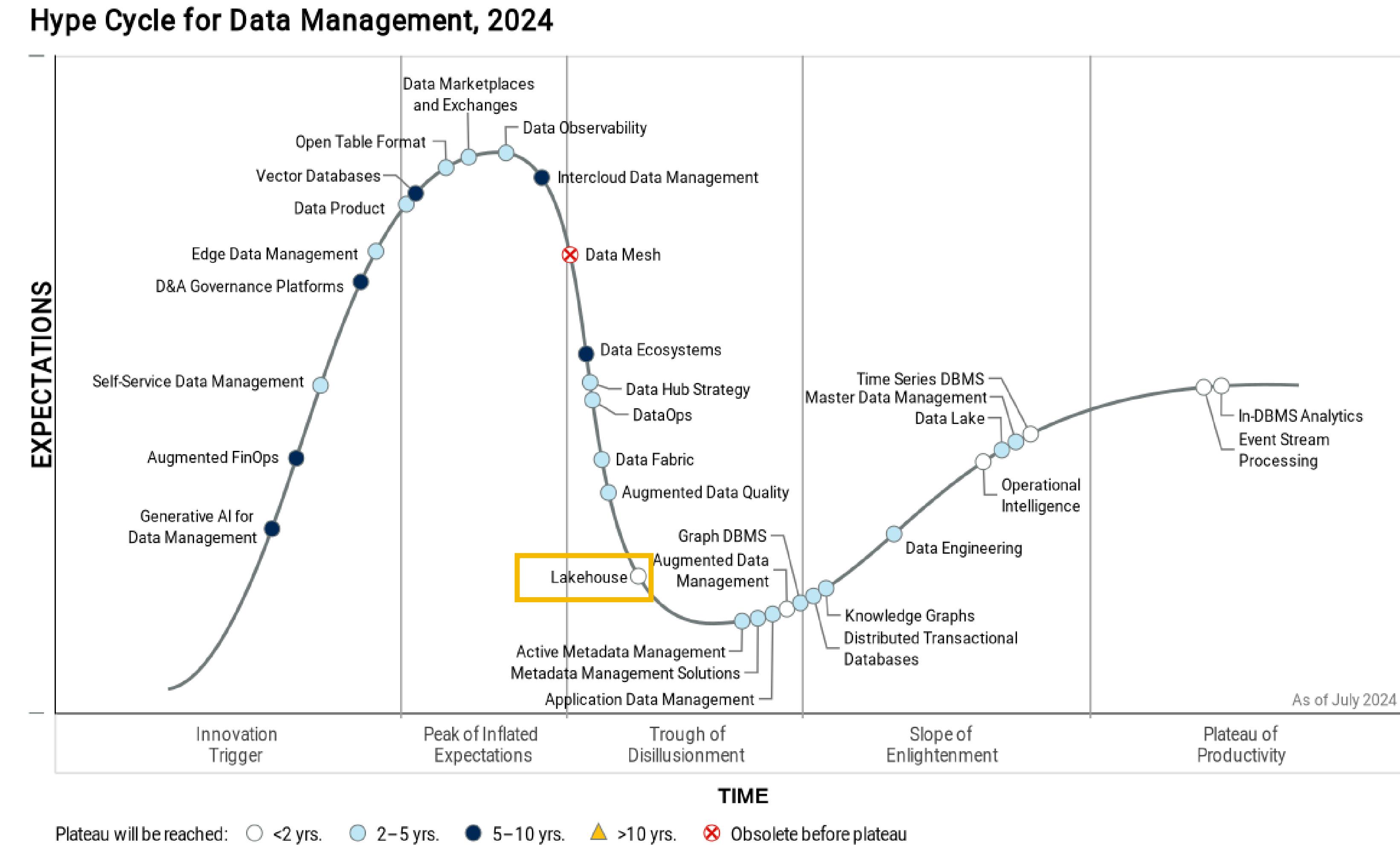
¹Databricks, ²UC Berkeley, ³Stanford University



- **Decoupling Storage and Compute**
- Store all kind of data in the same storage (Blob Storage)
- Enable ACID properties, so you can do update, delete and merge operations such as a Database
- Enable streaming of data
- Querying by using SQL, Python and Java

Is the Data Lakehouse only a Hype?

Lakehouse is in the
“Trough of Disillusionment”
phase, will reach
the **plateau of productivity in 2 years**



Gartner

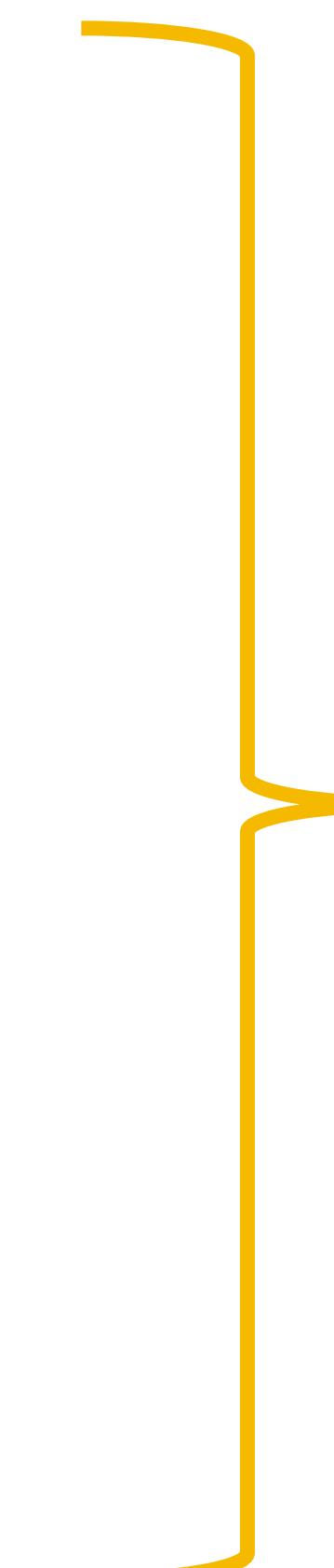
Iconsulting

But what means that storage is decoupled to Compute? Is there a real benefit?

The real benefit is that you can scale up and scale out compute resources as needed.

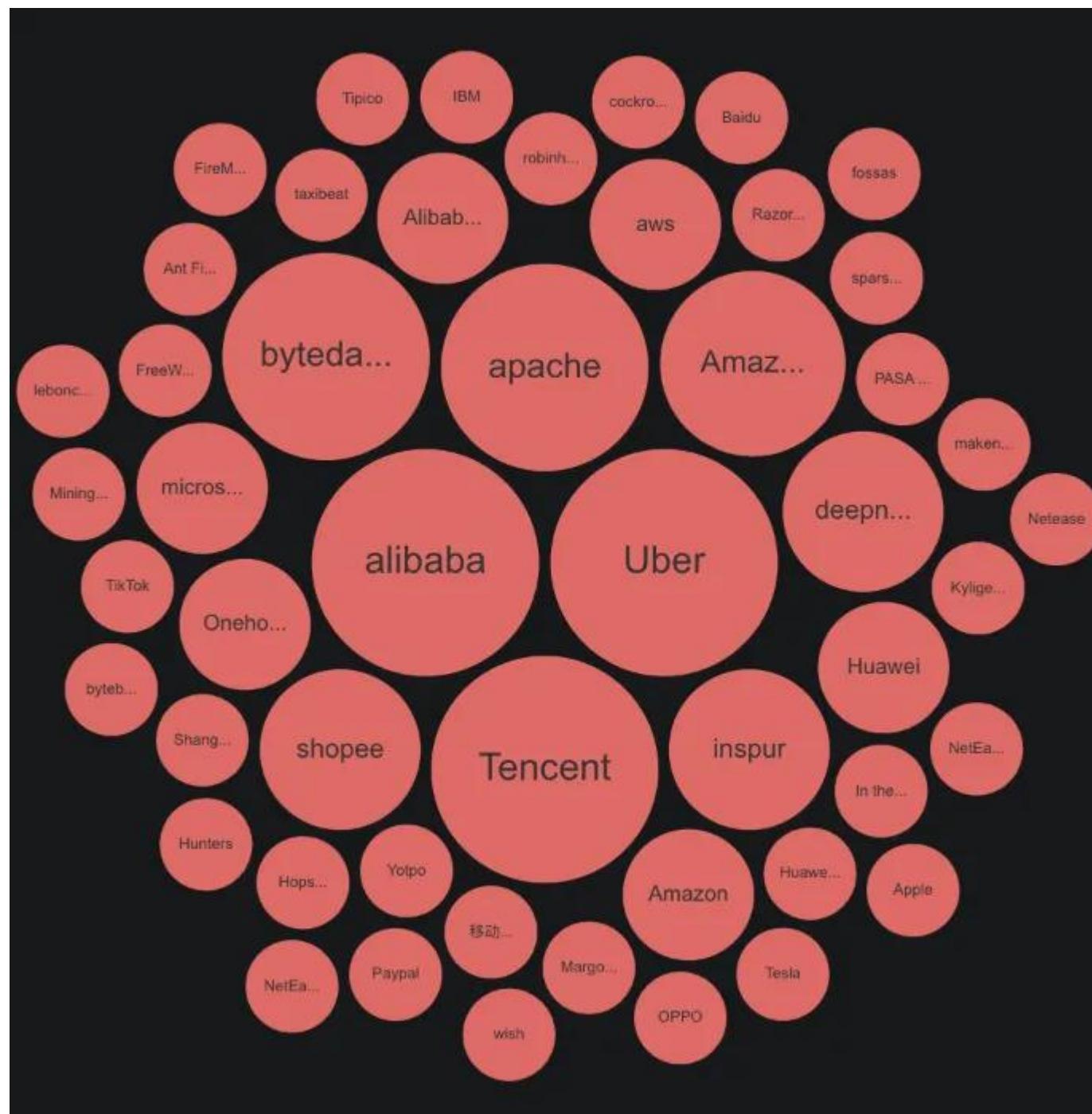
Often there are several types of workloads in a company:

- Data Engineering (ETL)
- Data Analysis (OLAP, Reporting)
- Data Streaming
- Data Science (ML & GenAI)

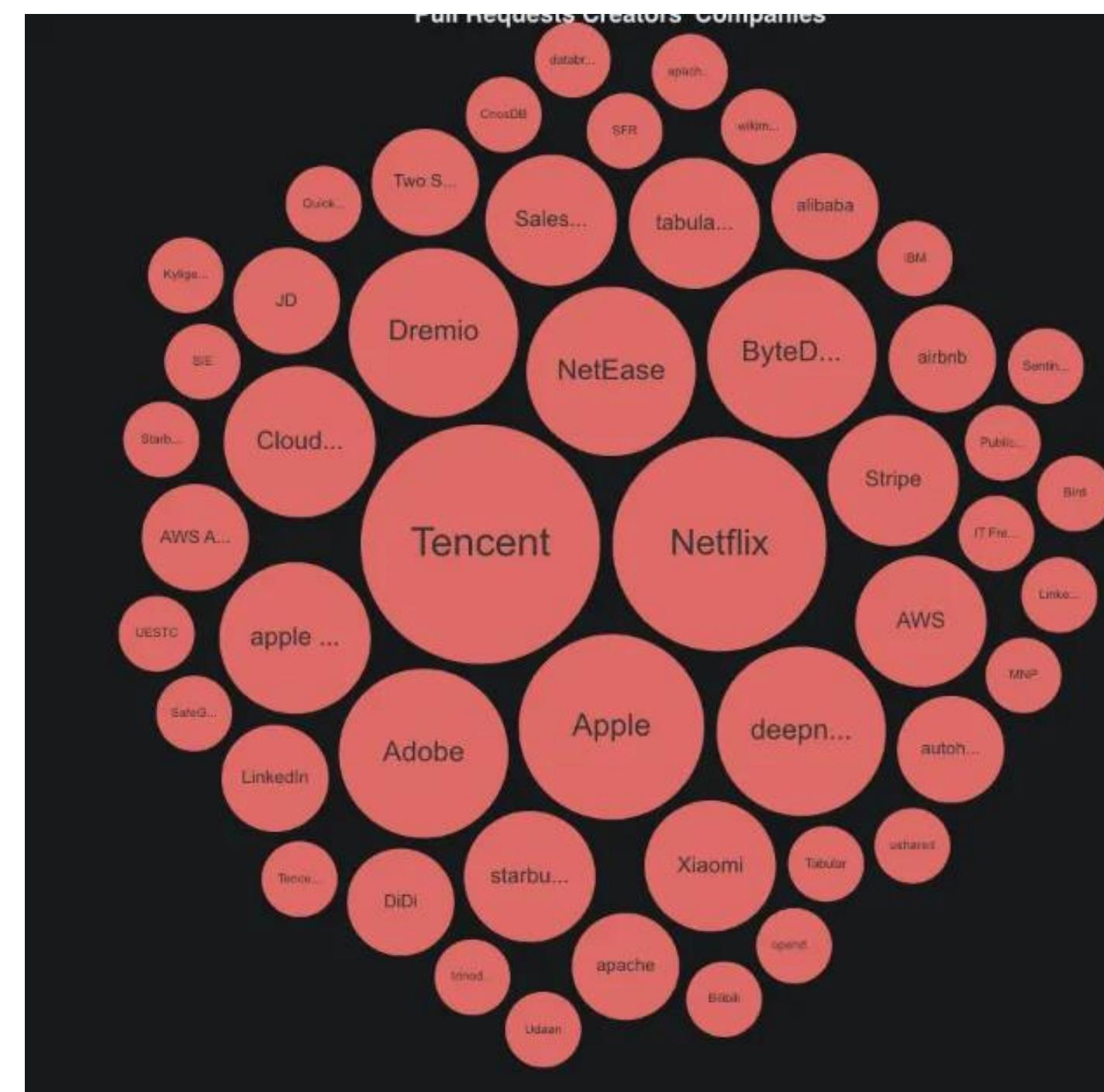


Each of them typically has dedicated compute resources (clusters) that share the same data stored in Blob Storage (e.g., Azure Data Lake Storage, AWS S3).

