

A DATA LAKEHOUSE STORY

09 April 2024

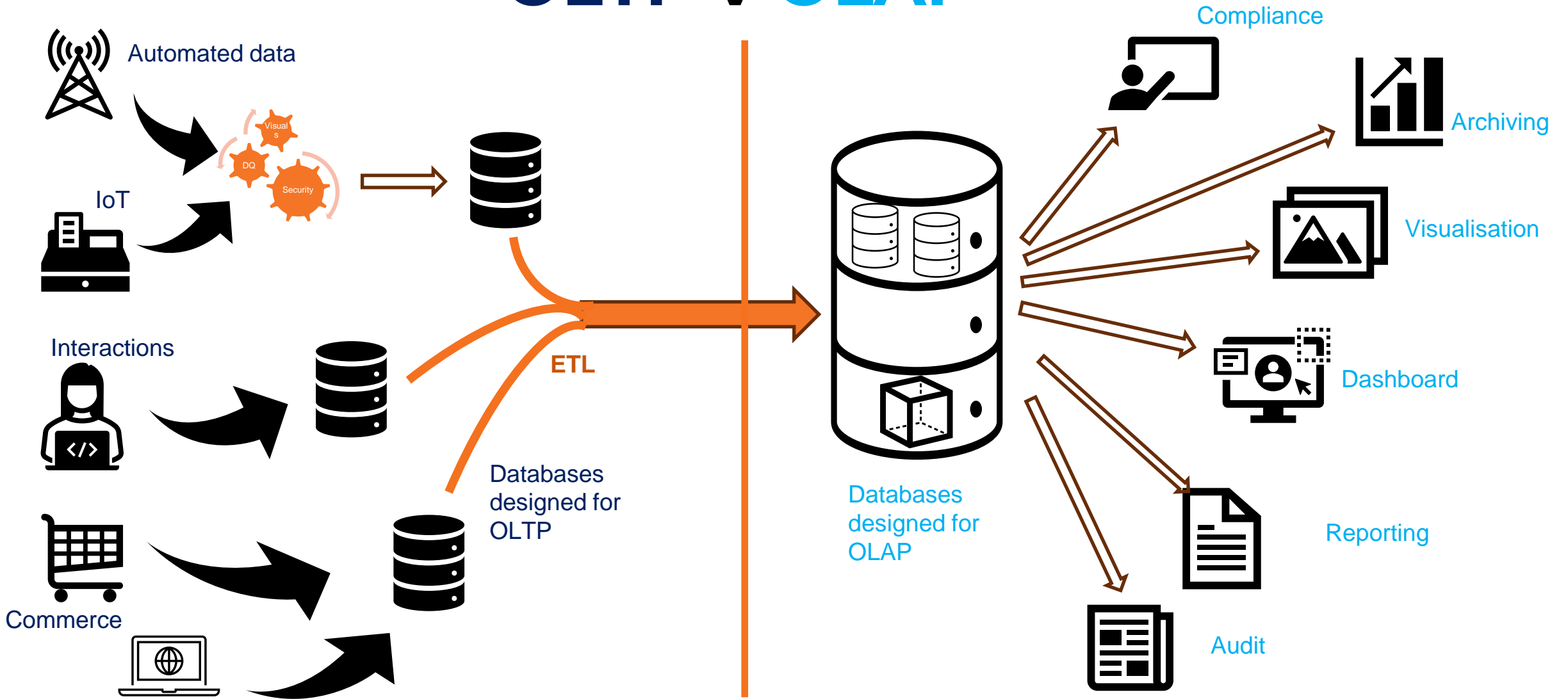
Data Architecture Event, Corinium

Michel Gehin, Data Architect

Case Study - how new technology has enabled:

- **the deployment of a Data Lakehouse**
- **Its multipurpose**
- **Business benefits and savings**

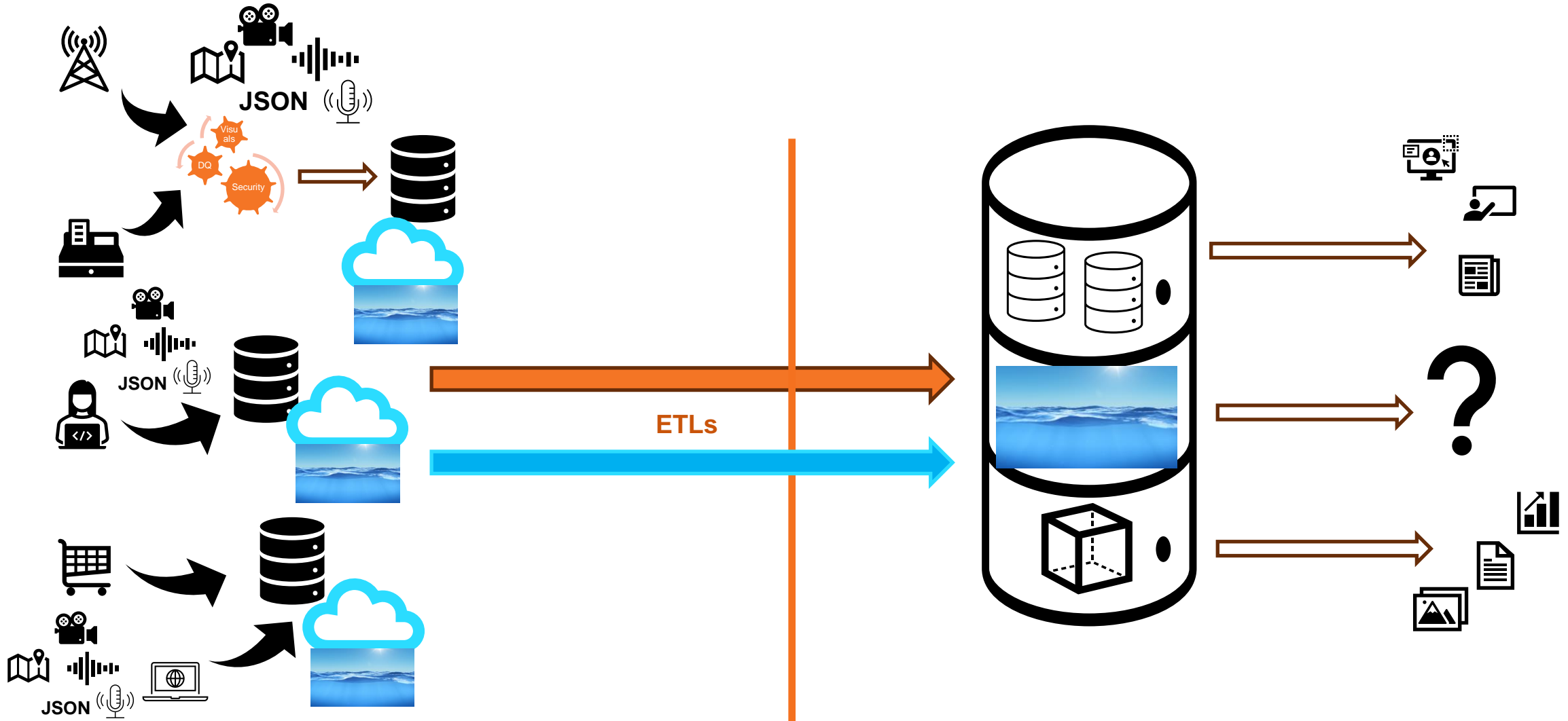
OLTP v OLAP



Fast transactional capabilities
Creation and updates of individual records
Lots of narrow transactions on a small amount of data

