

Datalake y DataWarehouse

BASES DE DATOS PARA IA

QUE ES UN DATAWAREHOUSE?

- CONTIENE INFORMACIÓN DE BASES DE DATOS TRANSACCIONALES
- SUMA INFORMACIÓN EXTERNA PARA VISTAS 360
- ORIENTADO A BIG DATA
- DATA ESTRUCTURADA DEL TIPO RELACIONAL
- ACID COMPLIANCE ES MENOS ESTRICTO
- AYUDA A CONSEGUIR INSIGHTS DE NEGOCIO EN INTELIGENCIA DE NEGOCIOS O **BUSINESS INTELLIGENCE (BI)**



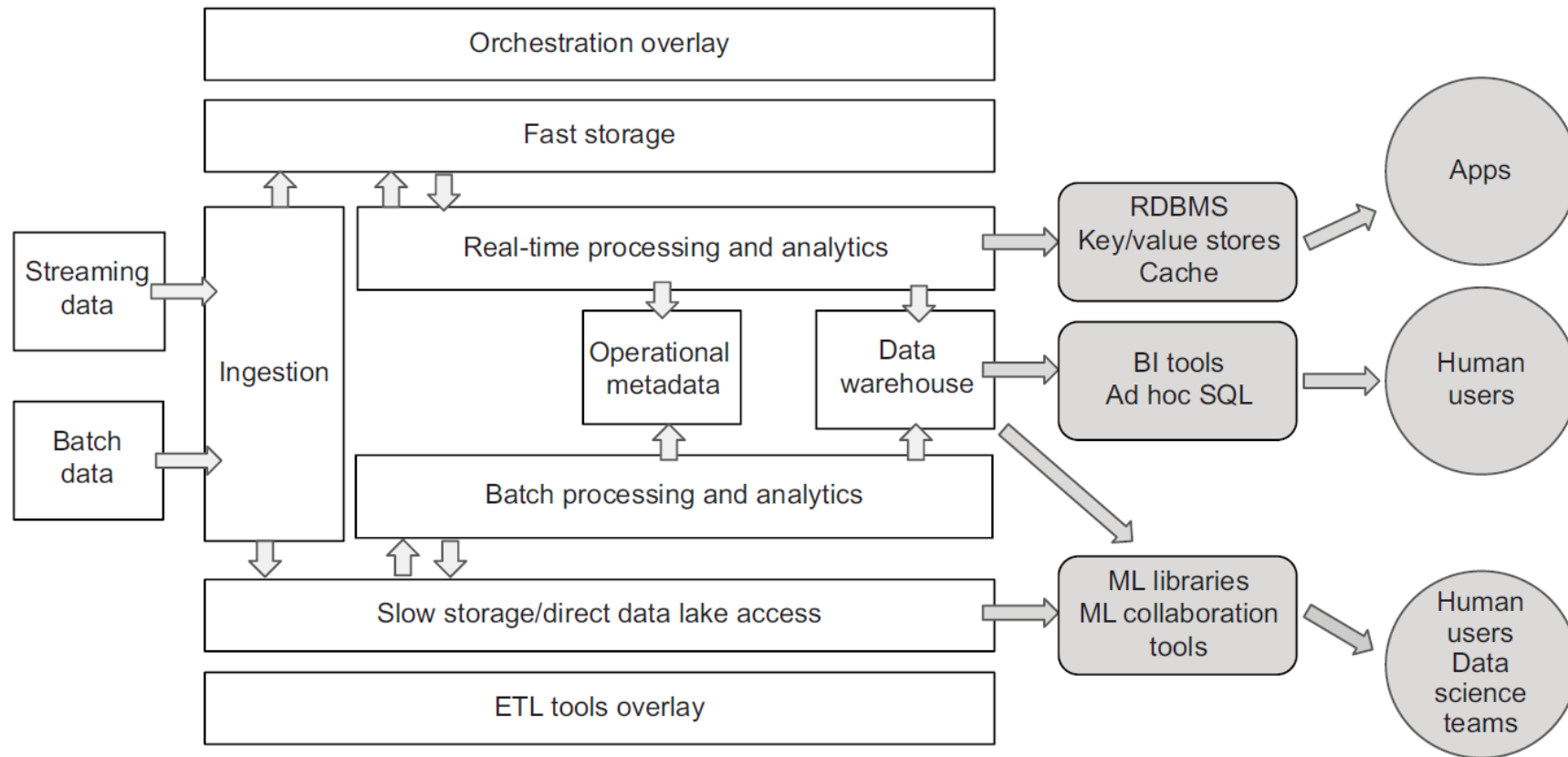