Three Lakehouse storage framework contended the market,



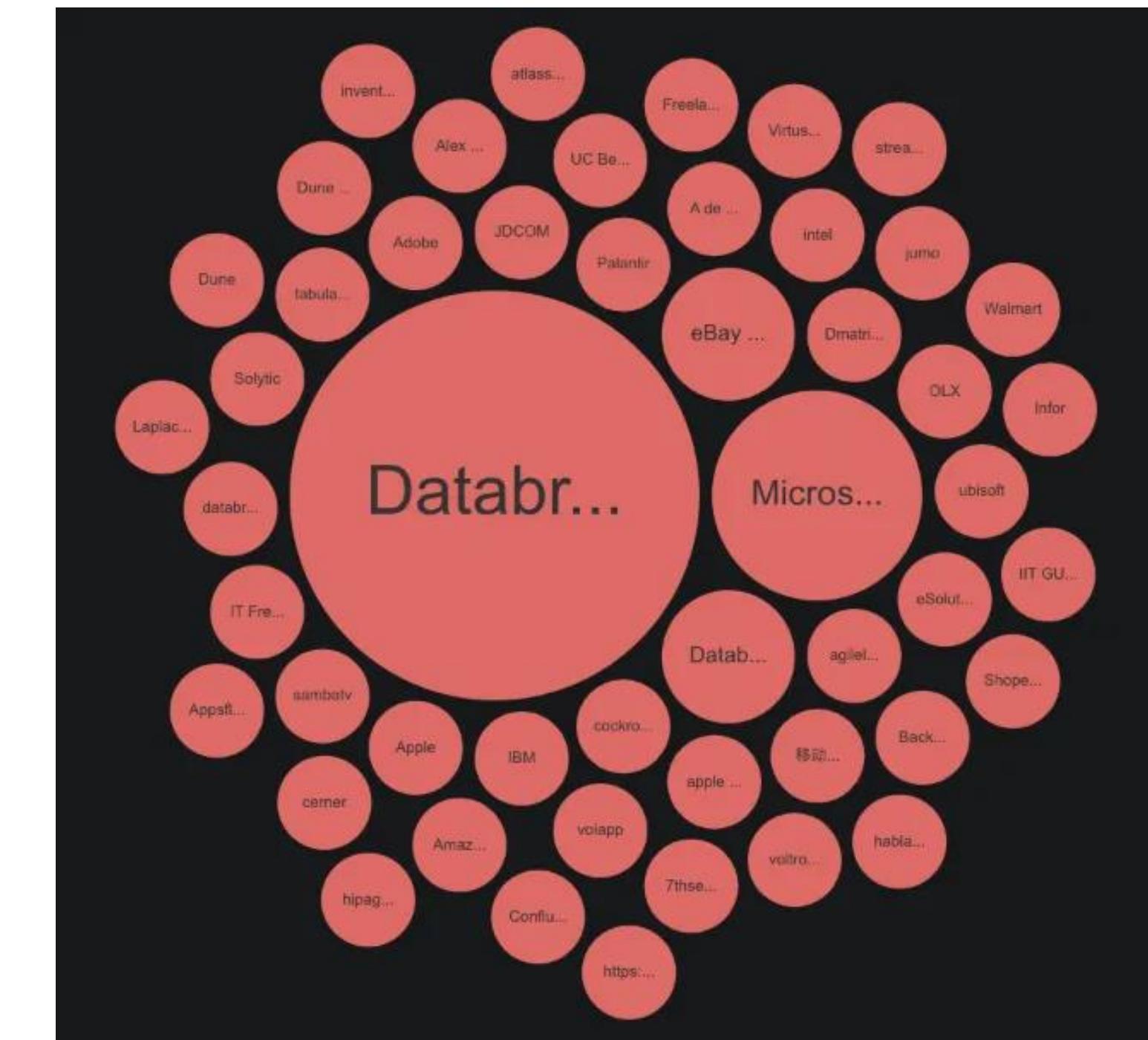
Contributors to Apache Hudi



Contributors to **Iceber**



DELTA LAKE



Contributors to Delta Lake

...But what a Data Lakehouse format really is?

A **Data Lakehouse** format is a metadata structure layered on top of **Parquet** files, enabling **ACID** properties.

Different formats have distinct metadata configurations.



```
s3_bucket/my_table/  
|- .hoodie/  
|  |- hoodie.properties  
|  |- metadata/  
|- file_1.parquet  
|- file_2.parquet  
|- file_N.parquet
```



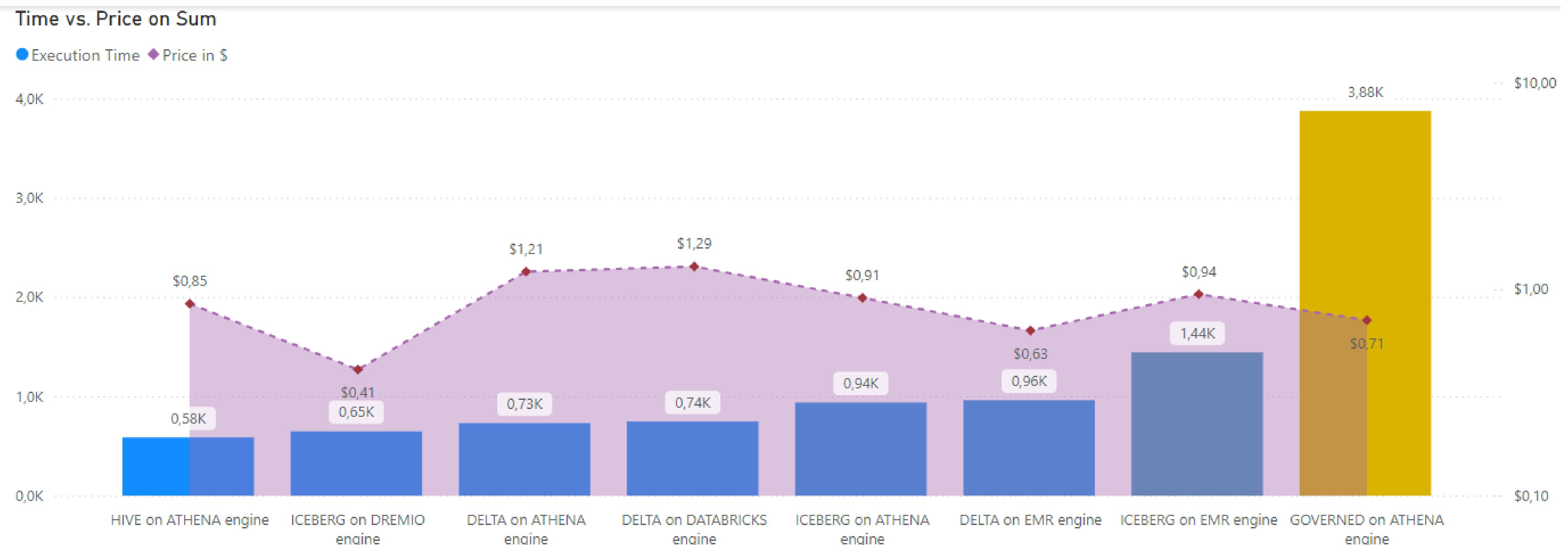
```
s3_bucket/my_table/  
|- _delta_log/  
|  |- 000000.json  
|- file_1.parquet  
|- file_2.parquet  
|- file_N.parquet
```



```
s3_bucket/my_table/  
|- metadata/  
|  |- v1.metadata.json  
|  |- snap-9fa1-2-16c3.avro  
|  |- 0d9a-98fa-77.avro  
|- file_1.parquet  
|- file_2.parquet  
|- file_N.parquet
```

...But which one is the best performer? Is there any benchmark?

Yes! There are several benchmarks, but performance vary depending on the compute engine that consume the format (e.g. Apache Spark) so it's not easy to compare performances.





Delta Lake Deep-Dive (with Practical Session)

So what is Delta Lake?

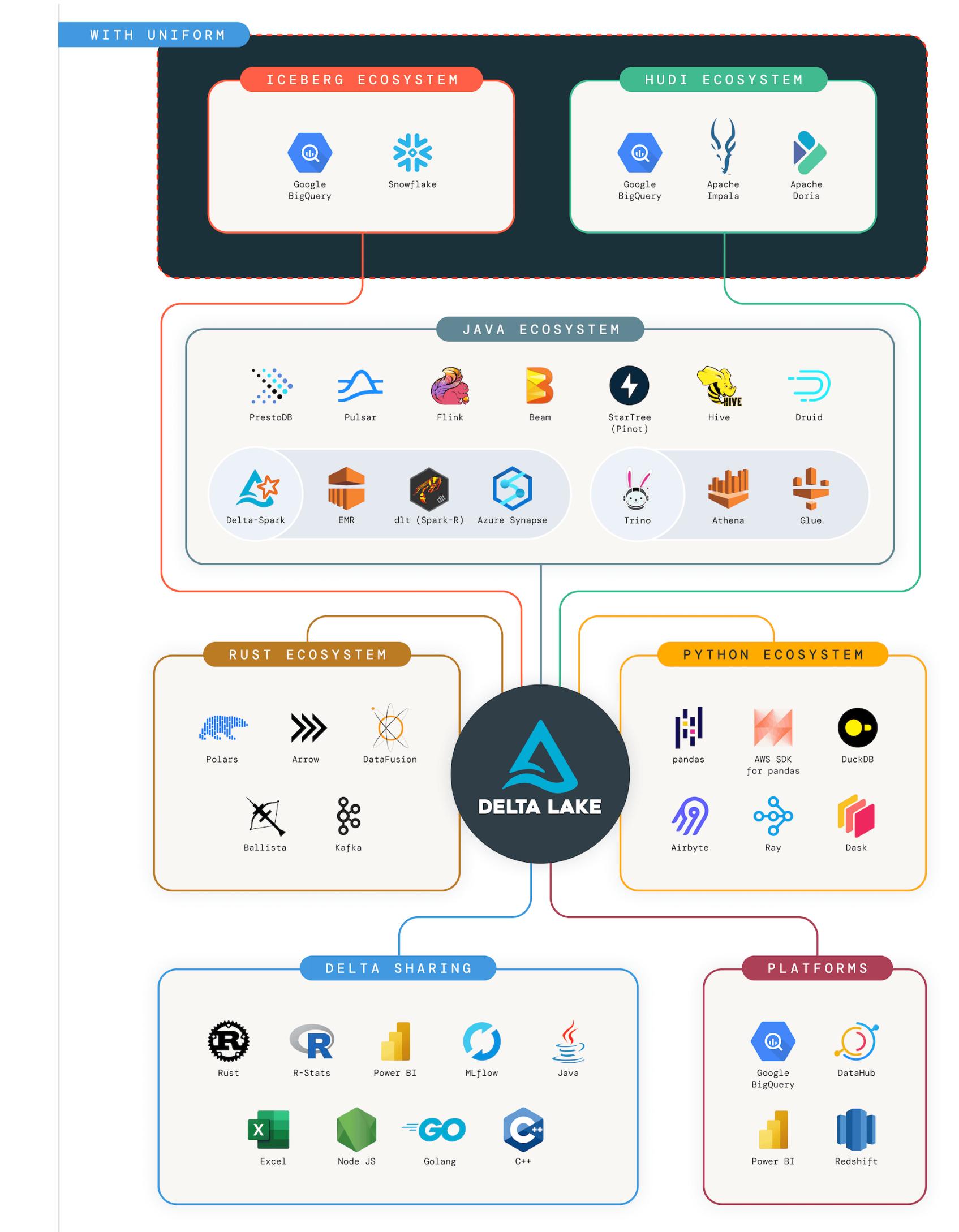
Delta Lake is an open-source **storage framework** that enables building a format agnostic **Lakehouse architecture** with compute engines including:

Spark, PrestoDB, Flink, Trino, Hive, Snowflake, Google BigQuery, Athena, Redshift, Databricks, Azure Fabric

With **Delta Universal Format** a.k.a. UniForm, you can read now Delta tables with Iceberg and Hudi clients.

Databricks Agrees to Acquire Tabular, the Company Founded by the Original Creators of Apache Iceberg

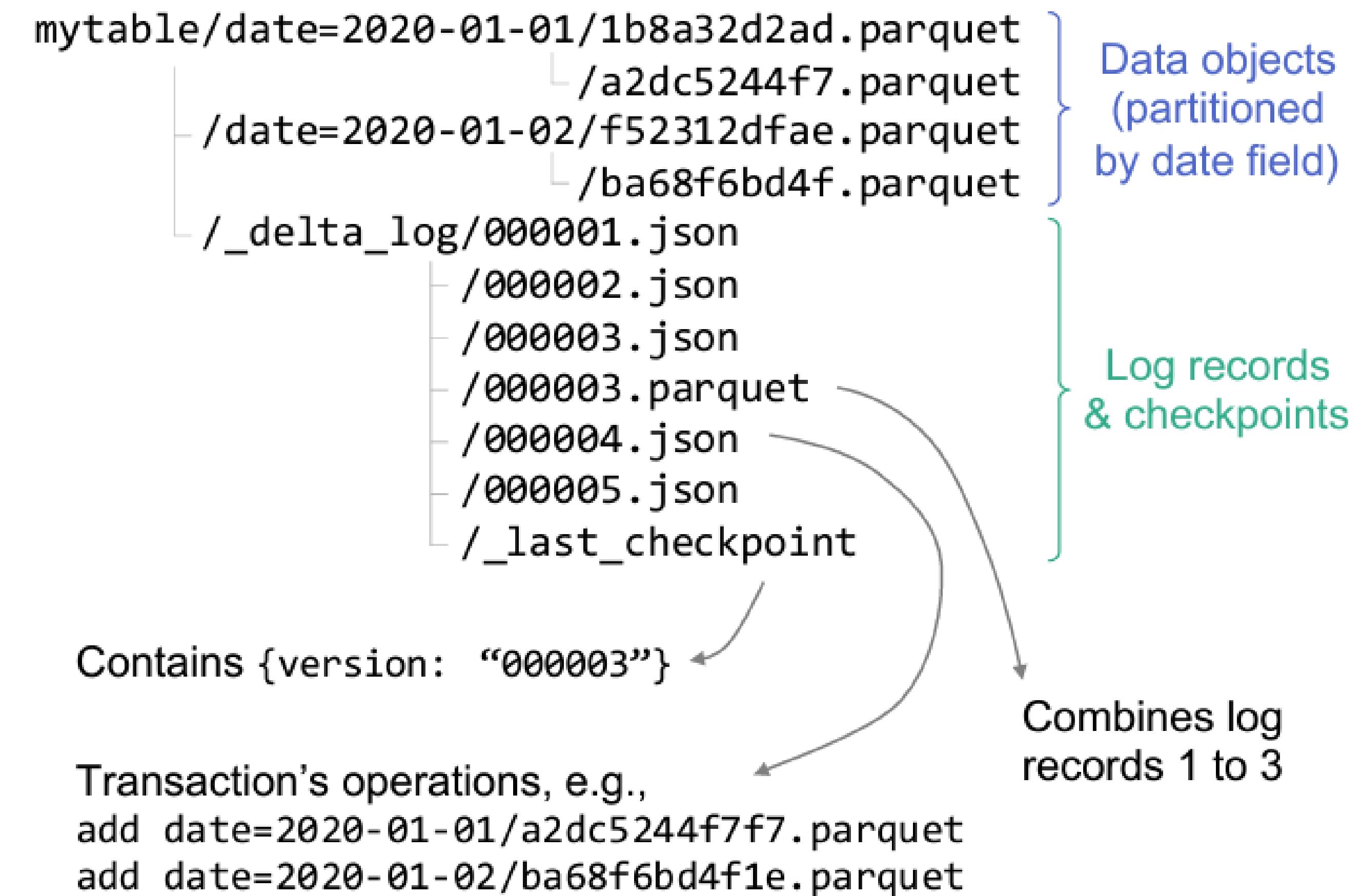
June 4, 2024



But how data are stored in Delta Lake Tables? And how it's possible to preserve ACID properties?

Delta Lake's approach is to store a transaction log and metadata directly within the **cloud object store**, using a set of protocols over object store operations to achieve serializability and preserve ACID properties.

The data within a table is then stored in Parquet format.



Ok, understood how data are stored, which feature it enables?

Delta Lake's transactional design enables a wide range of data management features similar to those in RDBMS, overcoming some limitations of Data Lakes:

- Efficient Update, Delete, and Merge operations
- Time Travel and Rollbacks

Delta Lake also supports **Streaming Ingestion and Consumption**, as well as **Schema Evolution and Enforcement**, but we won't focus on these.

Thanks for the explanation, but it seems to be hard to understand and to focus on...could we try it?