Fast querying of large datasets
Smaller number of attributes over aggregated stores and cubes.
Paired with a BI tool, semantic layers

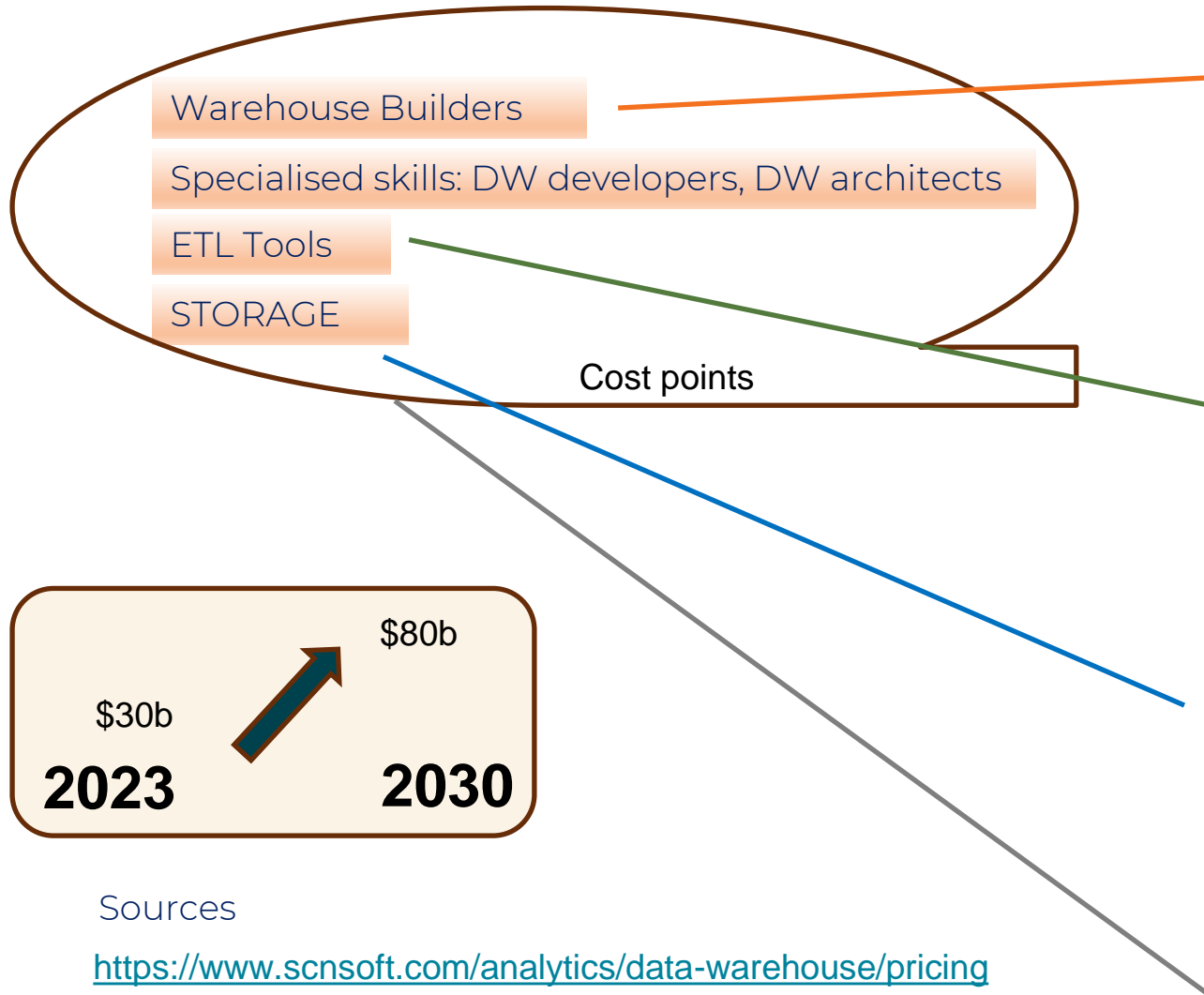
New data sources and formats



+ Unstructured data
+ Semi-structured data

Traditional DW
+ Data Lake

OLAP-specific Spend



Implementation Costs

A typical NZ company, with 5 internal sources of data (ERP, CRM etc), with just structured data batch-processed, hybrid of cloud and on-premise DW stores, manual DQ tools and limited ML/AI capabilities will **spend between \$50,000 and \$200,000 building a data warehouse**, not including BI tools

ETL Costs

will depend on organization, between \$800 and \$8,000 per month

Running Costs

Current estimations:
ON-PREMISE: \$1,000 per terabyte per month

CLOUD: up to \$84 per terabyte per month in the cloud (hot storage + basic compute – cold is cheaper).

The average size of a DW worldwide is around **100 terabytes**.