Google BigQuery



CARACTERÍSTICAS

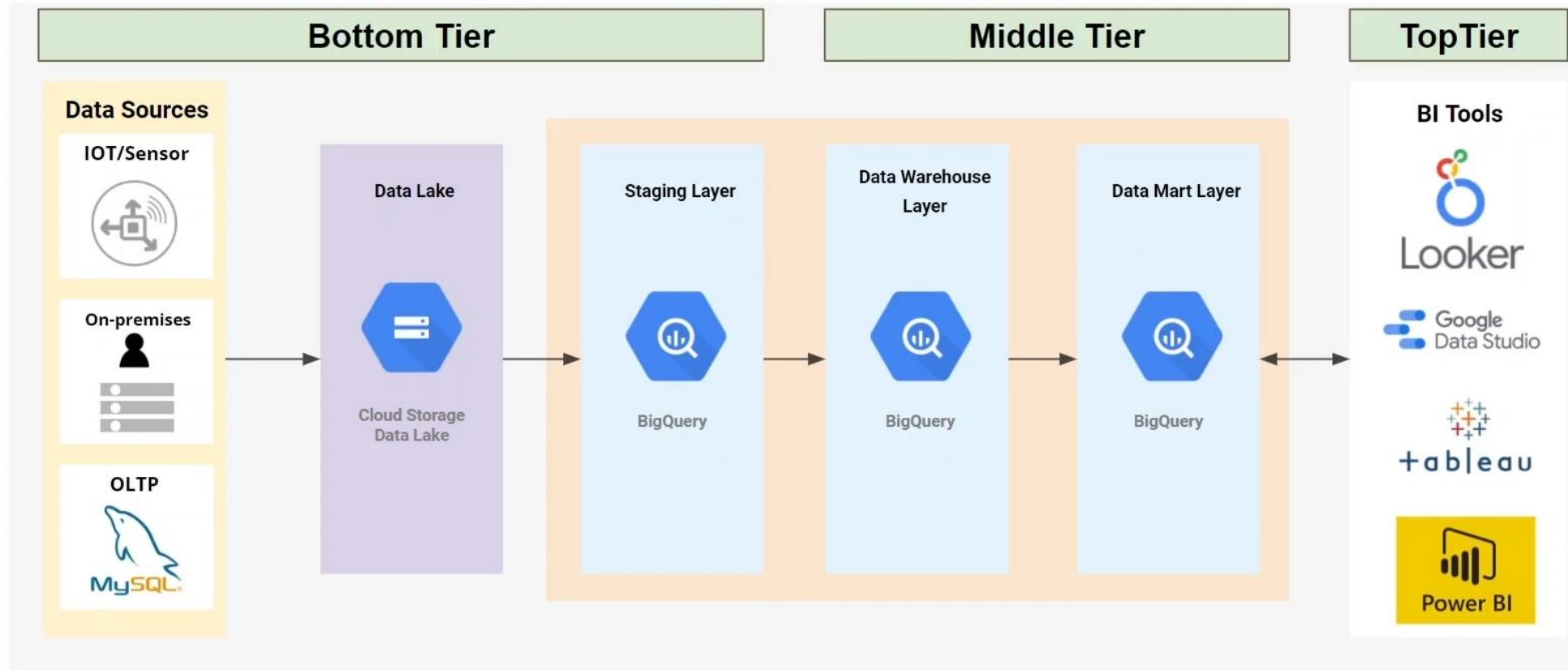
- SUBJECT ORIENTED:
SE ORIENTA A SOLUCIONAR PROBLEMAS DE NEGOCIO CON DATOS
- INTEGRADO:
RECIBE INFORMACIÓN DE DIFERENTES LUGARES CON EL FIN DE INTEGRARLOS
- VARIANTE EN EL TIEMPO
SE GUARDA INFORMACIÓN CON CIERTA FRECUENCIA (DIARIO, SEMANAL) Y SE GUARDA LA INFORMACIÓN HISTÓRICA
- NO ES VOLATIL:
 - NUEVOS DATOS ENTRAN SIN QUE VIEJOS DATOS SE BORREN