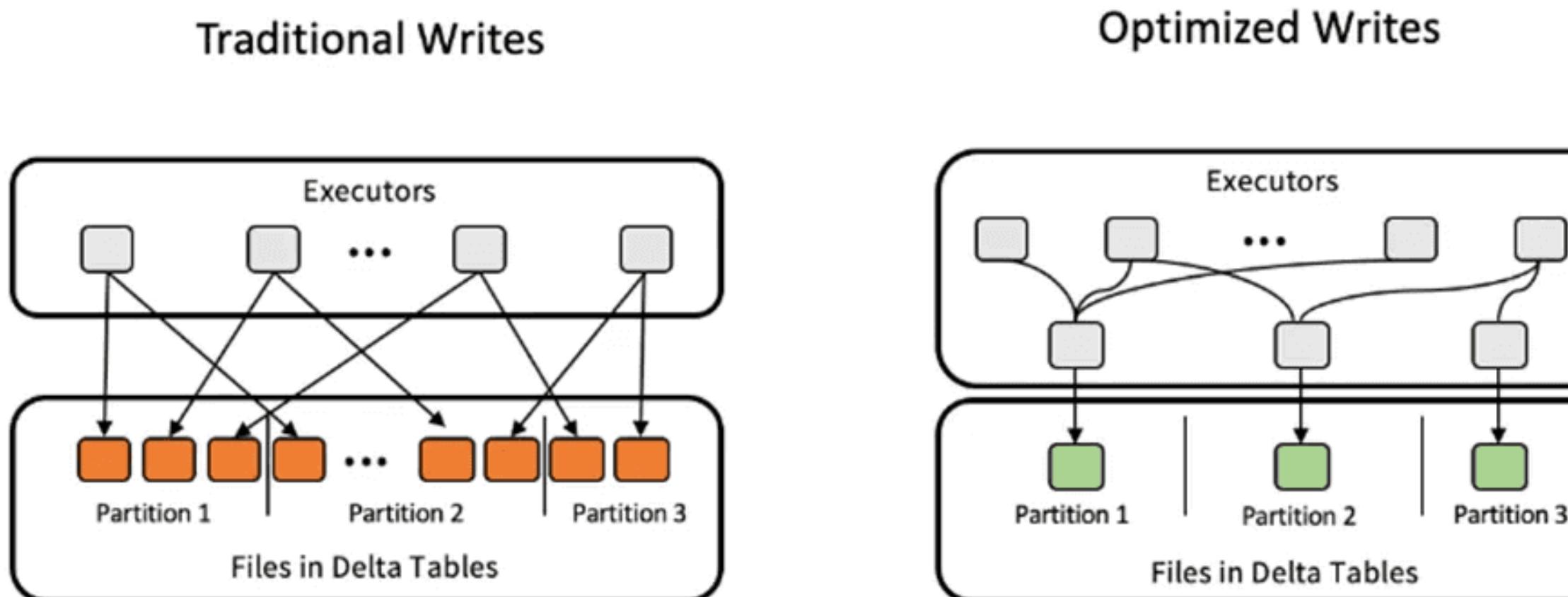
Well, let's test it in a real environment

[Test Environment](#)

...But what happens to the transaction log when several DML operations are done? Does it scale?

Frequent DML operations can fragment both the log and data, reducing efficiency. Small files can also be problematic, as they slow down query reads due to the overhead of listing, opening, and closing many files.

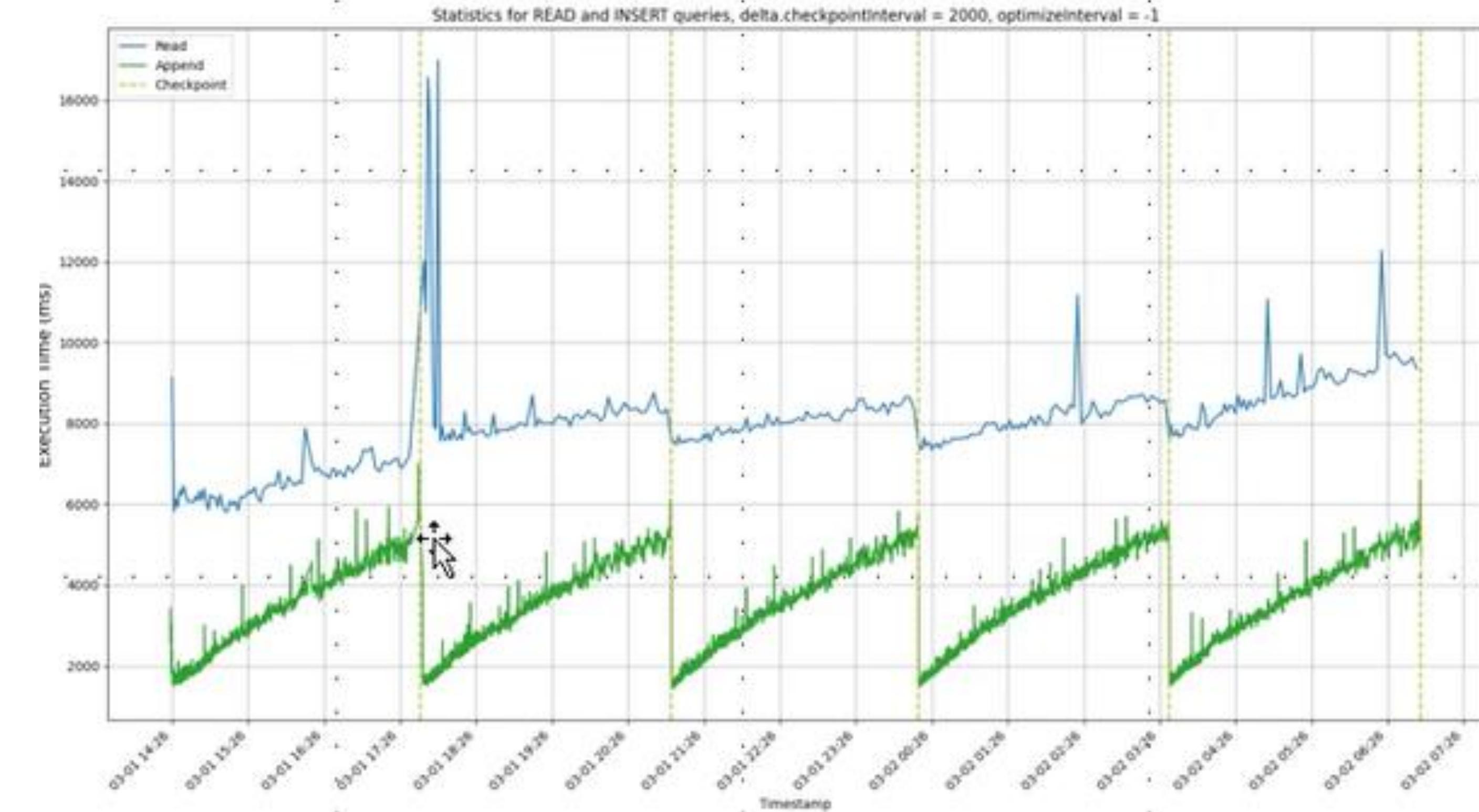
To maintain data read and write performance, **OPTIMIZE** operations are necessary to compact both the transaction log and data.



...ICONULTING and DISI worked together to understand better about Delta Lake performances

OPTIMIZE operation strategy is important to preserve performance over time.

Some tools, such as **Databricks**, performs automatic optimizations to preserve performance both read and write.





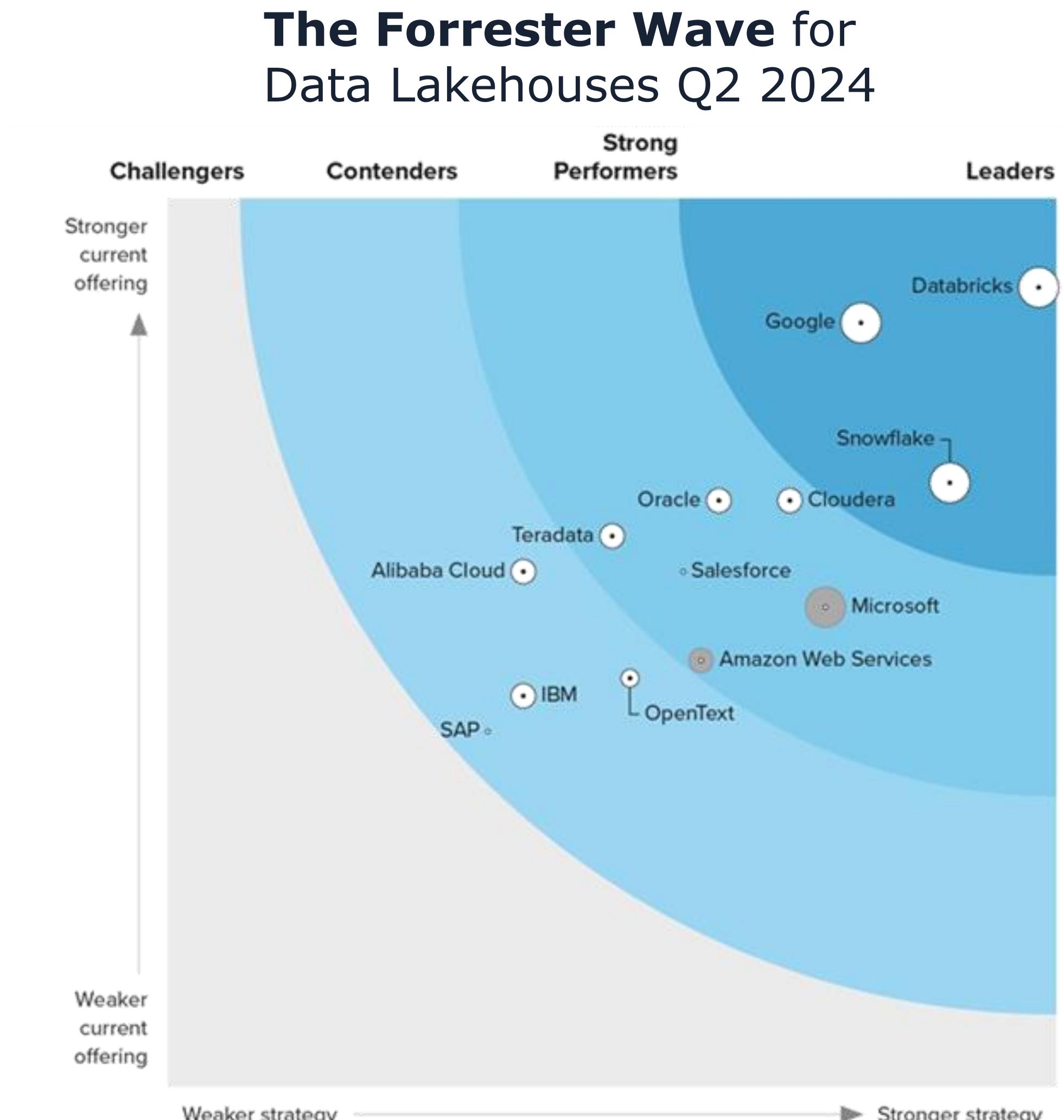
Data Lakehouse: Vendors

Which are the main player here? Who could win the competition?...

More **vendors** have focused on the new **Data Lakehouse paradigm**.

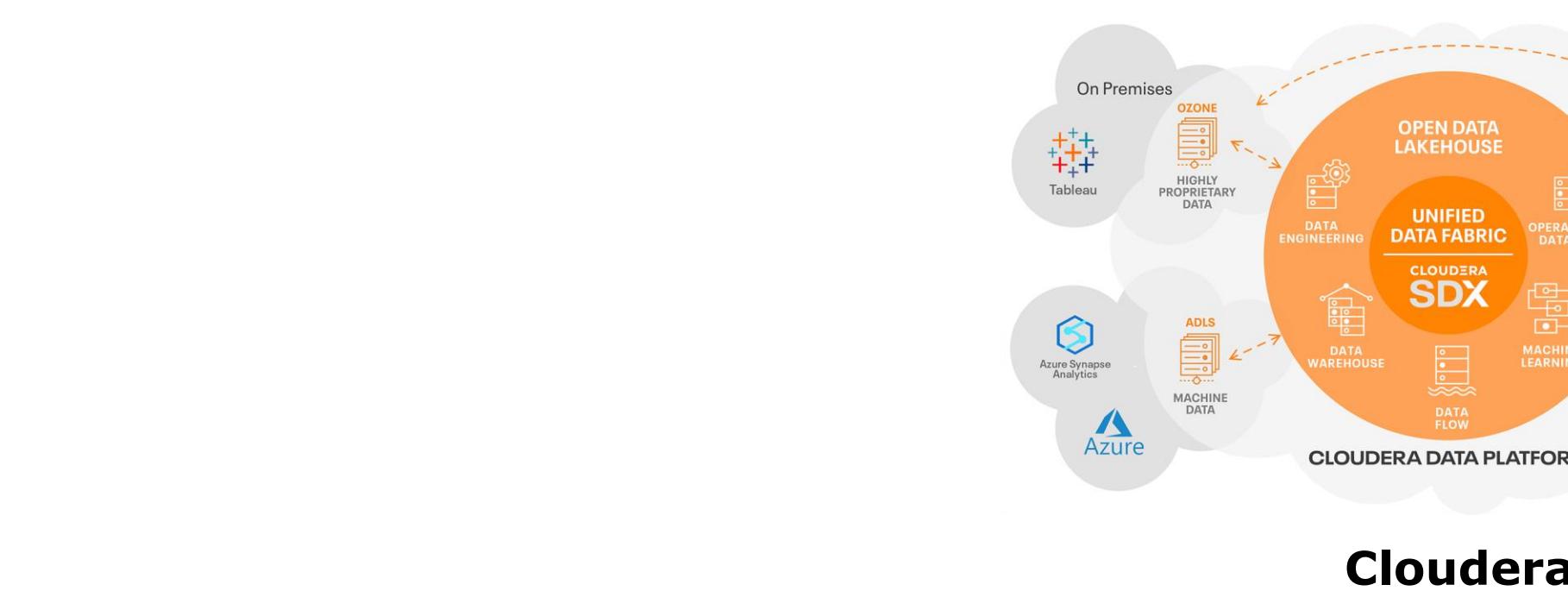
Currently the leaders in the market are **Databricks**, **Google** and **Snowflake**.