Specialised Staff Costs

\$???k per month

Oracle Warehouse Builder

Amazon REDSHIFT

Snowflake

WhereScape RED

IBM Infosphere

Talend

SAS Data Management

Google Big Query

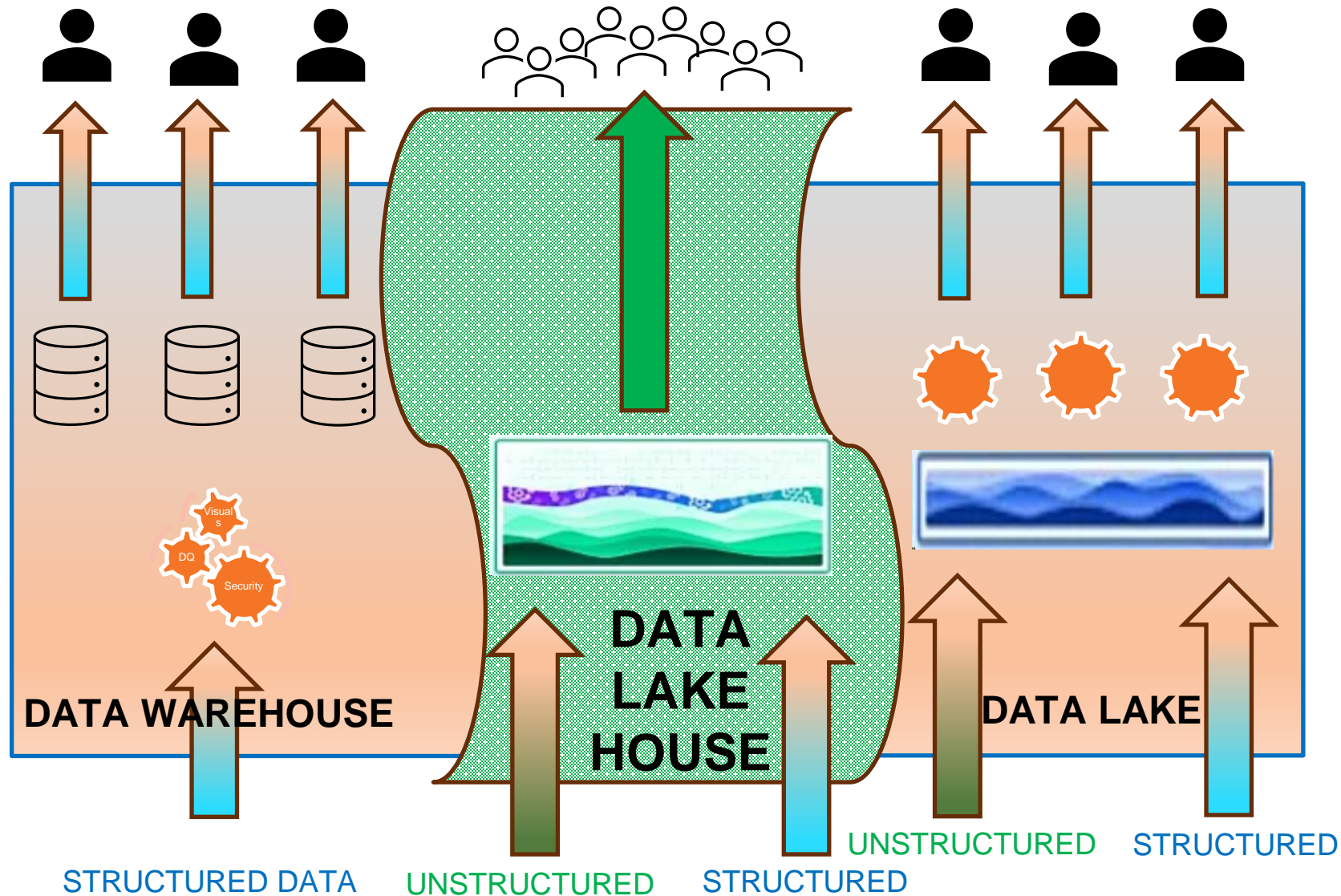
Microsoft SSAS

Sources

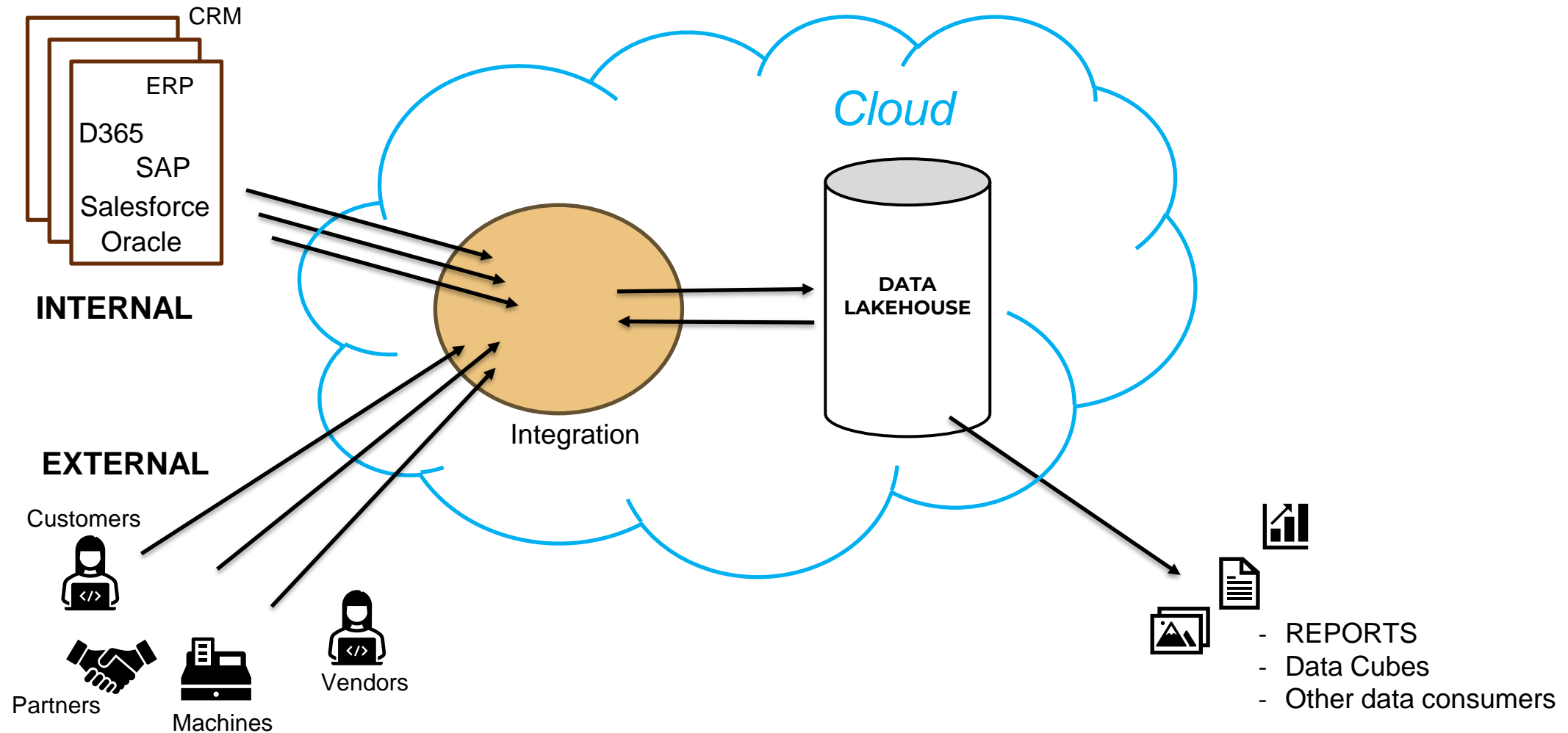
<https://www.scnsoft.com/analytics/data-warehouse/pricing>