Arquitectura general de un DataWarehouse

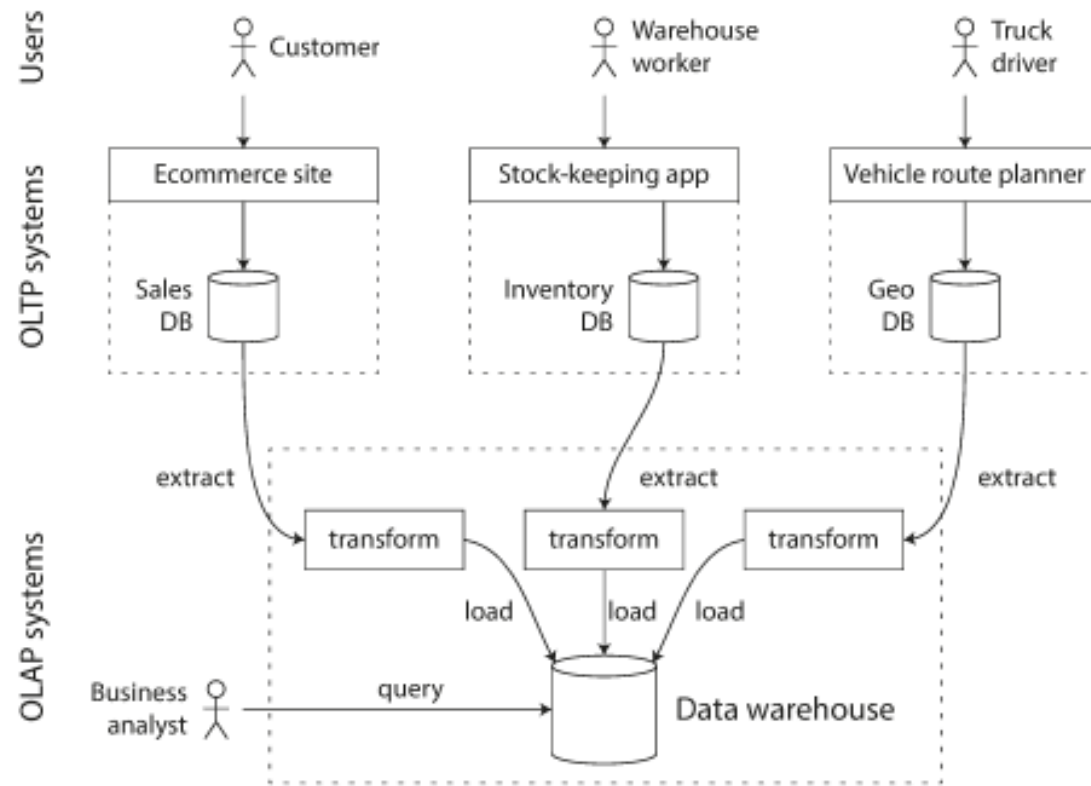


COMO SE VE UN DW?