But also **AWS** and **Microsoft** are strong performers

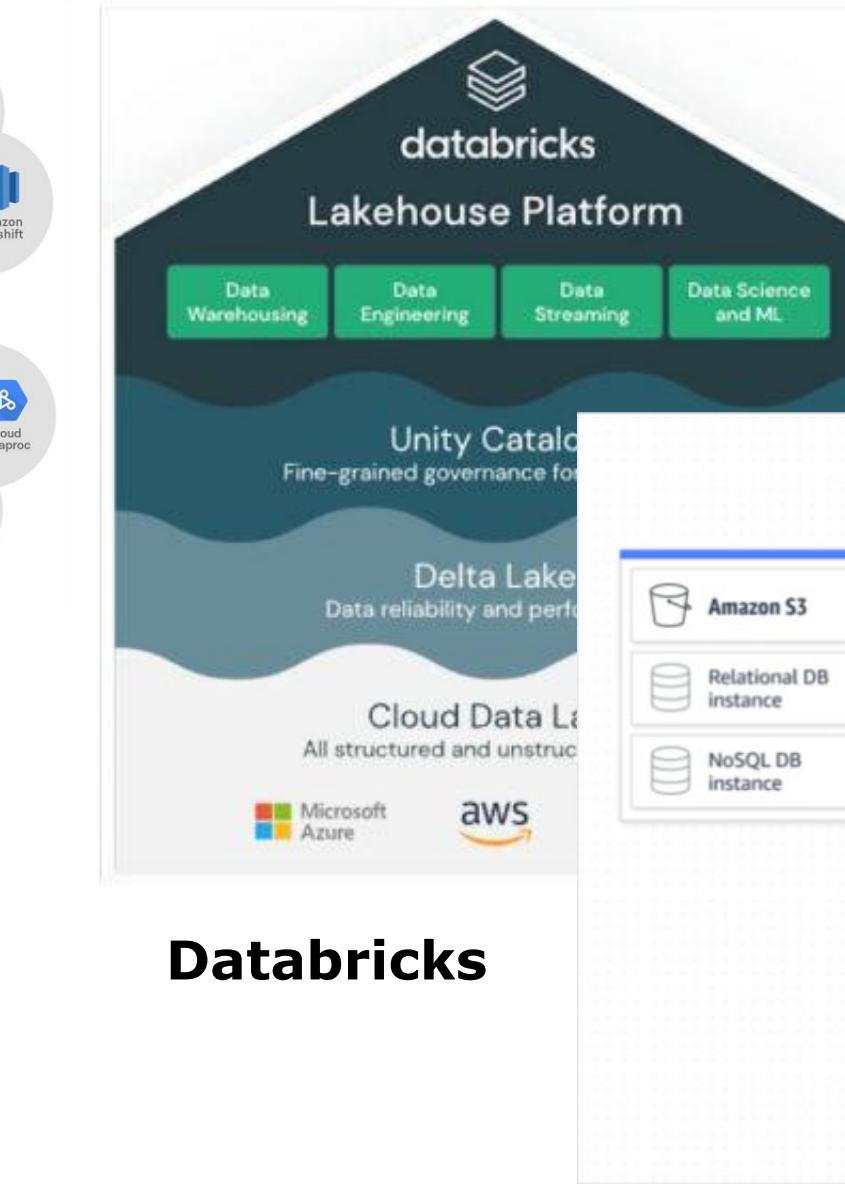


*A gray bubble or open dot indicates a nonparticipating vendor.

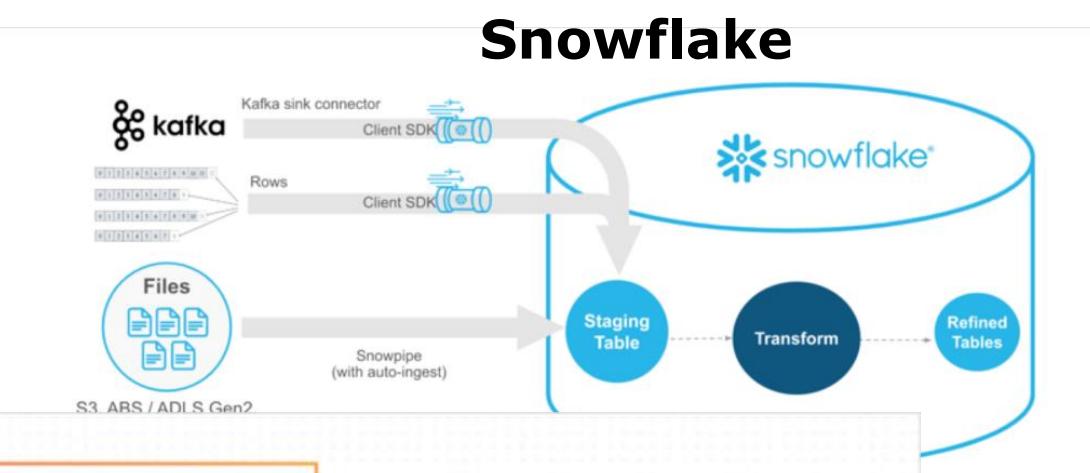
What's differs a vendor to another?



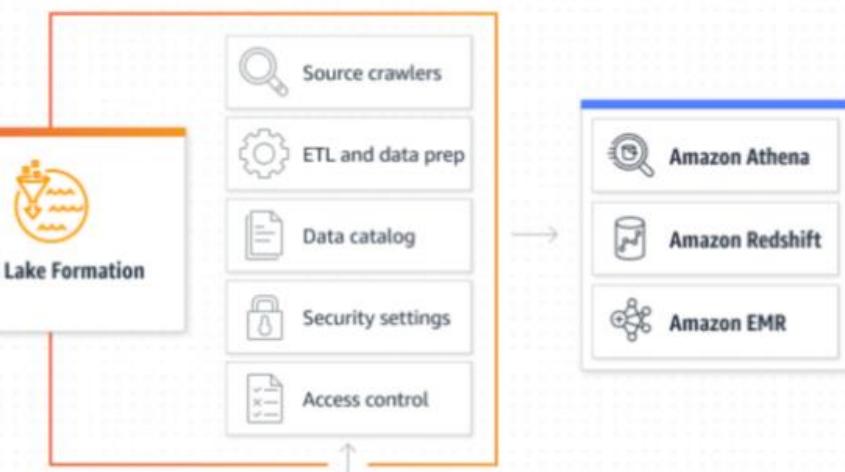
Cloudera



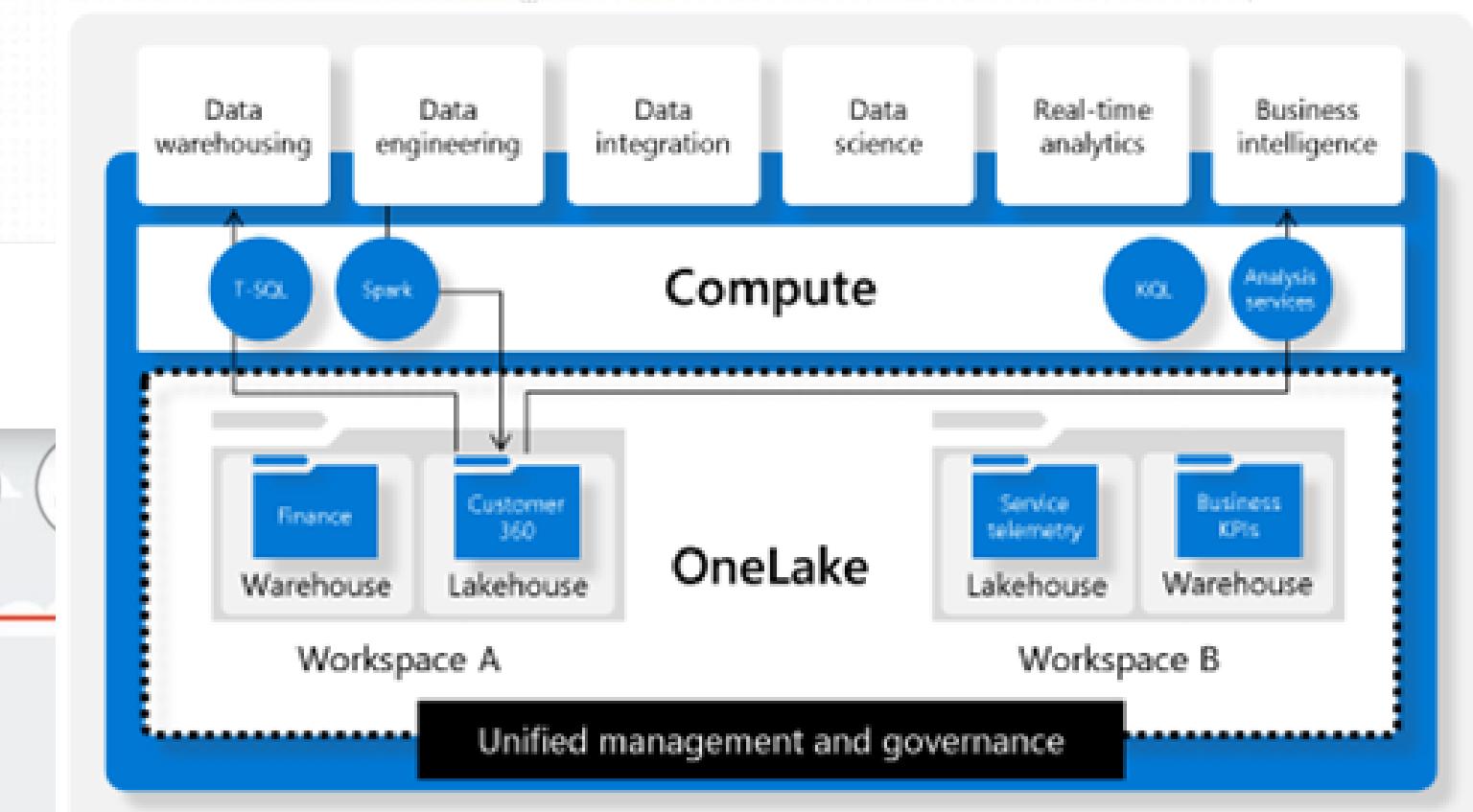
Databricks



AWS



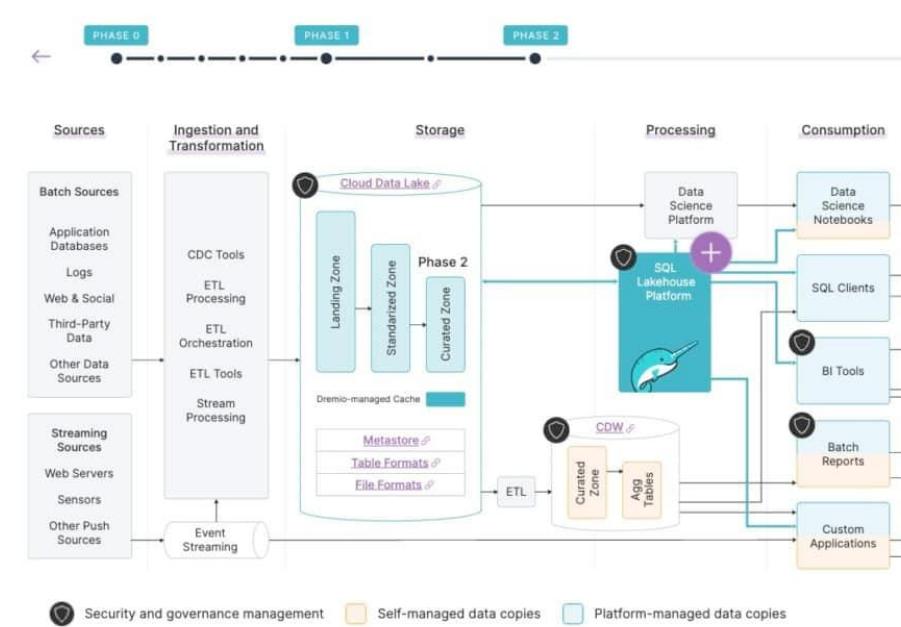
MS Fabric



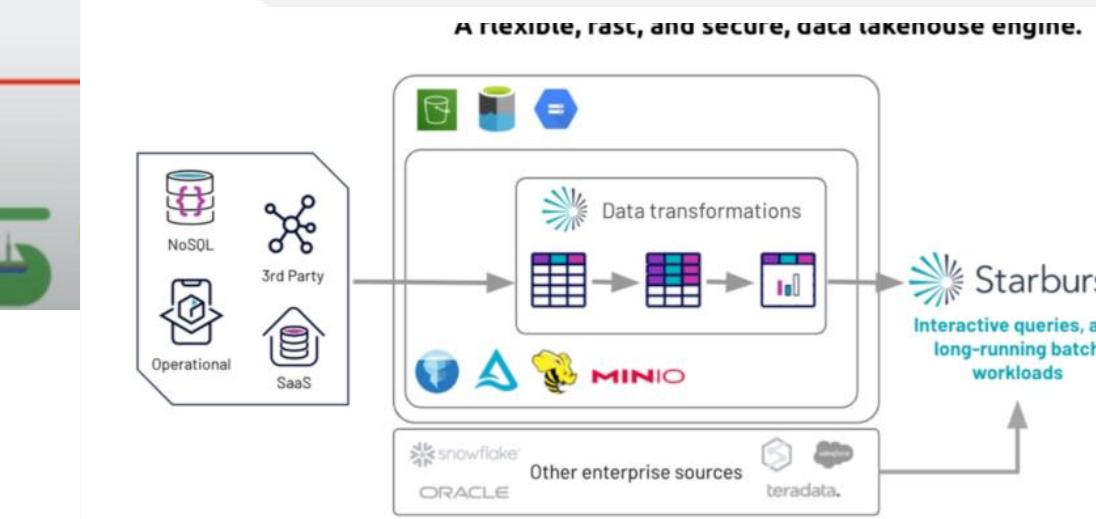
Google Big Query



Dremio



Security and governance management Self-managed data copies Platform-managed data copies



Starbursts

Consulting

Each vendor has its own **Architectural Landscape** for Data Lakehouse with their own building blocks.

Building blocks are similar but each has a different configuration.

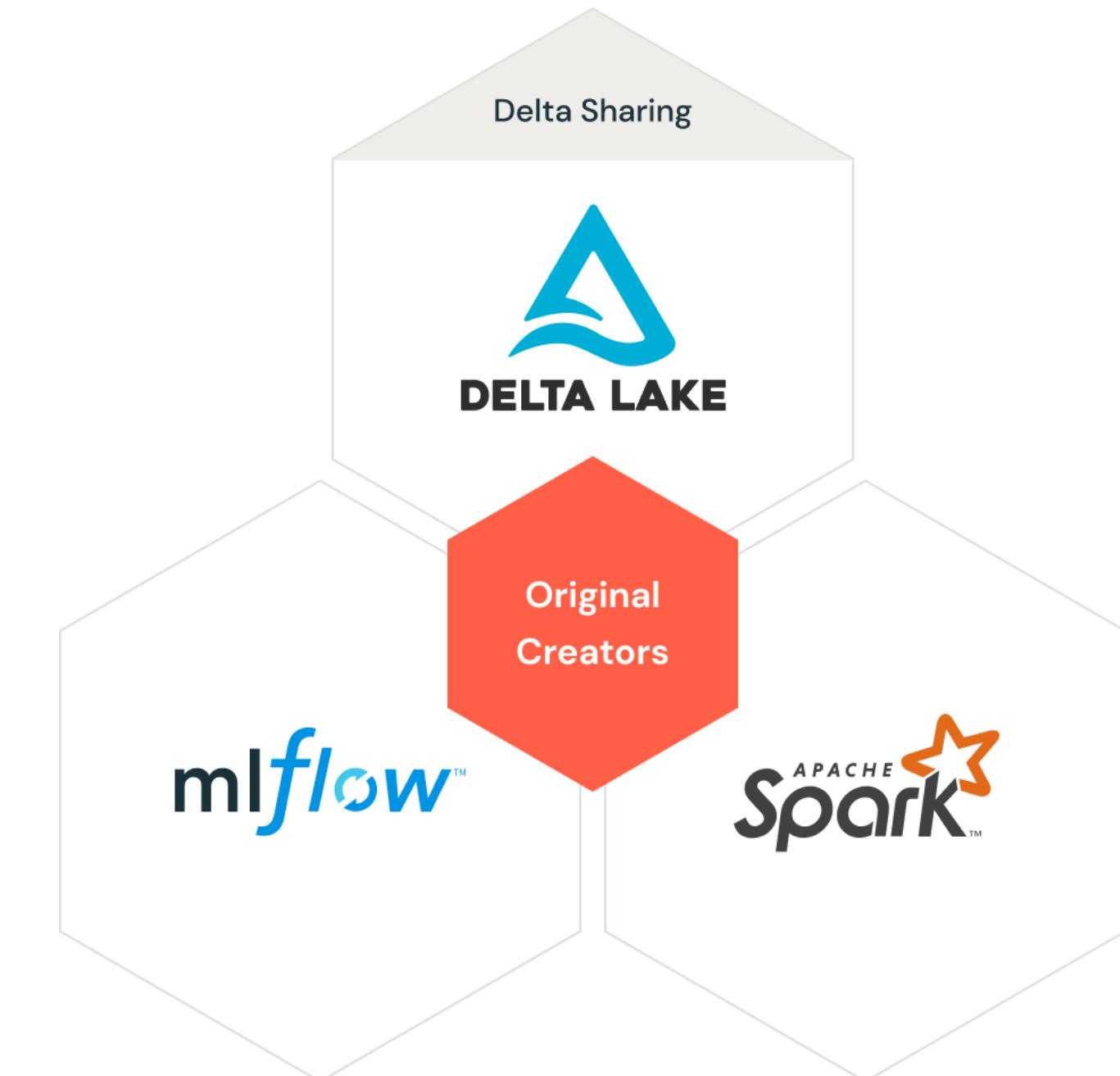
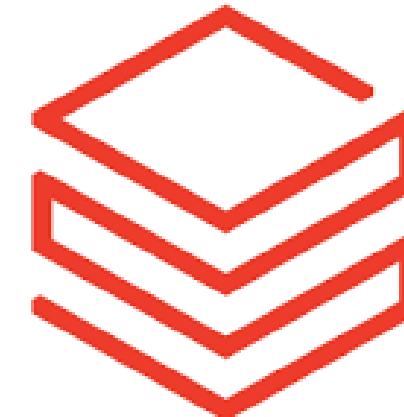
...So which one is the best? And which one do you prefer?