<https://www.integrate.io/blog/the-true-cost-of-a-data-warehouse/>

The Data Lakehouse



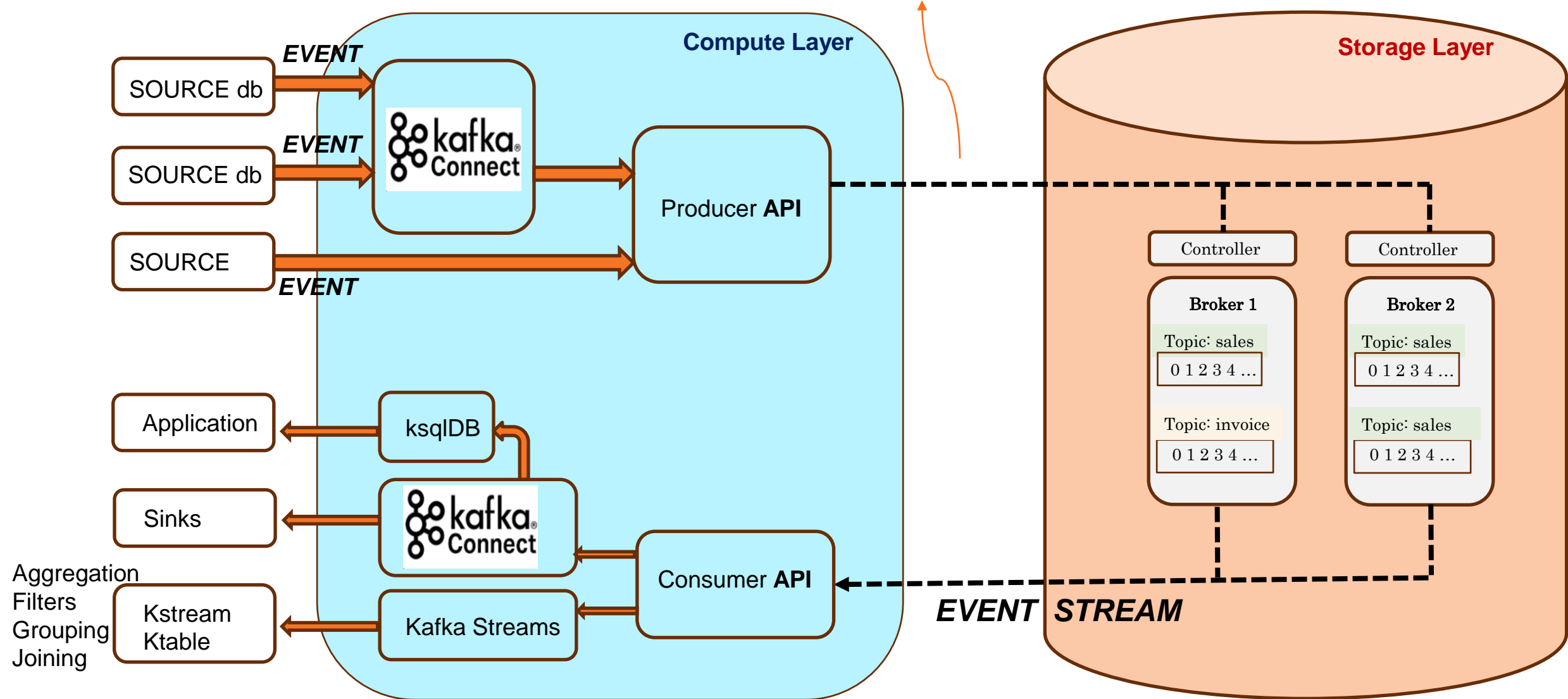
Integration and Lakehouse



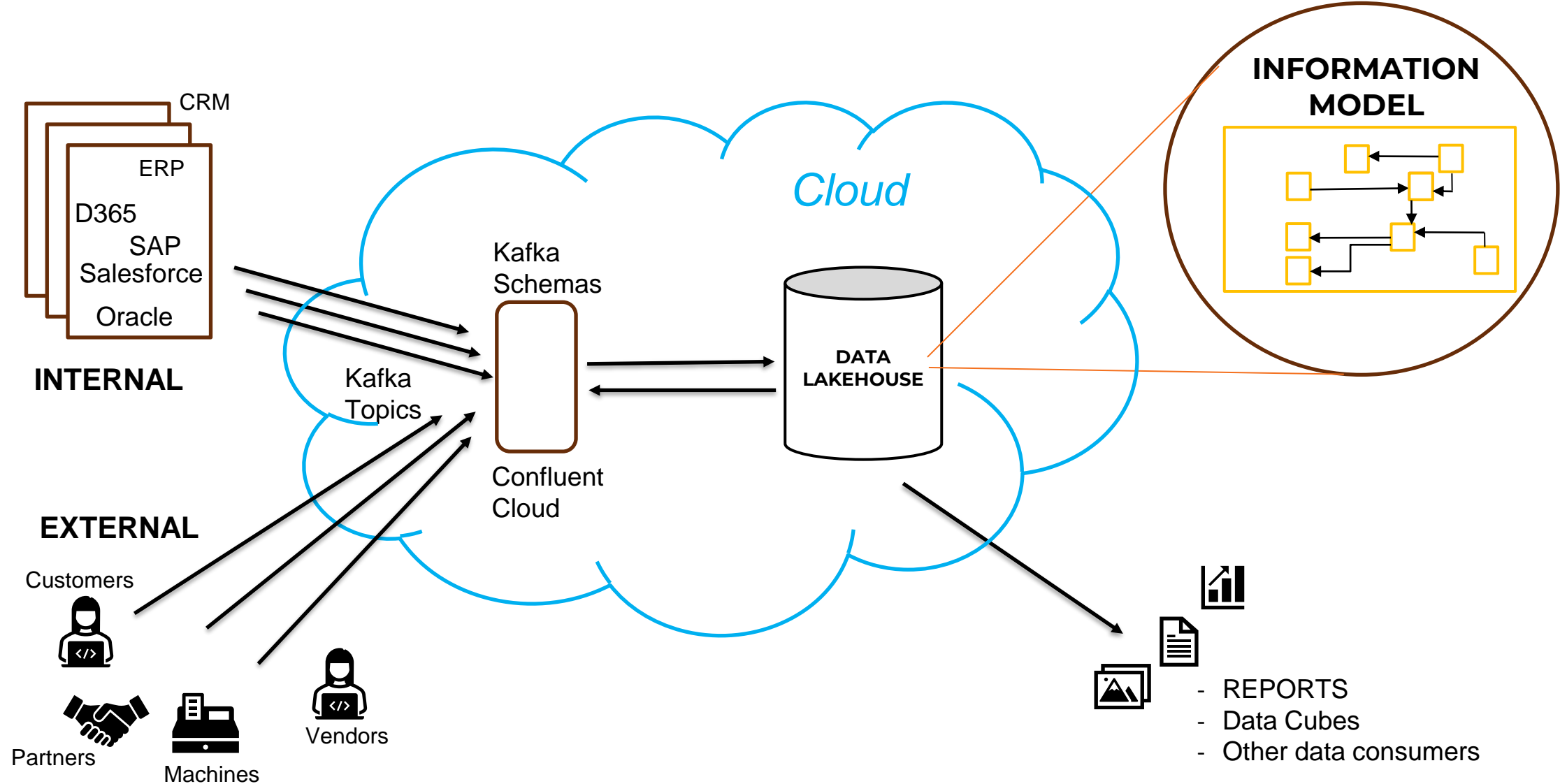
Kafka: fast integration

RECORD

| timestamp | key | value | headers |
|-----------|-----|-------|---------|
|-----------|-----|-------|---------|