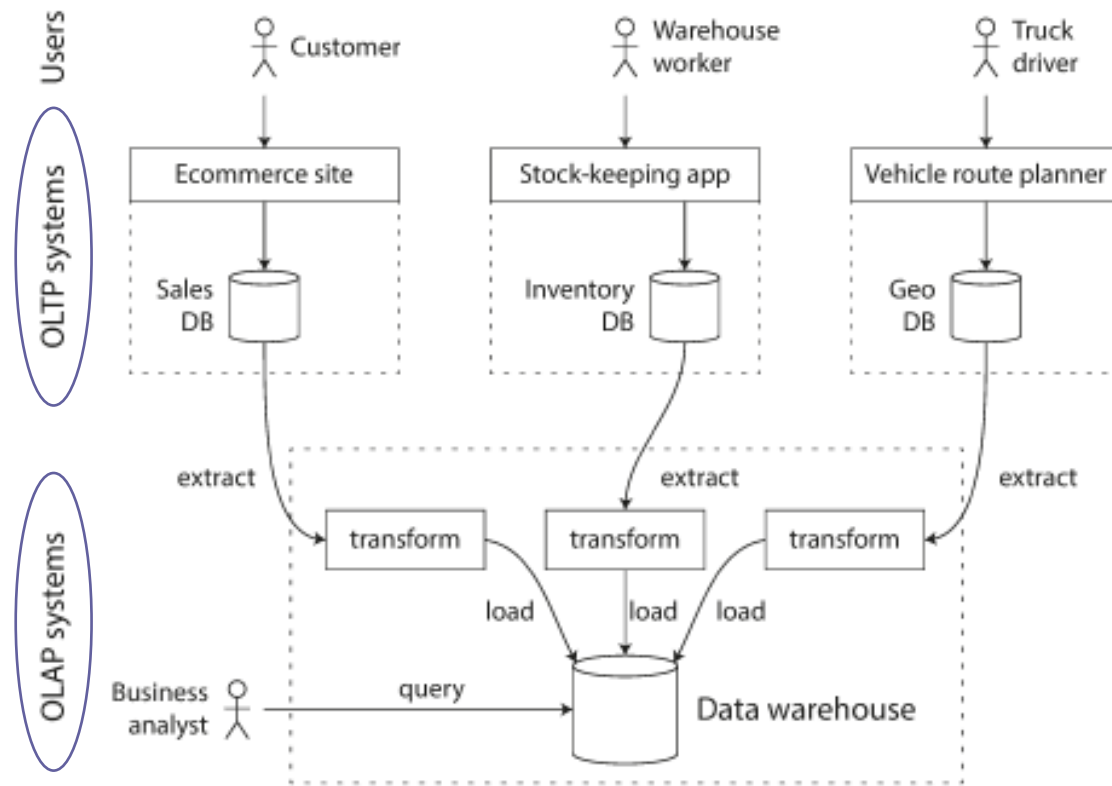
Three-tier Data Warehouse Architecture



Proceso de ETL



Proceso de ETL



Row Oriented Storage vs Column Oriented Storage

OLTP

ONLINE
TRANSACTION
PROCESSING

Data in
logical table

X	Y	Z
x1	y1	z1
x2	y2	z2
x3	y3	z3

Row-oriented storage on disk



All columns of a given row are stored together