As of today, in my opinion, the best Platform to implement a Data Lakehouse is **Databricks**

Databricks founders also found **Delta Lake** format as well as **Spark** compute engine.

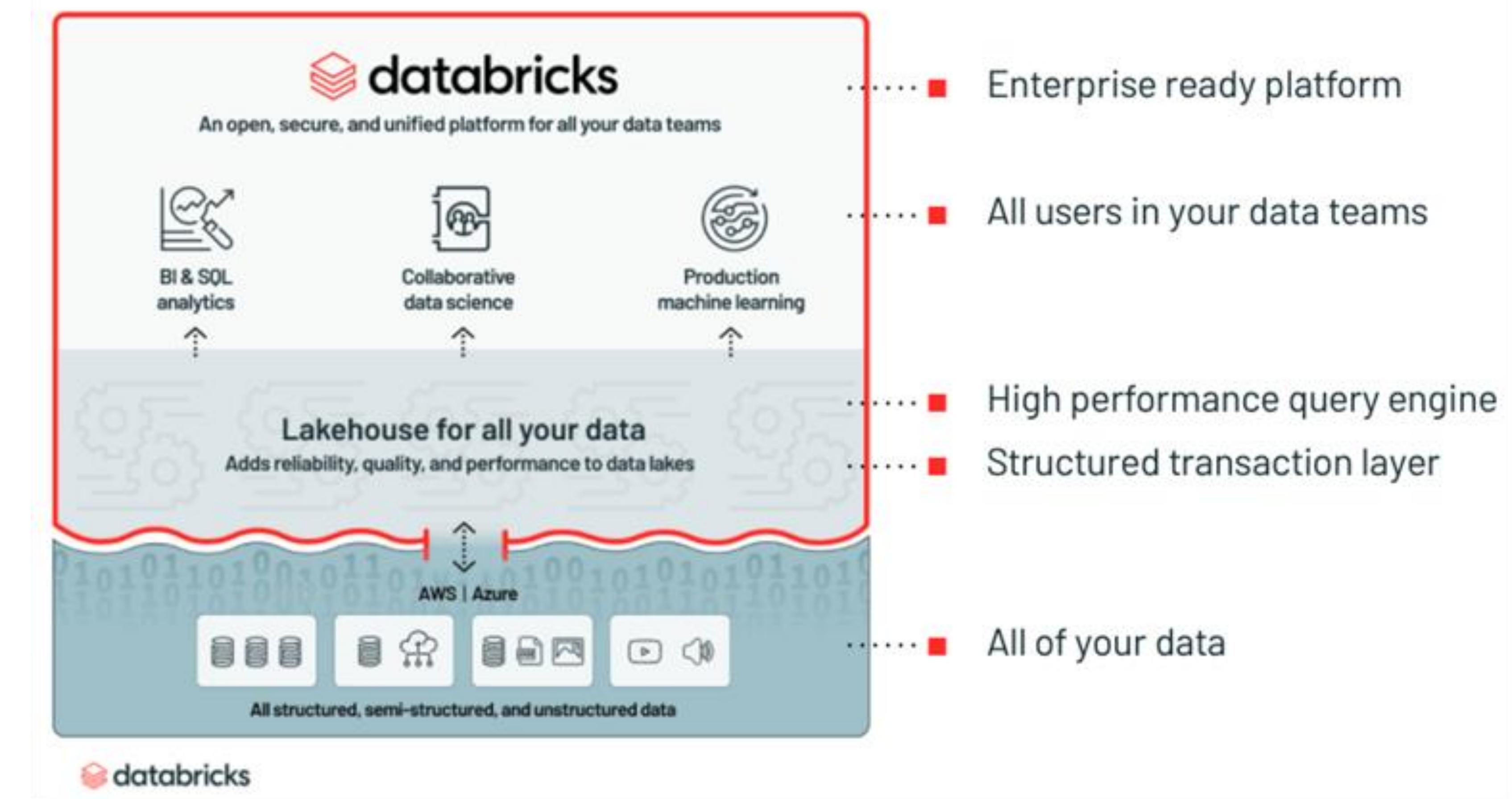
Understanding and deepen **Databricks** is a good step toward understanding how a Data Lakehouse works.

Databricks is best-in-class platform also for **Data Governance, Machine Learning & GenAI**



...Ok, I got it. But give me more about Databricks!

Databricks is a modern Data Platform that leverages the **Lakehouse** to enable BI, Data Science, Machine Learning and AI



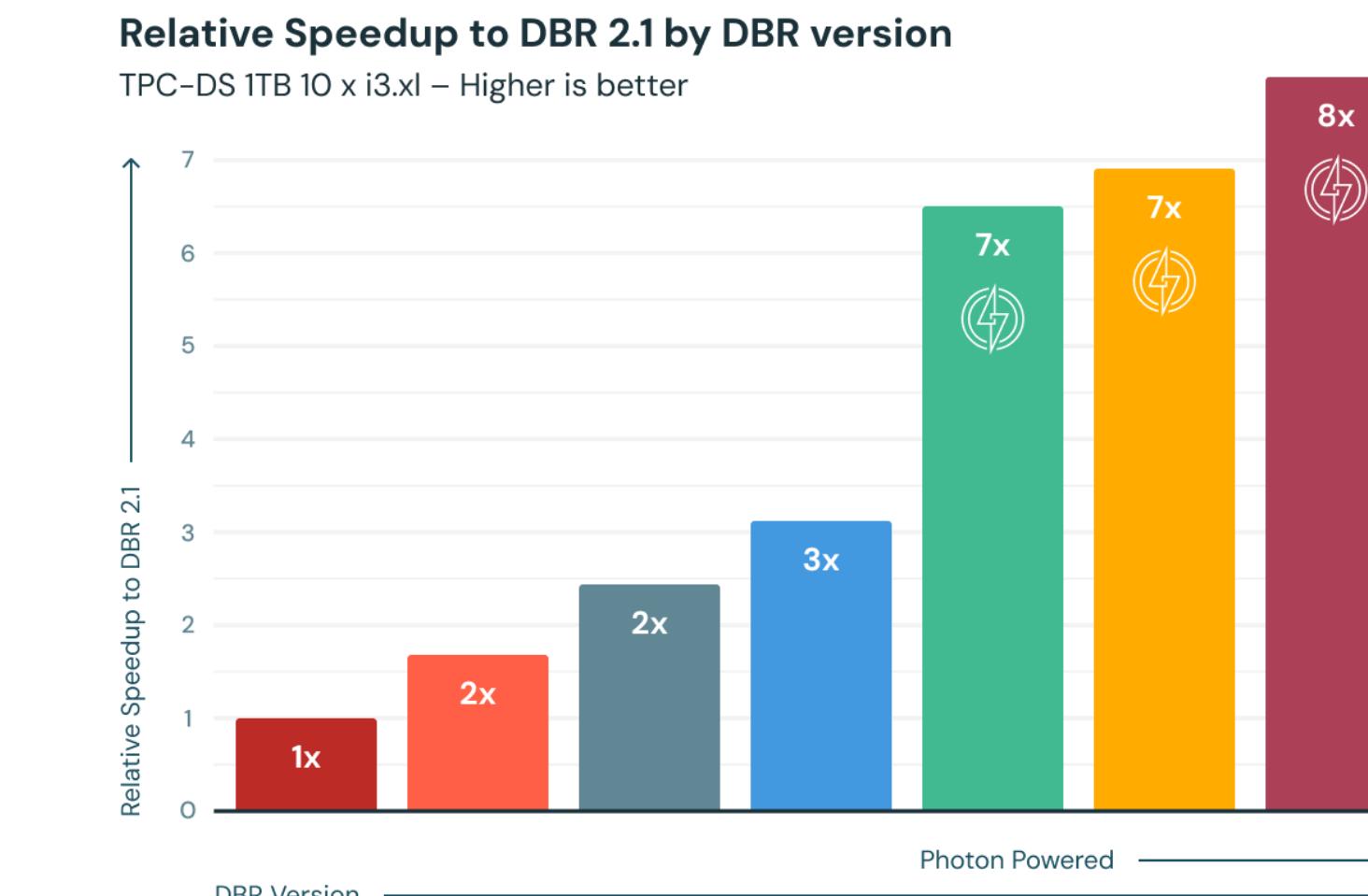
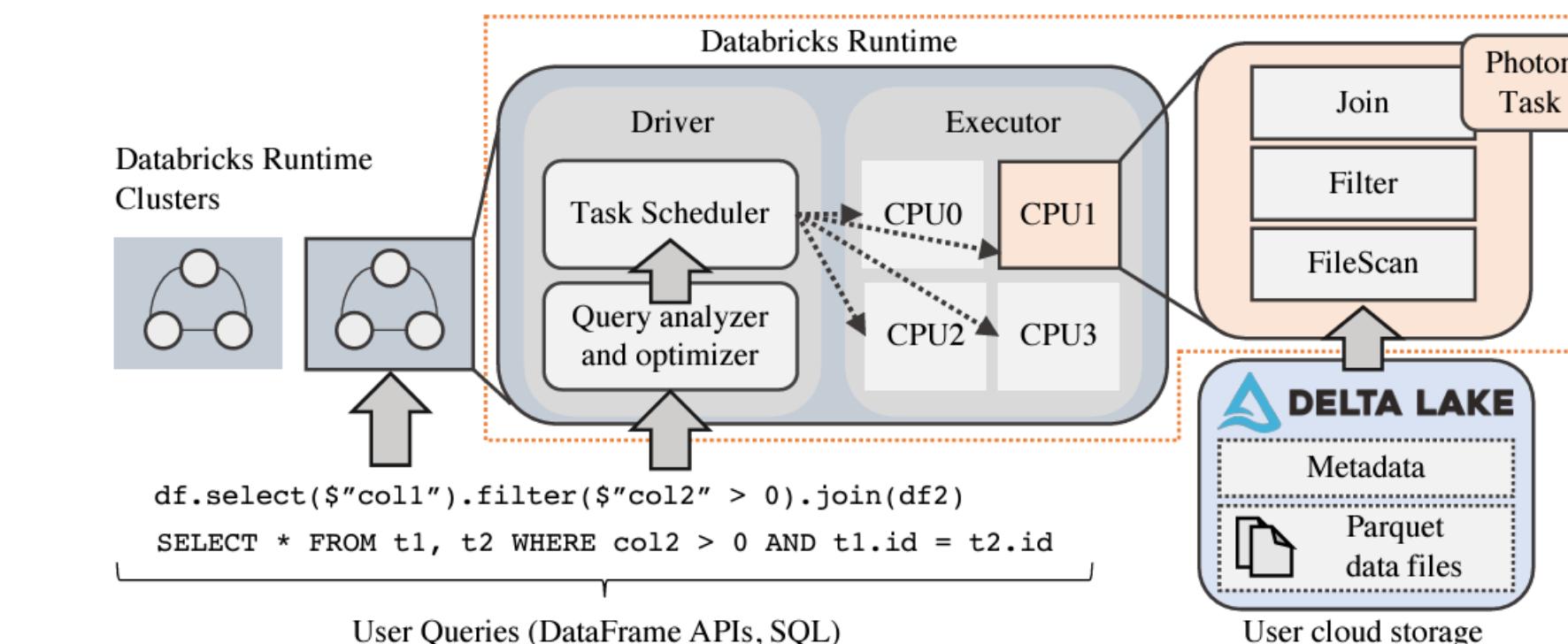
...But what about Lakehouse performances for the BI Workload? Is it better than classic RDBMS?

Building a structured data warehouse on top of an unstructured data lake **presents new challenges for SQL query execution engines**.

Significant efforts have been made to enhance performance, and it can now be compared to that of in-memory relational DBMSs, thanks to **new optimizations** like Databricks' **Photon engine**.

Photon: A Fast Query Engine for Lakehouse Systems

Alexander Behm, Shoumik Palkar, Utkarsh Agarwal, Timothy Armstrong, David Cashman, Ankur Dave, Todd Greenstein, Shant Hovsepian, Ryan Johnson, Arvind Sai Krishnan, Paul Leventis, Ala Luszczak, Prashanth Menon, Mostafa Mokhtar, Gene Pang, Sameer Paranjpye, Greg Rahn, Bart Samwel, Tom van Bussel, Herman van Hovell, Maryann Xue, Reynold Xin, Matei Zaharia
photon-paper-authors@databricks.com
Databricks Inc.





**Govern the Data Lakehouse implementation
in real-world scenarios**

Before starting to implement a Data Lakehouse, we took the time to define a Data Strategy and a Tech Blueprint...

Several questions arise before starting the implementation that need accurate answers:

What?

Identify **use cases** and **data subjects** from business users' analytical needs.

When?

Define an **implementation roadmap** for use cases and data subjects.

Who?

Setup an **interdisciplinary team** with professionals and business owners.

Setup a **Control Tower** to control the initiative.

How?

Define strong standards and adopt **DevOps**, **DataOps** and **ModelOps** principles.

Setup a **development methodology** (e.g. Agile).

How much?

Simulate a **5 Year TCO** to evaluate the Data Lakehouse Initiative in terms of **CAPEX** and **OPEX**.

Monitor ongoing **cloud consumption costs**.

Why?