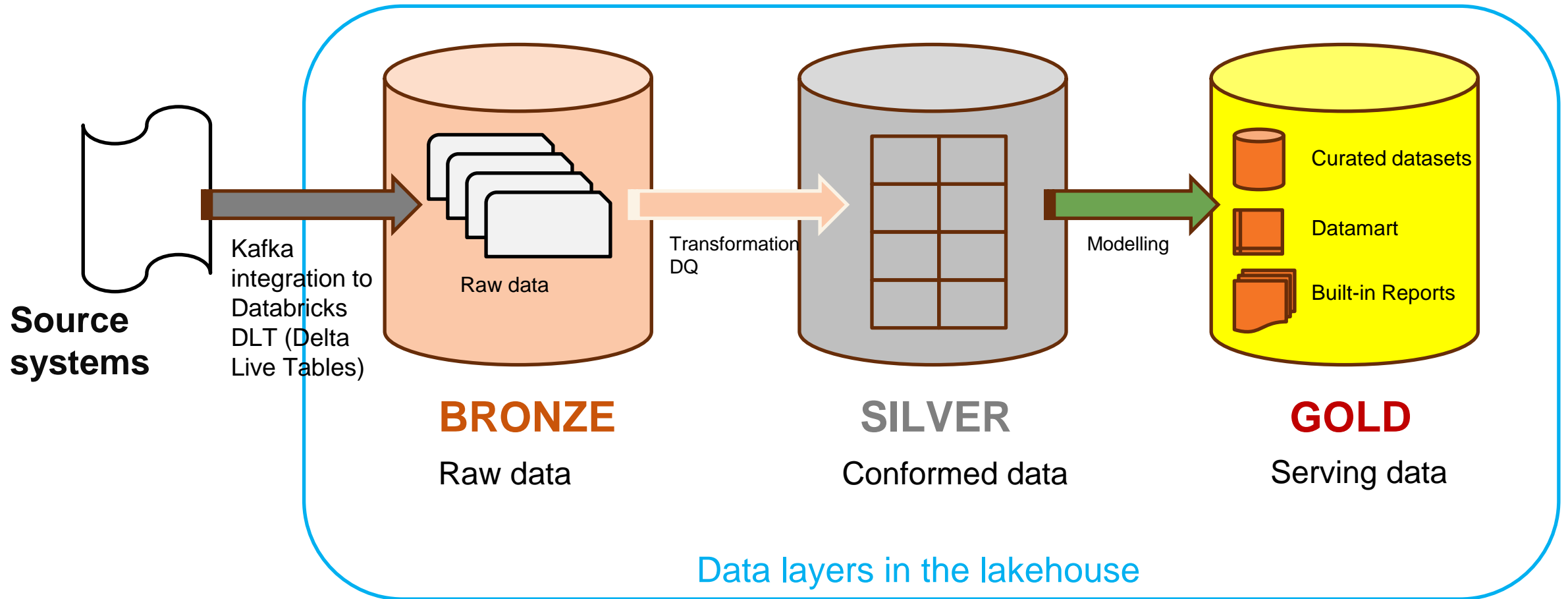


Data Architecture Overview

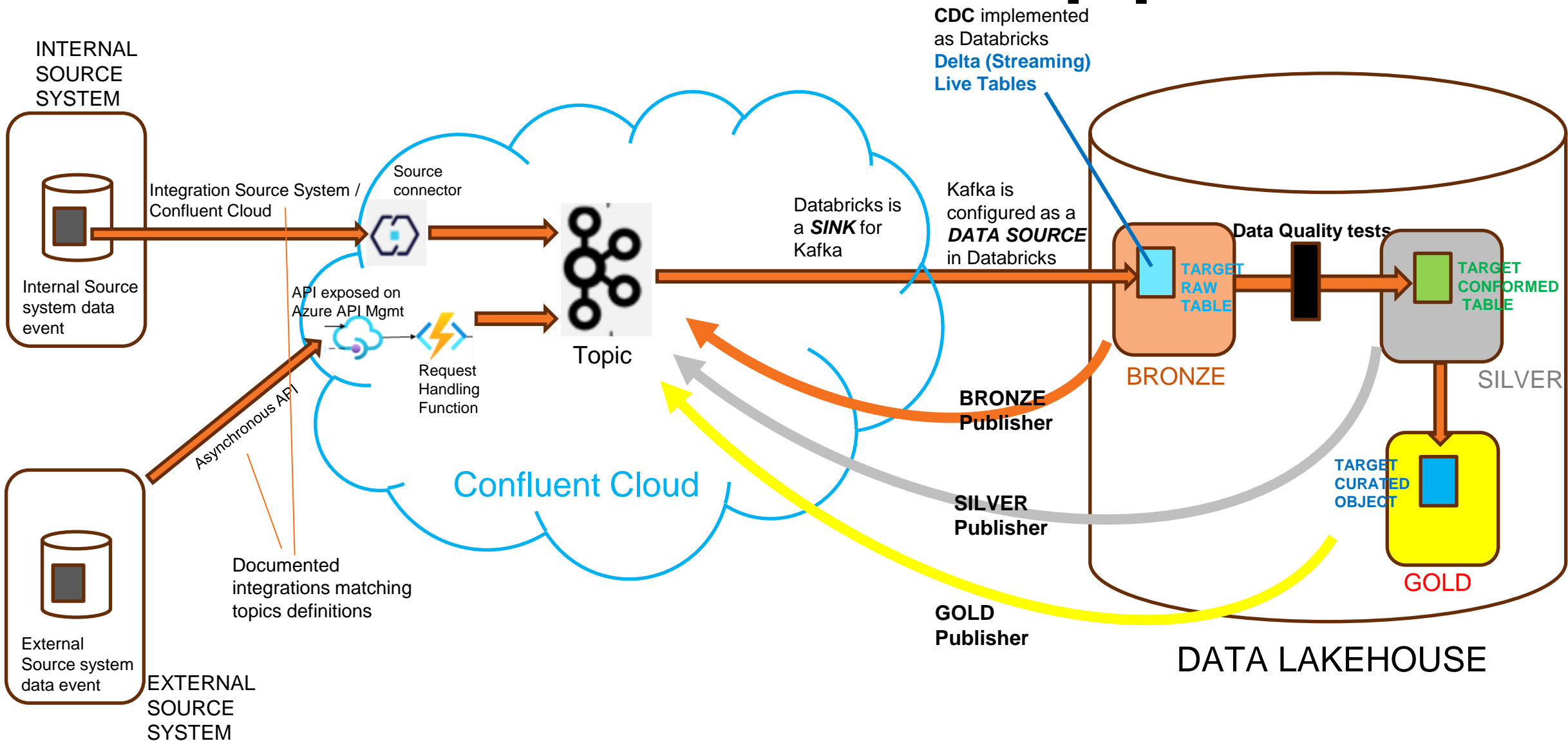


The Data Lakehouse – the physical view

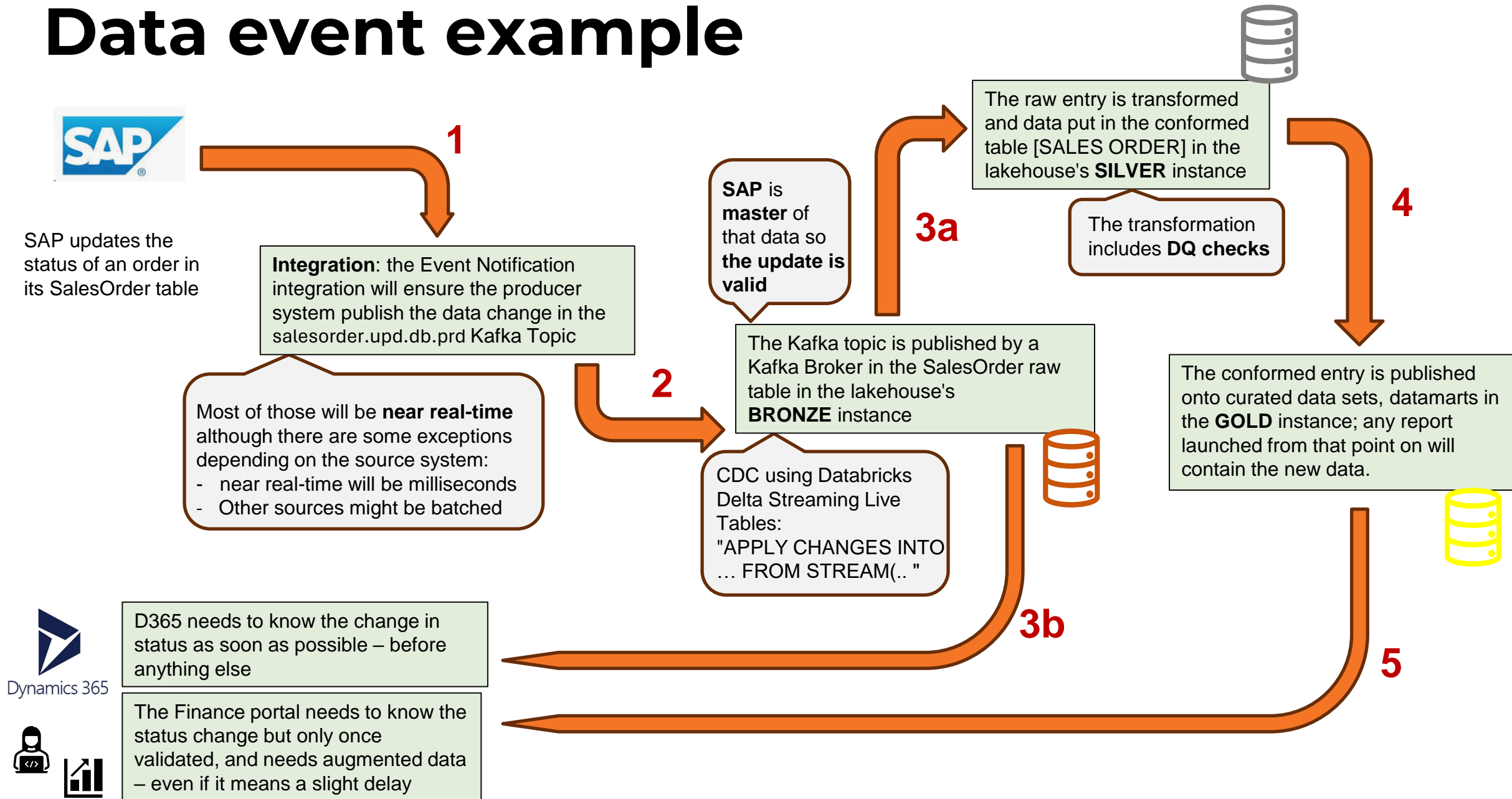
The Data Lakehouse is composed of several physical asset layers, as follows:



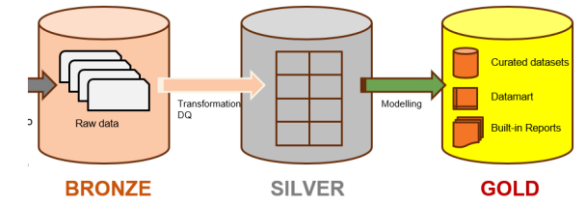
The Data Lakehouse – data pipelines



Data event example



How to consume data from the GOLD layer of the Lakehouse



CURATED DATA SETS