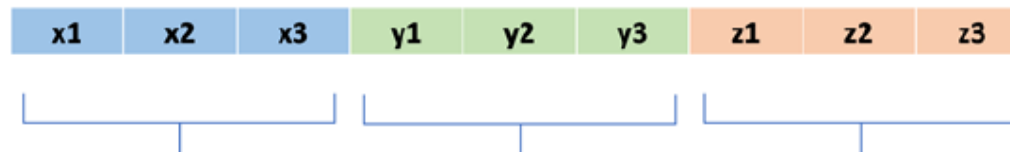
OLAP

ONLINE
ANALYTICAL
PROCESSING

Data in
logical table

X	Y	Z
x1	y1	z1
x2	y2	z2
x3	y3	z3

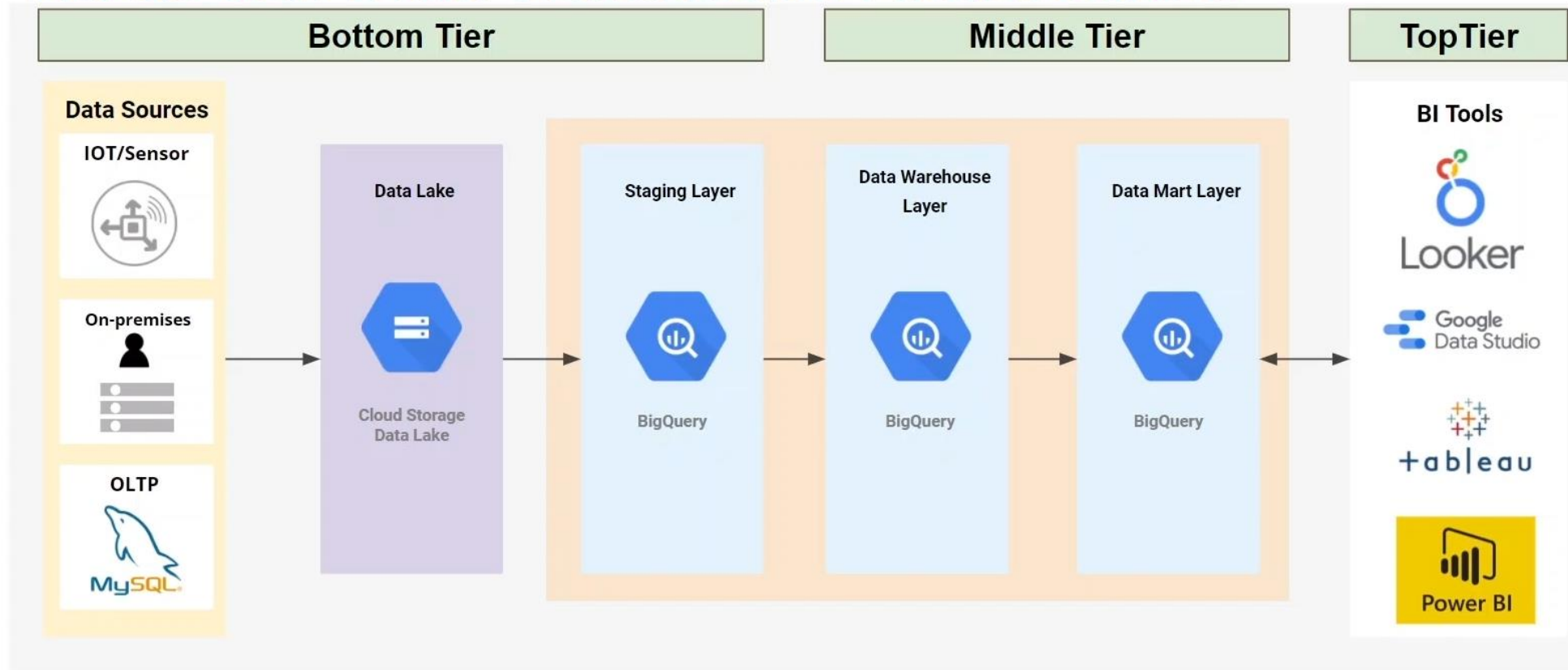
Column-oriented storage on disk



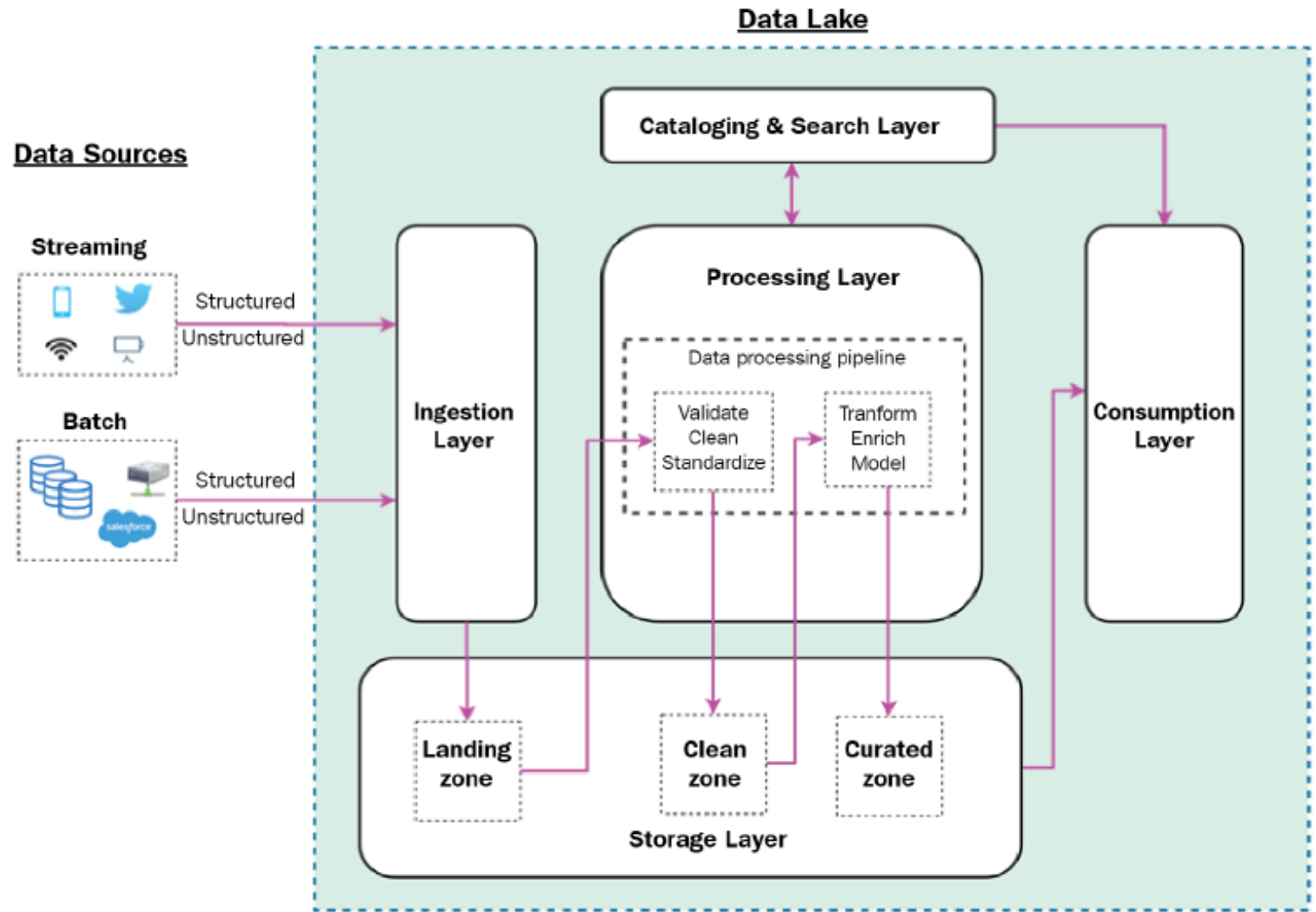
For a set of rows (a.k.a. "chunk"), all values of per column are stored together