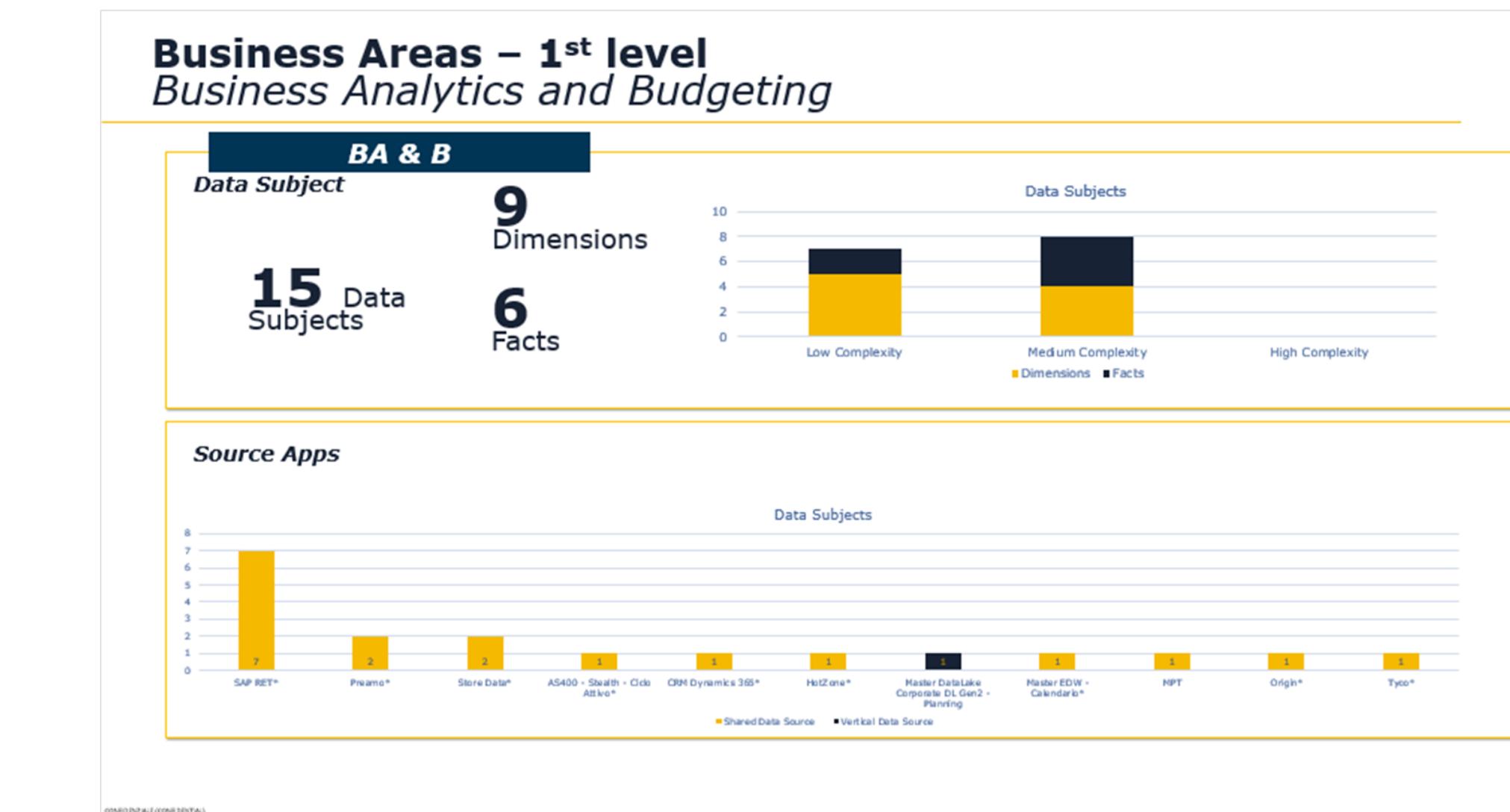
Define **clear objectives** of the initiatives in terms of benefits both for IT and Business Stakeholders.

Where?

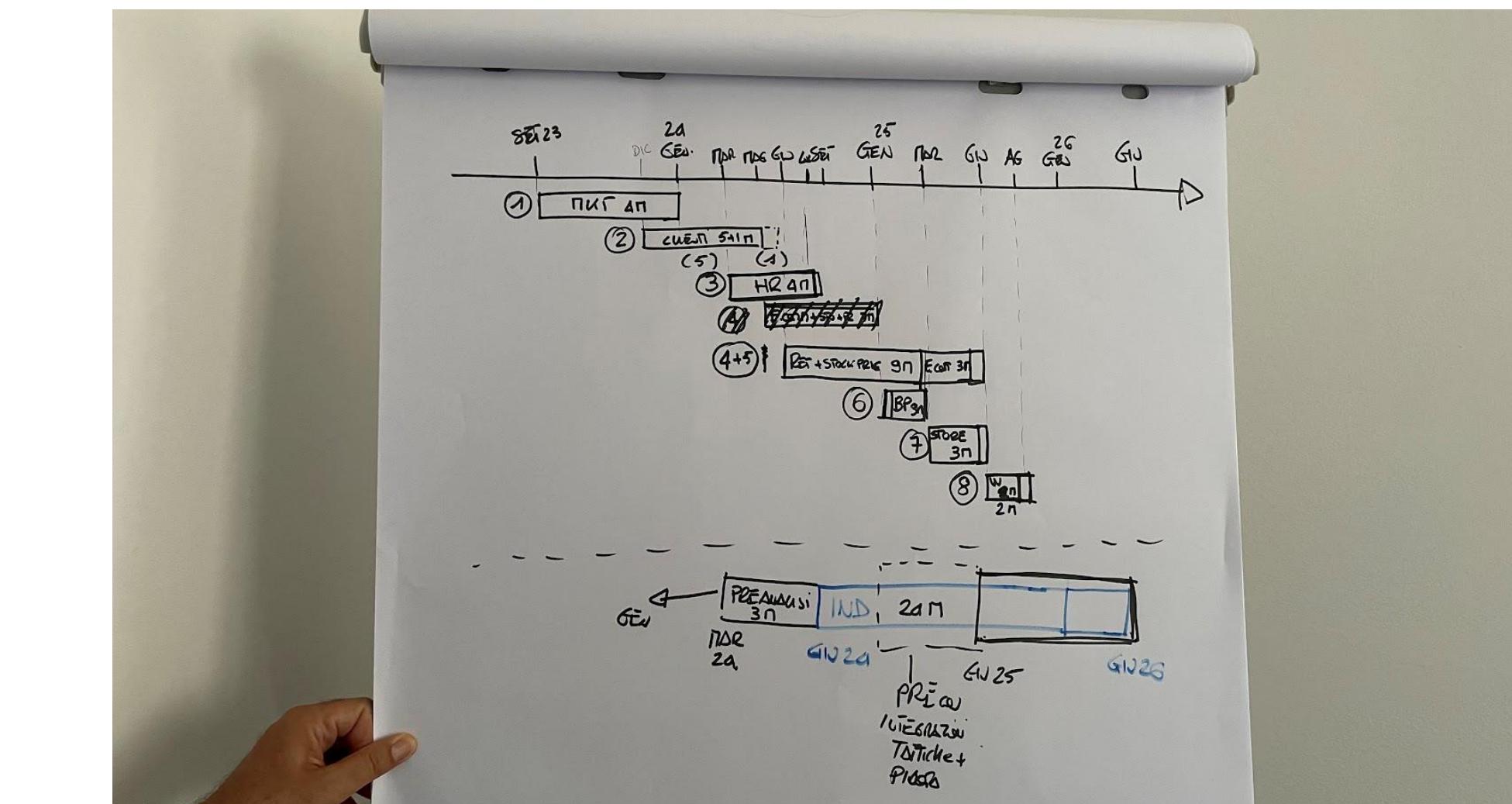
Design and Setup the **Architectural Landscape** in terms building blocks and interaction between these.

...We identified Data Subjects and design a high-level roadmap...

For each business area, we identify the main **data subjects** like **facts** and **dimensions** with their data sources (e.g. ERP system).



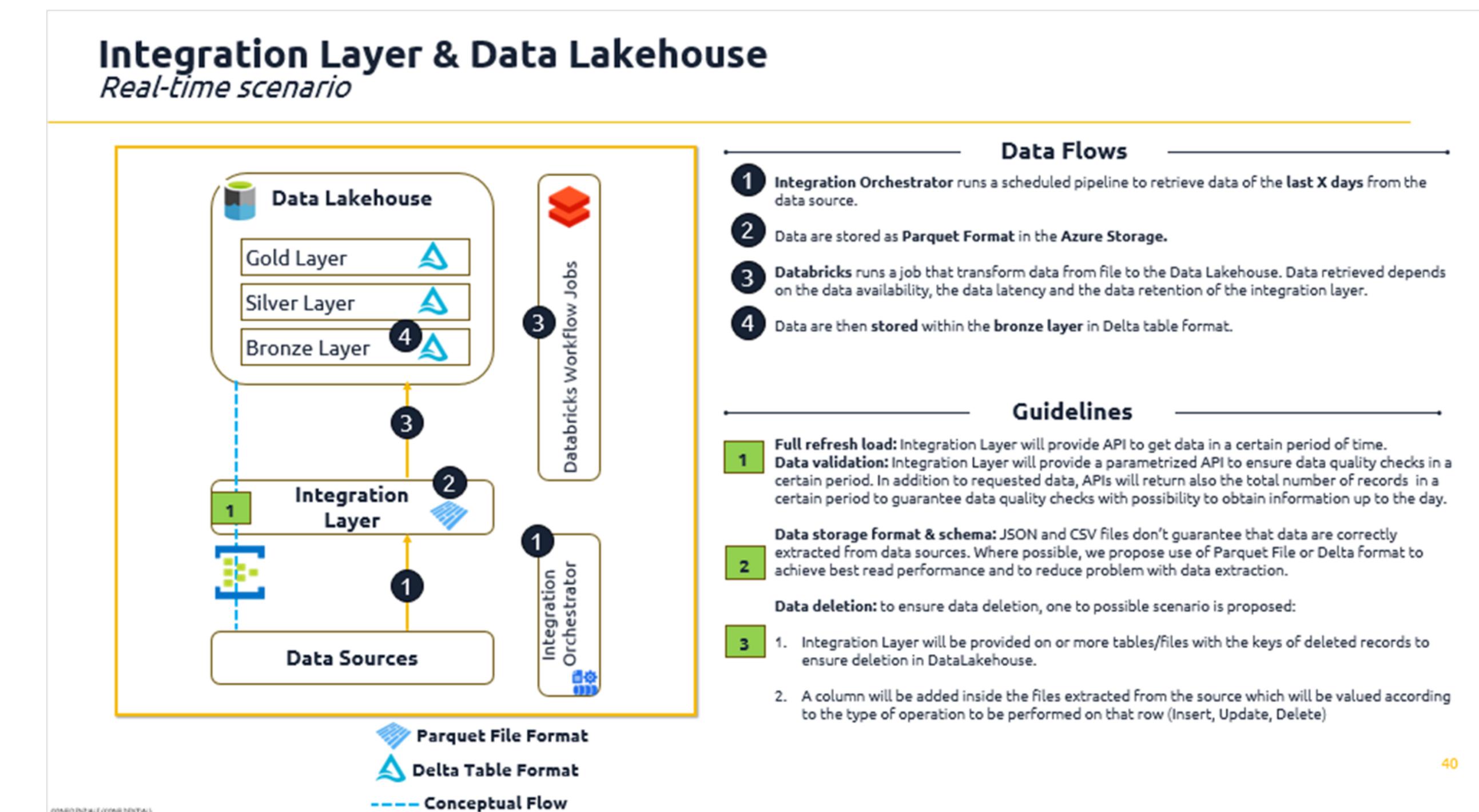
High-level Implementation Roadmap



consulting

...We identified the SLA required for each data subject and then we define an integration strategy to their data sources...

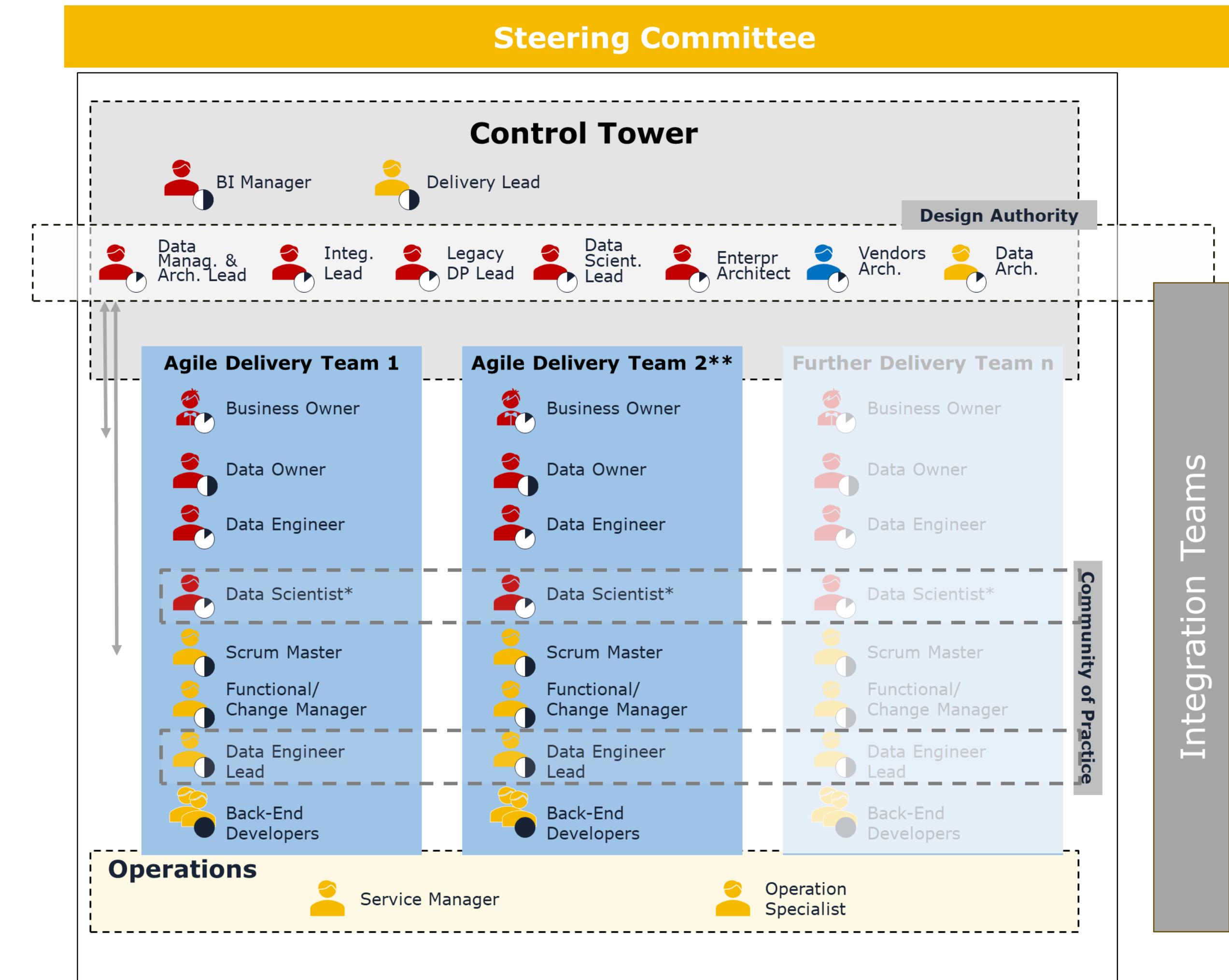
For each **data subjects** we identify the SLA required, where the most important one is the **data latency** (batch, near-real time, real-time). Latency needs to identify the best interface to the data source.



...We engaged the multidisciplinary team that will drive the Data&AI transformation with the Lakehouse implementation

Defining the **right team** with appropriate **business involvement** is a fundamental step for the successful implementation and adoption of the Data Lakehouse.

Below is an example of a customer team structure and framework.