The GOLD layer will include curated data sets designed and built according to core business needs – for example, total sales with stores physical locations together.

DATAMART ACCESS

The data in the conformed information model will also be available, with data consumers able to join tables together with knowledge of the documented Information Model.

BUILT-IN REPORTS

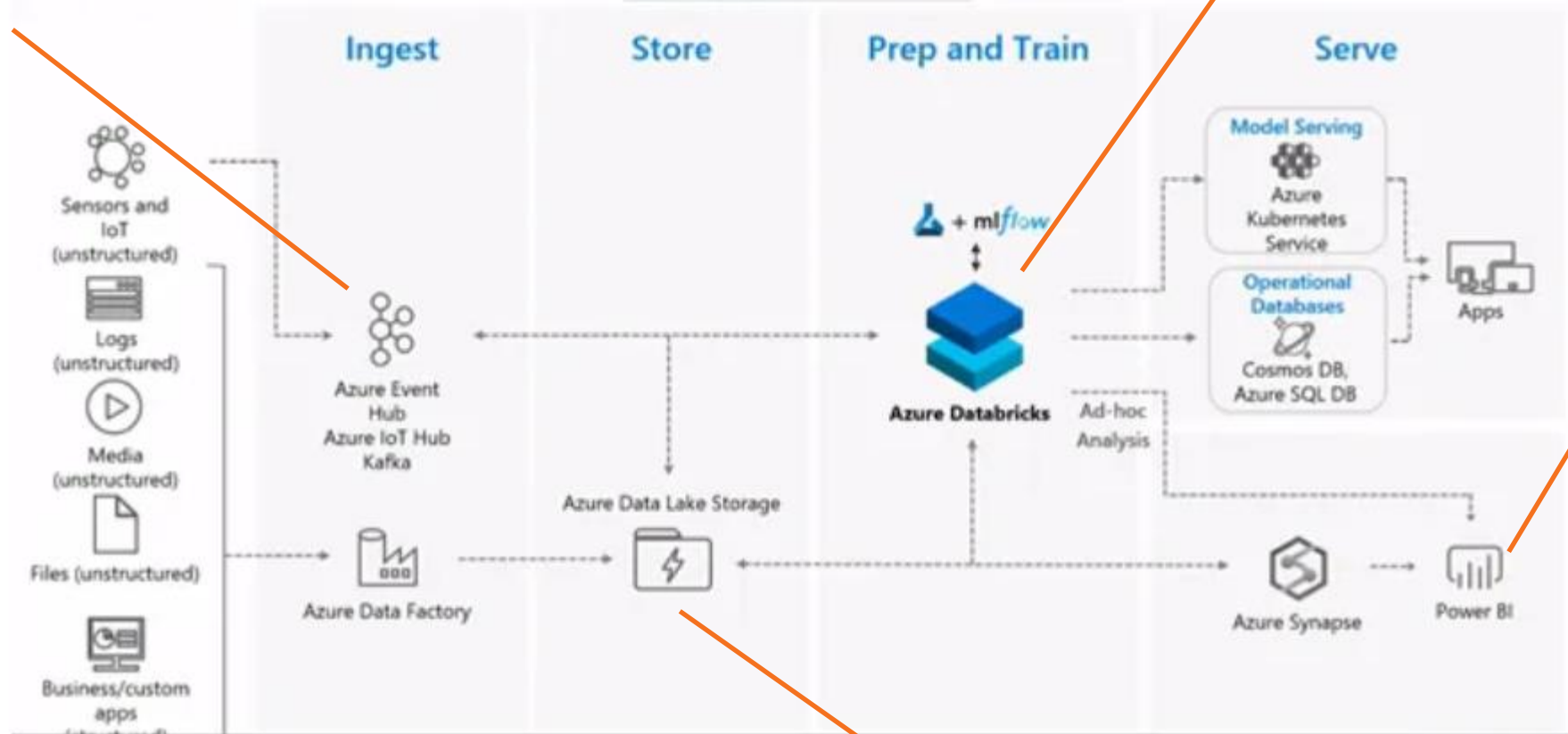
The GOLD layer will also contain ready-to-use reports

The data is DQd

The data has been classified and tagged, with security classification, privacy, sensitivity, domain, version, etc.