COMO SE VE UN DW?

Three-tier Data Warehouse Architecture

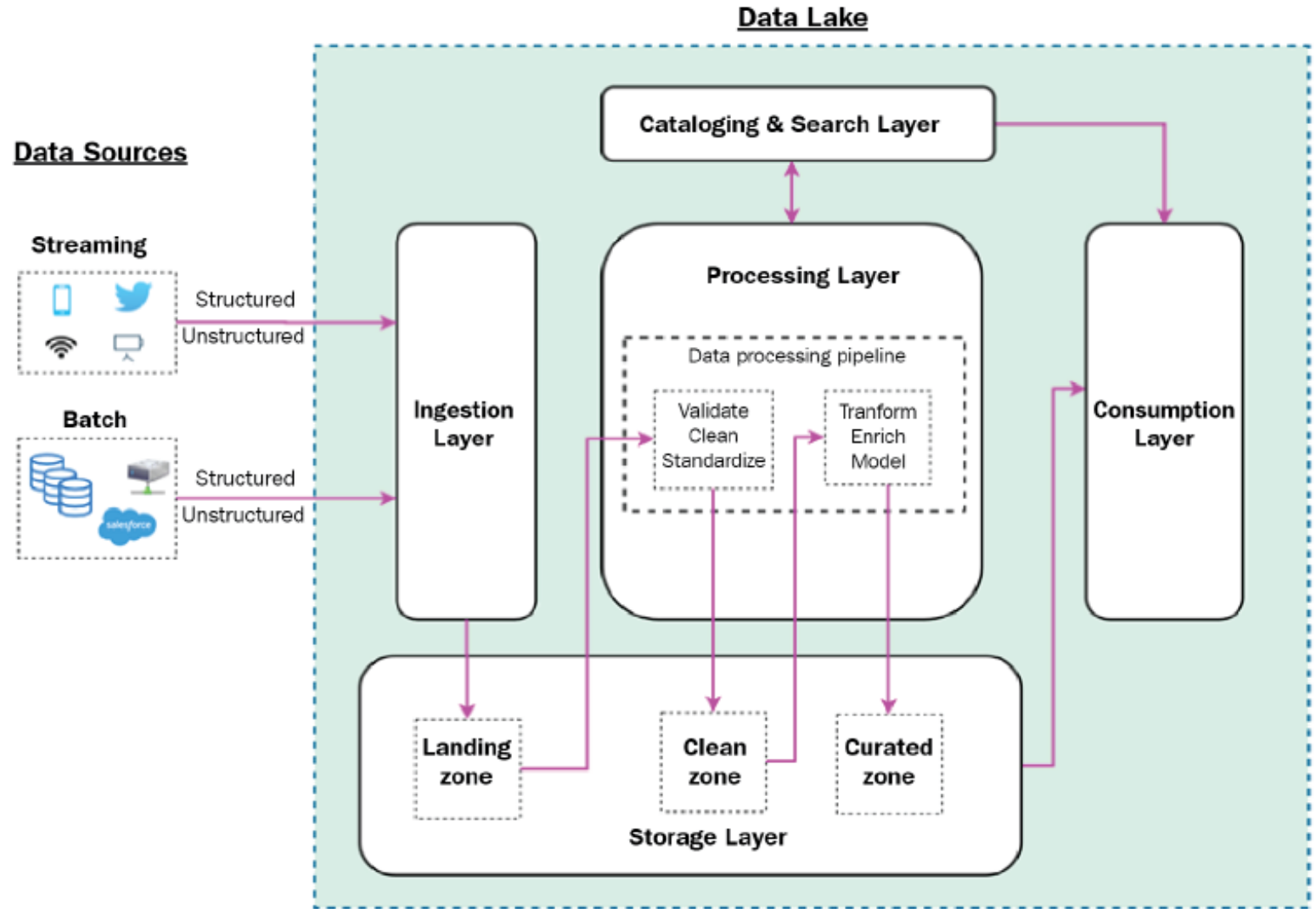


Data Lake Architecture