Business Customer



D&A / IT Team Customer



Transformation Partner

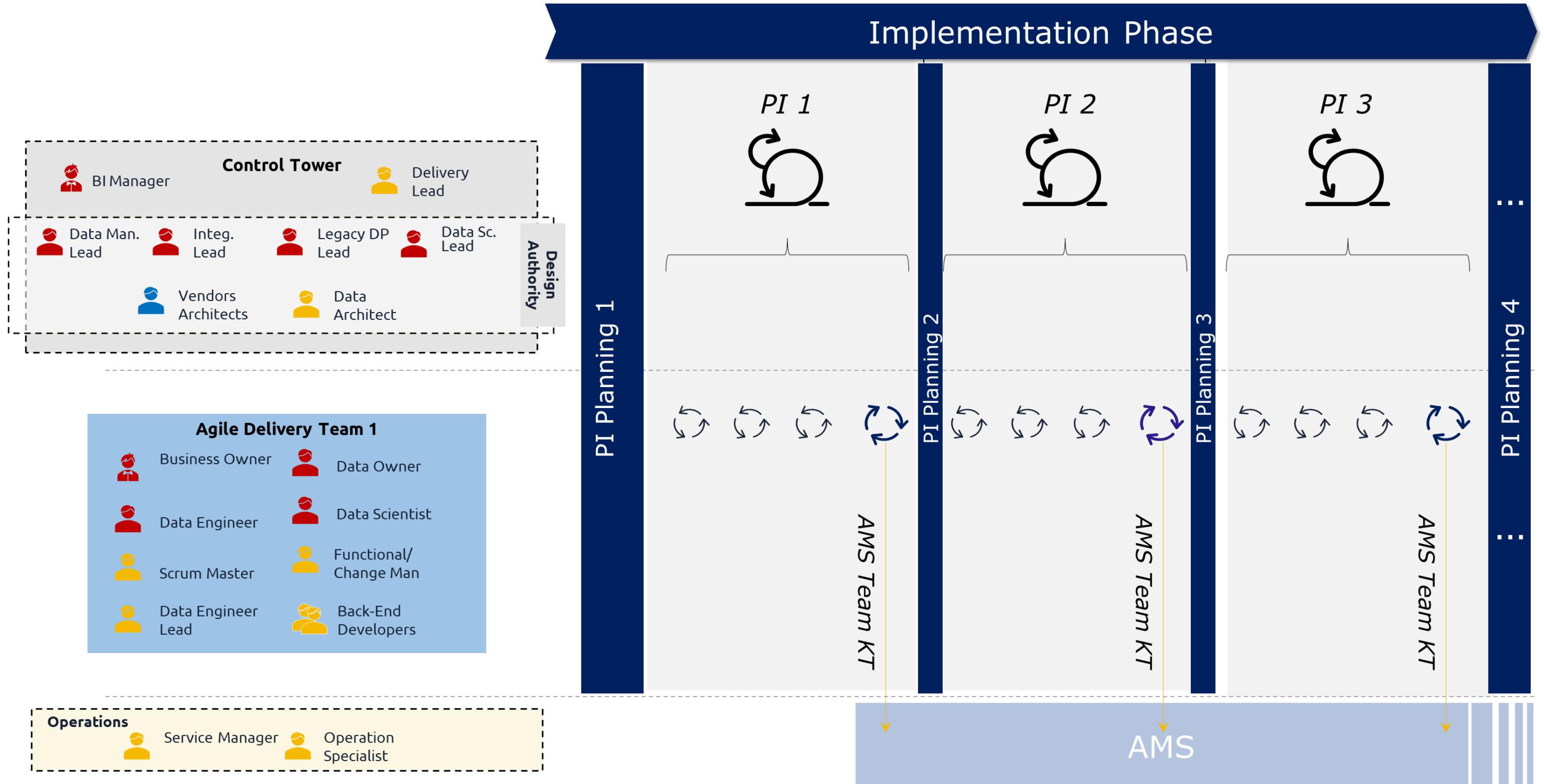


Vendor (i.e. Databricks, Microsoft)



...We defined the delivery process; in this case we adopted the scaled agile (SAFe) framework...

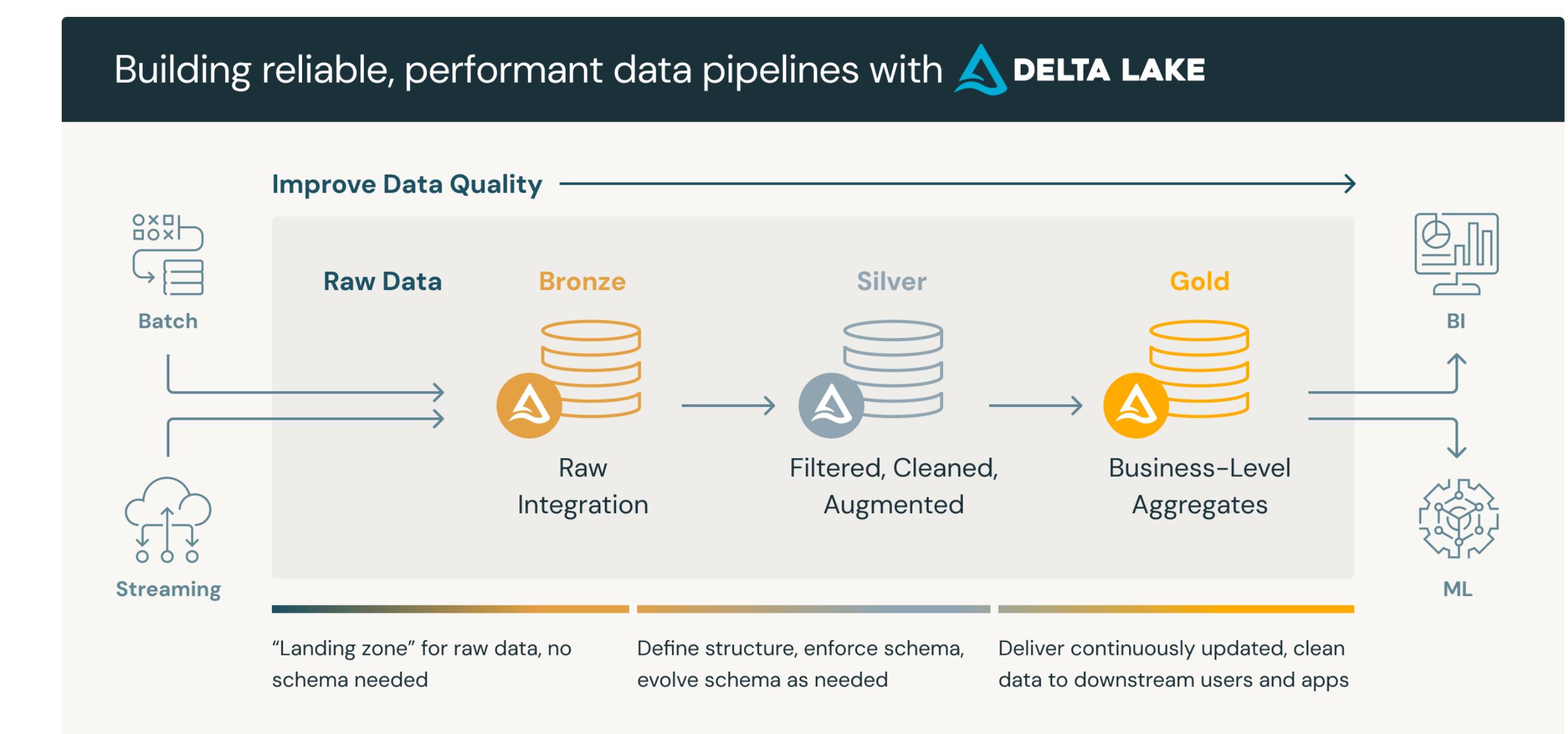
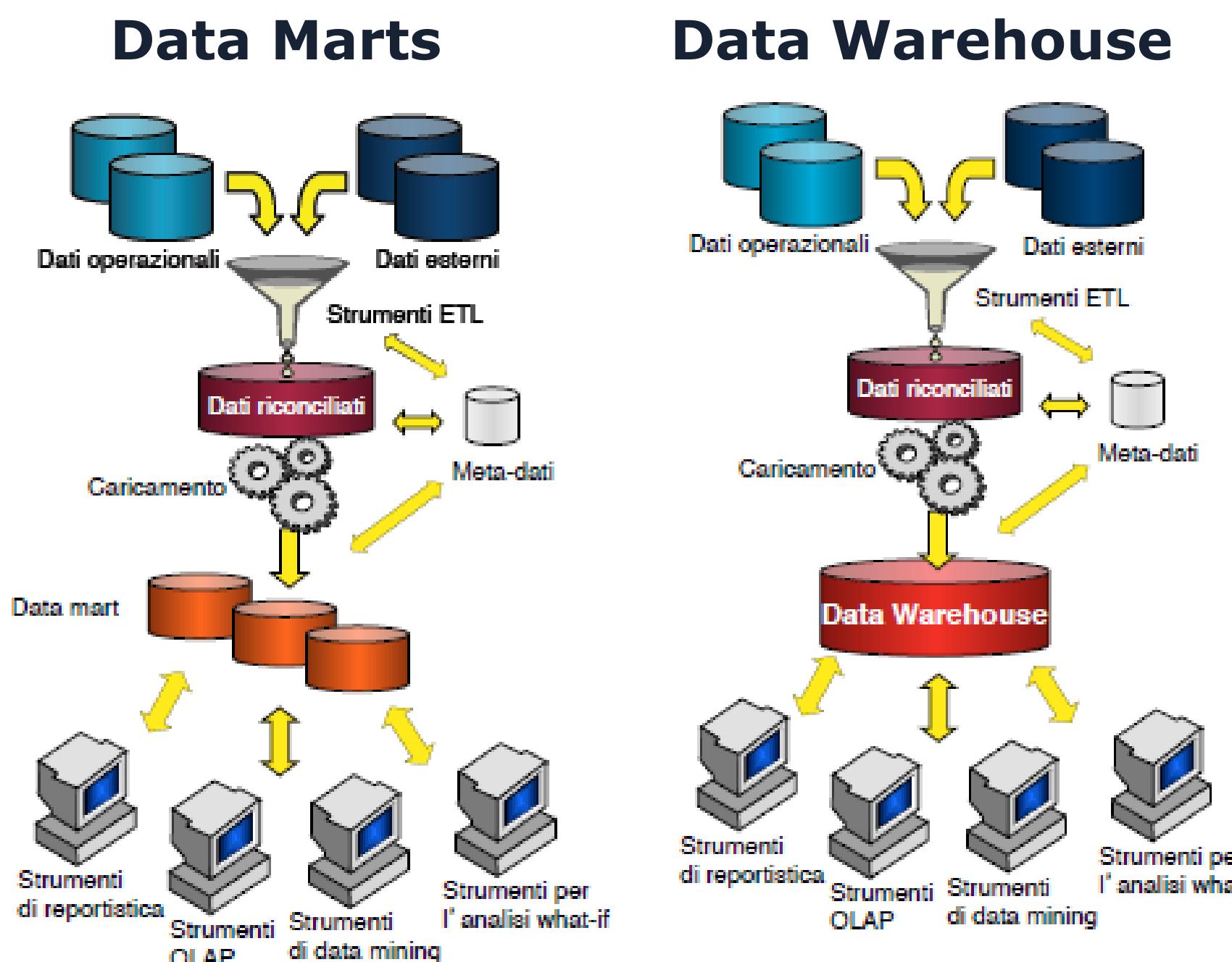
Scaled Agile (SAFe) framework enables us to define high-level program increments that are then delivered through one or more agile teams



31

...We adopted a standard how to organize data and transformations within the Data Lakehouse...

Medallion architecture (also known as multi-hop architecture) is a data design pattern used to logically organize data in the Lakehouse. It is essentially a rebranding of the well-known Data Warehouse design principles, which remain valid and recommended.

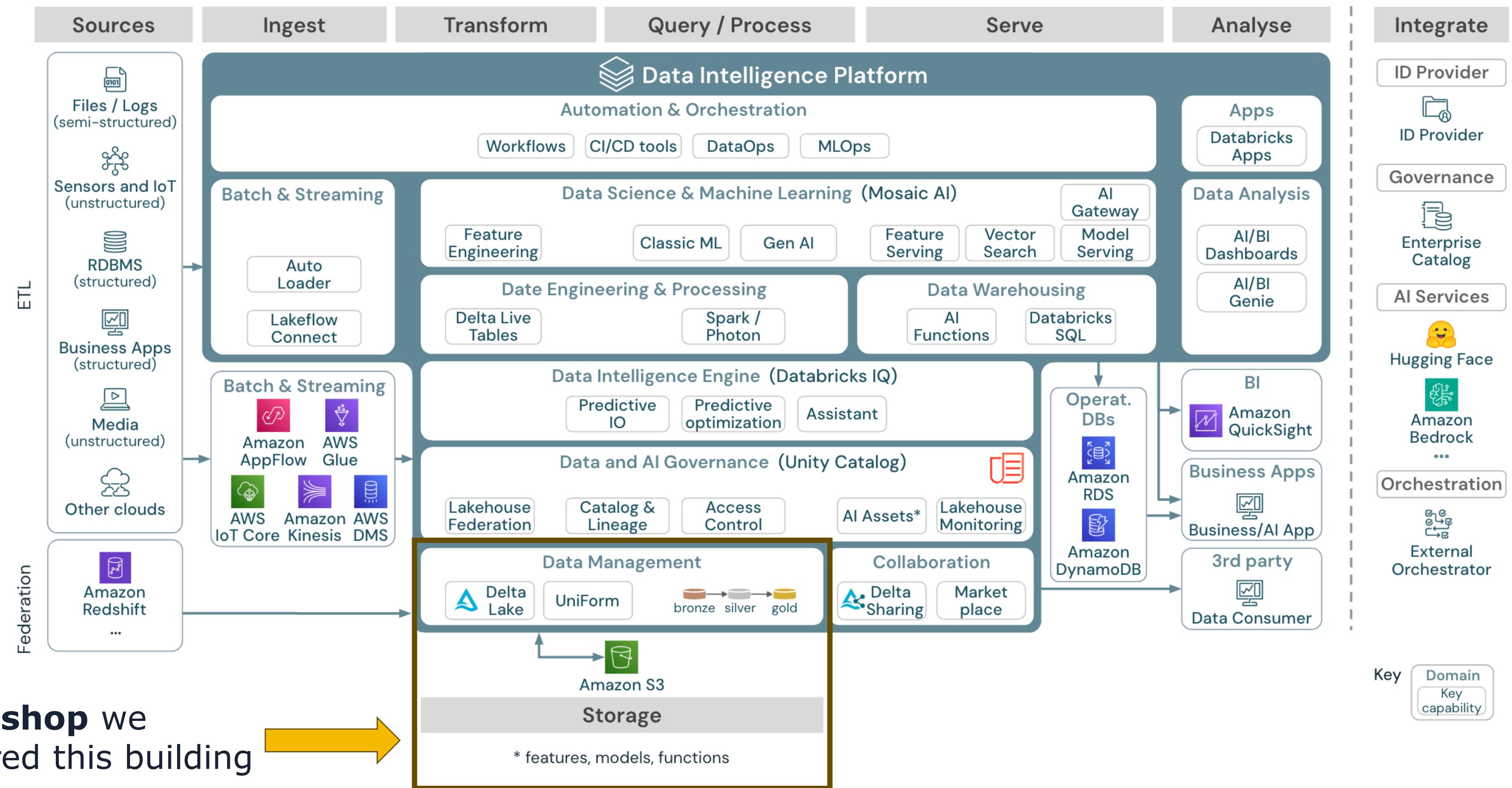


...We selected the architectural building blocks to use to satisfy the analytical business needs (use cases)...

Today's **real data platforms** are very complex and require a solid understanding of **data architecture**.

Below is an example based on **Databricks** and **AWS**, closely resembling one implemented for our customers.