Databricks

Kafka integration –
Kafka set up as
source, while
Databricks set up
as sink in Kafka



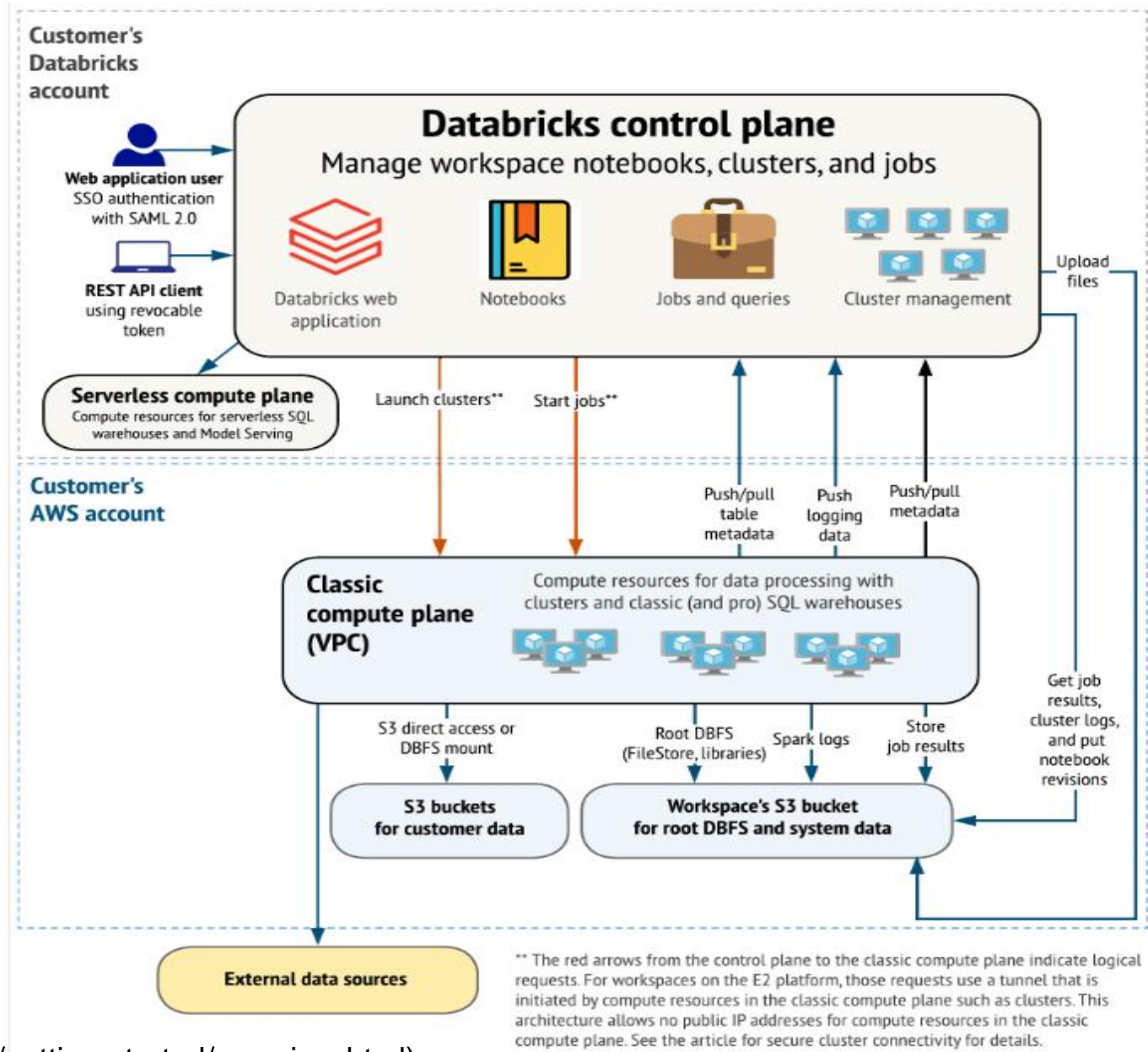
Includes tools for analytics, ML, AI, ..

Connectivity for
reporting tools
such as Power BI

Physical Data Lake: Azure Gen 2 Data Storage

Includes backup, recovery, DR

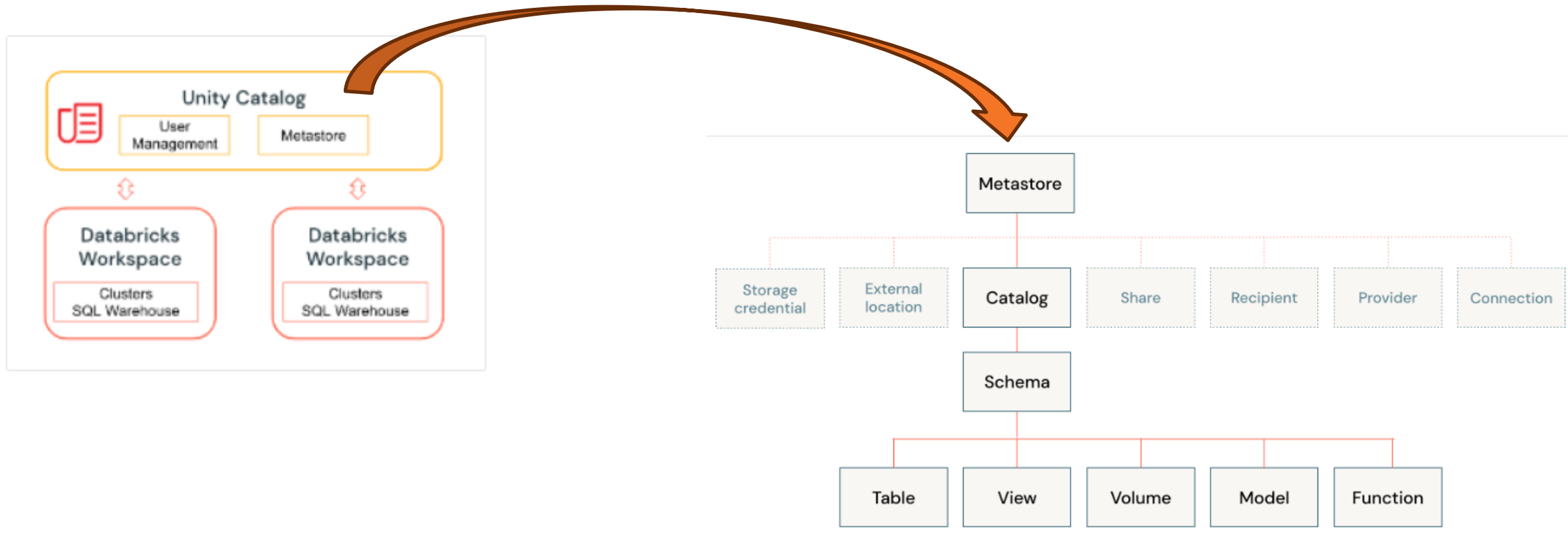
Databricks Architecture



(SOURCE: <https://docs.databricks.com/en/getting-started/overview.html>)