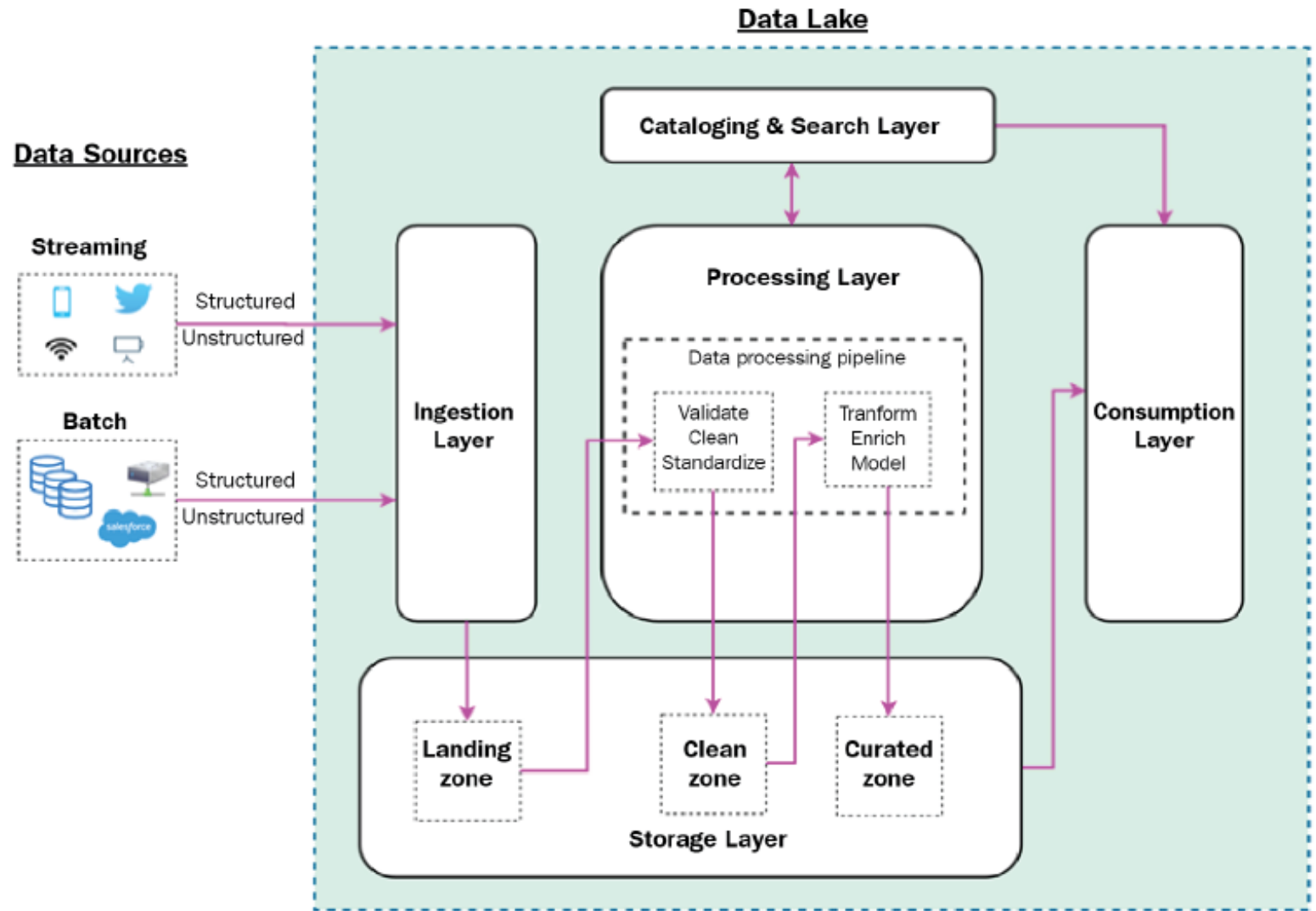
Data Lake Architecture

- LANDING ZONE:
 - Ingesta de datos
 - Se guarda “como vino”
 - Política de mantenimiento subjetiva