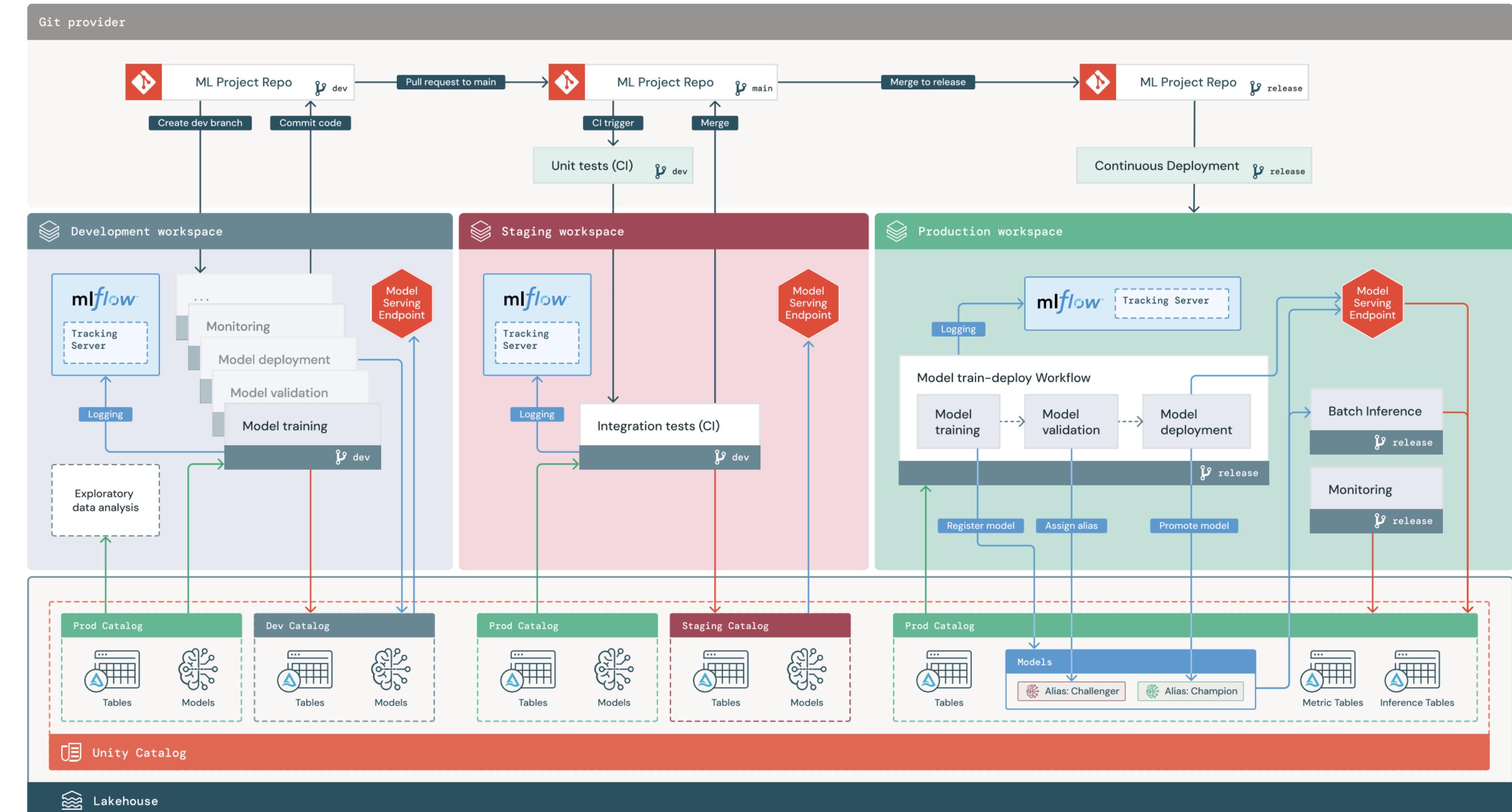
In this **workshop** we mainly covered this building block



...We defined a development strategy for all the Data Lakehouse artifacts: Data, Code and Models...

In a real scenario, it is very important to establish a development strategy to develop robust artifacts such as:

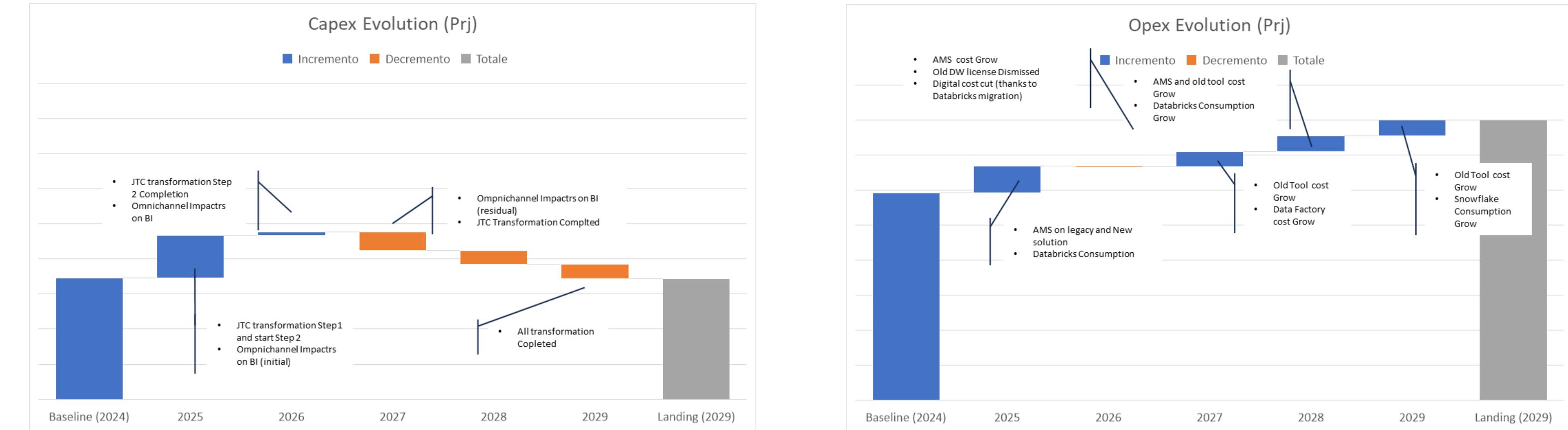
- **Code (DevOps)**: data pipelines (ETL) and model training code
- **Data (DataOps)**: structured (table), semi-structured (csv, text) and unstructured (pdf)
- **Models (ModelOps)**: both ML and LLM ones



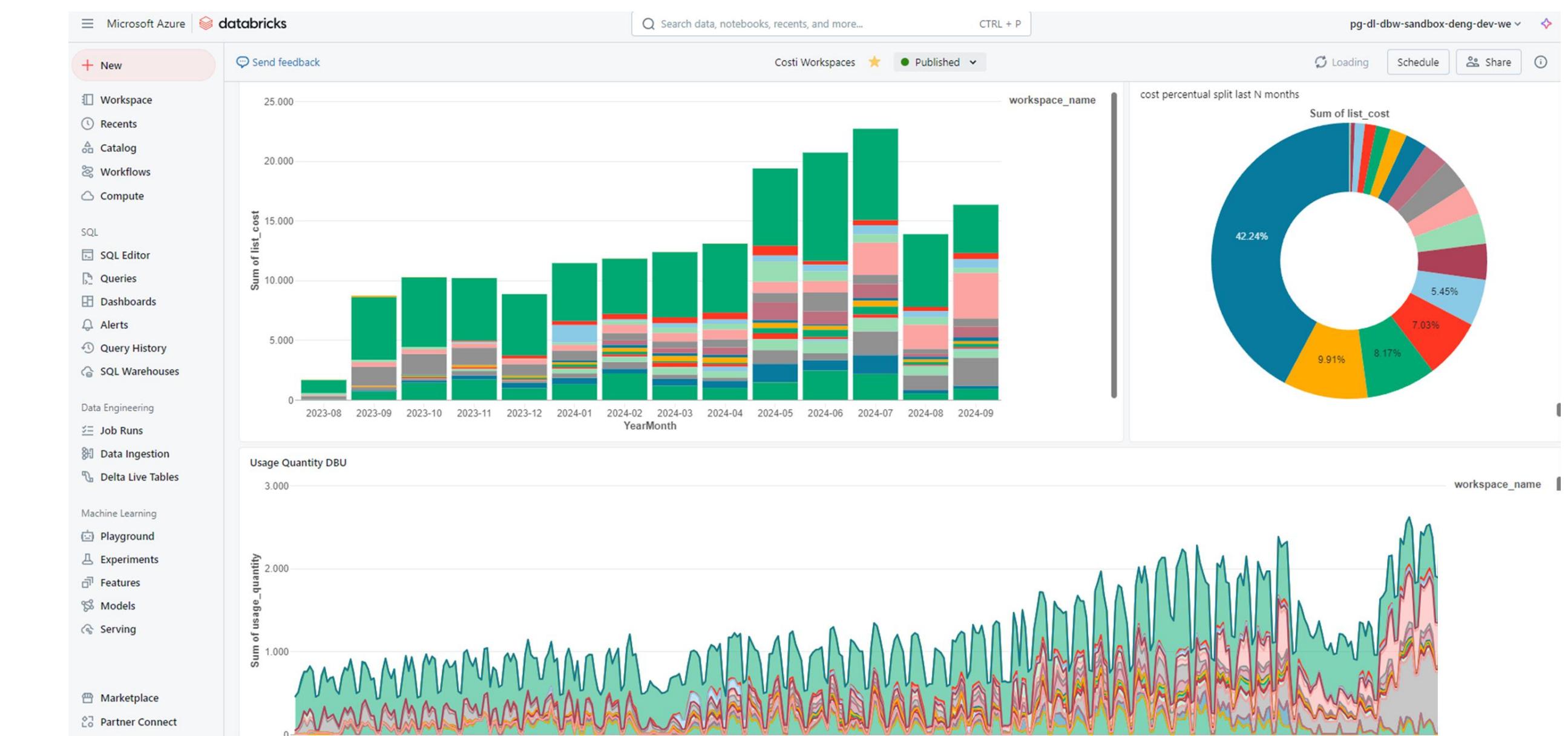
This is an example of Reference Architecture for **DataOps, DevOps and ModelOps** based on **Databricks**

...We simulated a 5 Year TCO and then monitored cloud cost consumptions...

Estimated Total Cost of Ownership (TCO) is important to evaluate the initiative **costs** and **Return On Investment (ROI)**



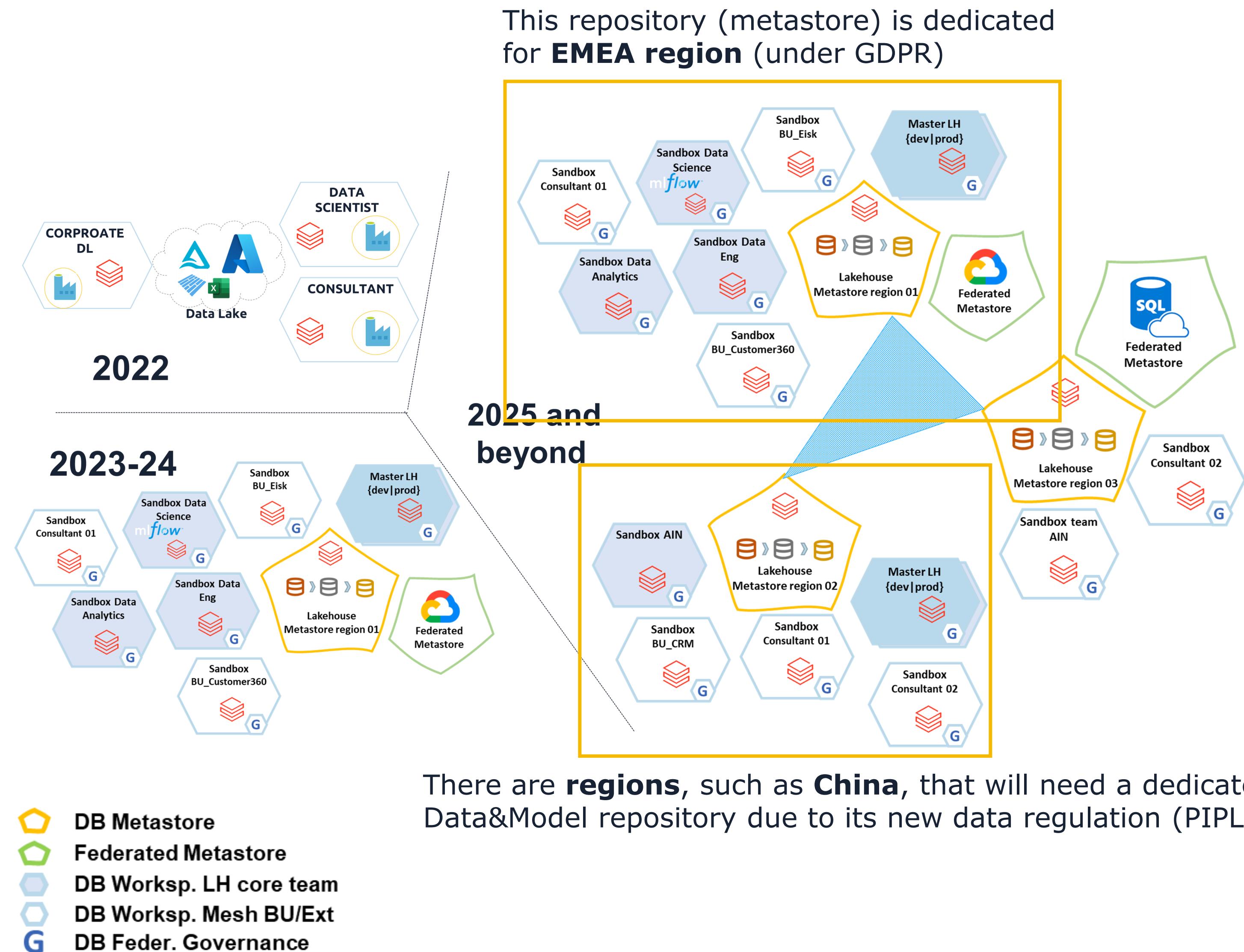
Cost controlling (FinOps) is an important pillar for Data Lakehouses, since these are mainly bases on **pay-per-use** pricing model and costs can be very high.



...We defined an incremental evolution roadmap for data and model consumptions to other regions and brands.

We designed an **evolution roadmap** to provide data to other regions and other brands.

In this case, several data and model repositories (also called metastores) have been enabled to cover multi region regulation requirements (GDPR, PIPL)

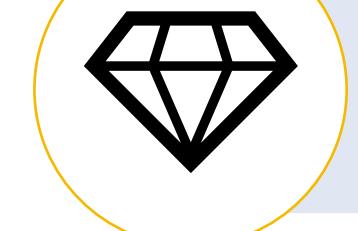


Takeaways



...What we learned in this lesson? What are the takeaways?

There are few takeaways from this lesson:

-  **Data Lakehouse** will be one of **the most important design paradigm** in the future
-  **Decoupling Storage and Compute** is one of the most **revolutionary feature** of Data Lakehouse, since enables seamlessly resource scaling
-  There are several formats, but **Delta Lake** will be probably be the **most common** one
-  There are several players that allow to implement a Data Lakehouse, where **Databricks** is the **most interesting** one
-  A **Data Lakehouse initiative** requires a first step to **design** and **plan** its **implementation** that could be very **hard** but it is also the **most stimulating part**



Q&A Session

iconsulting



Andrea Carmè

Senior Manager

a.carme@iconulting.biz

Thanks for your attention