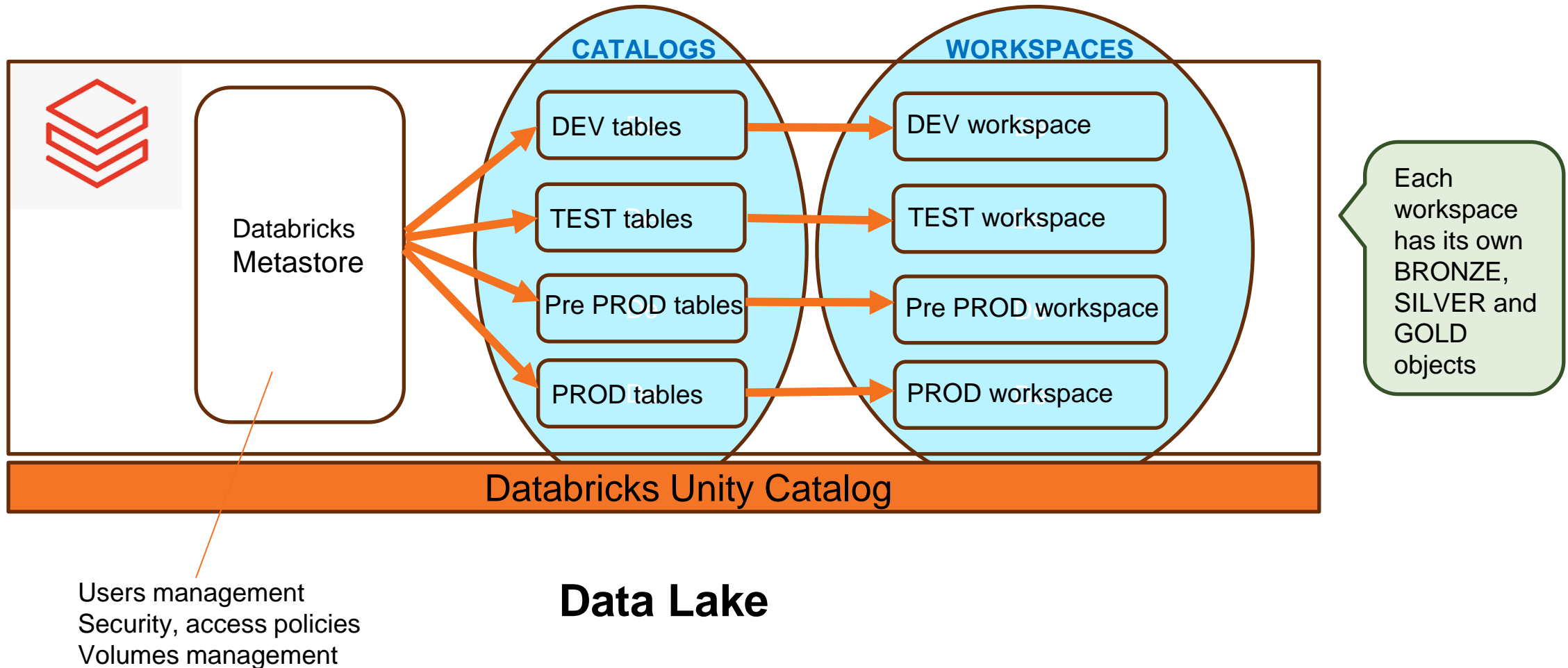
Databricks Unity Catalog

A one-stop, single-view management tool for all workspaces, for all layers, across all environments



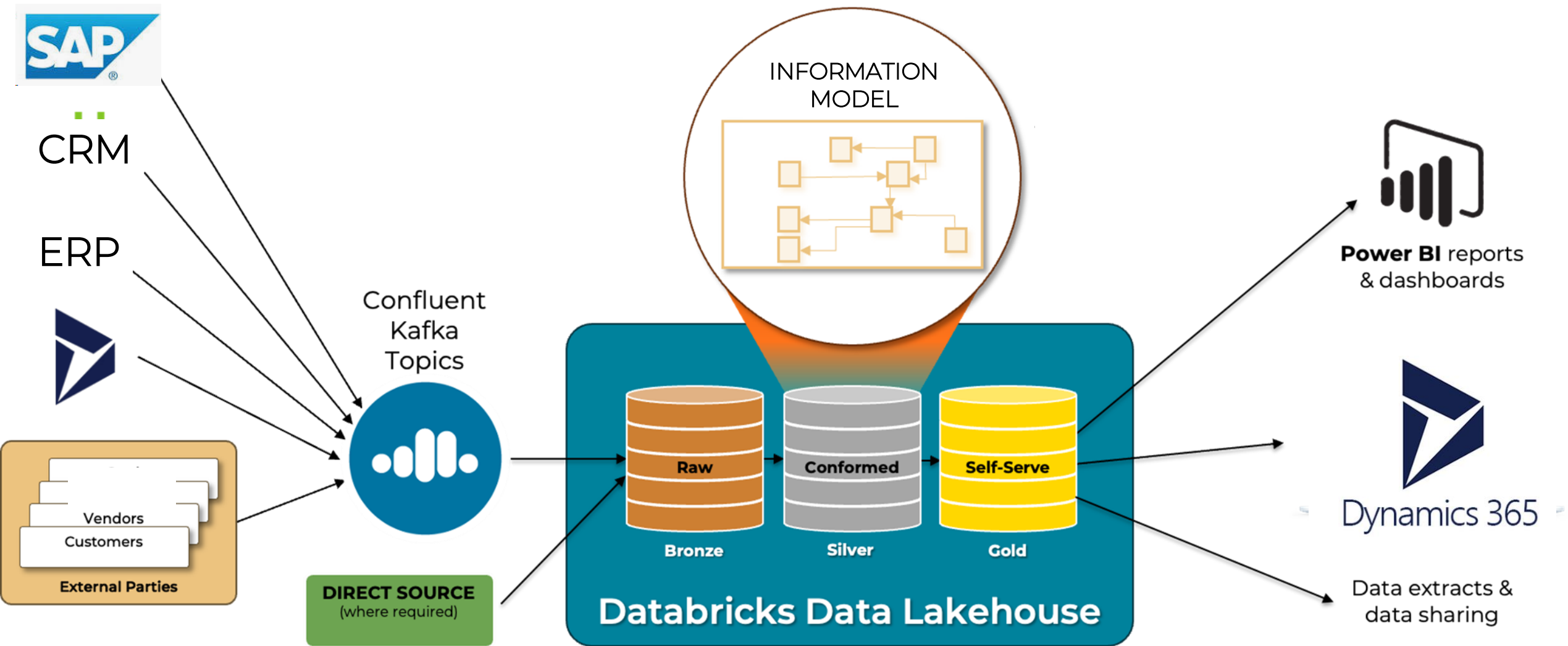
Environments on Databricks

Leverage Databricks Unity Catalog to manage all 4 environments, each with 3 layers, with one Data Lakehouse / one Catalog
Unity Catalog enables the separation of data objects across catalogs and workspaces, isolating catalogs from one another.



THE LAKEHOUSE

FOR BOTH **OLTP** AND **OLAP**



Summary

Smart use of tech has enabled a more simple, coherent, data landscape:

- Each function carried by a **component that fits** the business needs
- **Integration** between each component is key
- Track **un-necessary duplication** of function and eliminate
- Pay for the **speed** and **volume** you need now, ensure scalability for future needs
- You don't necessarily need best of breed, just **best fit** to your business
- Today's offerings in **cloud** and **big data** enable you to pick each component separately, as long as the integration is at least "almost native"



Thank you