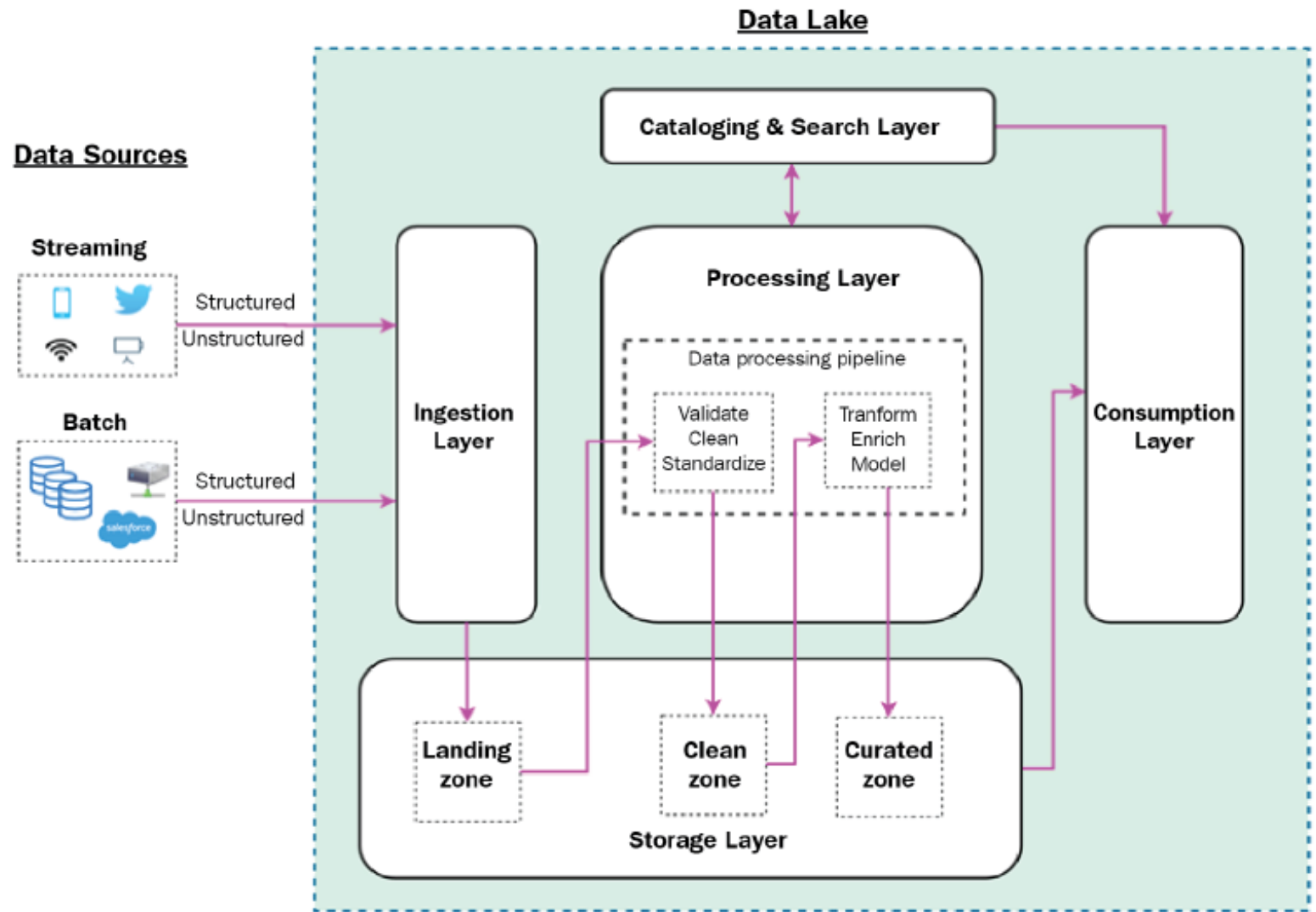
Data Lake Architecture

- CLEAN ZONE:
 - Se procesa de la landing zone
 - Se valida
 - Se limpia
 - Se optimiza el guardado (parquet)
 - Se enmascaran datos PII

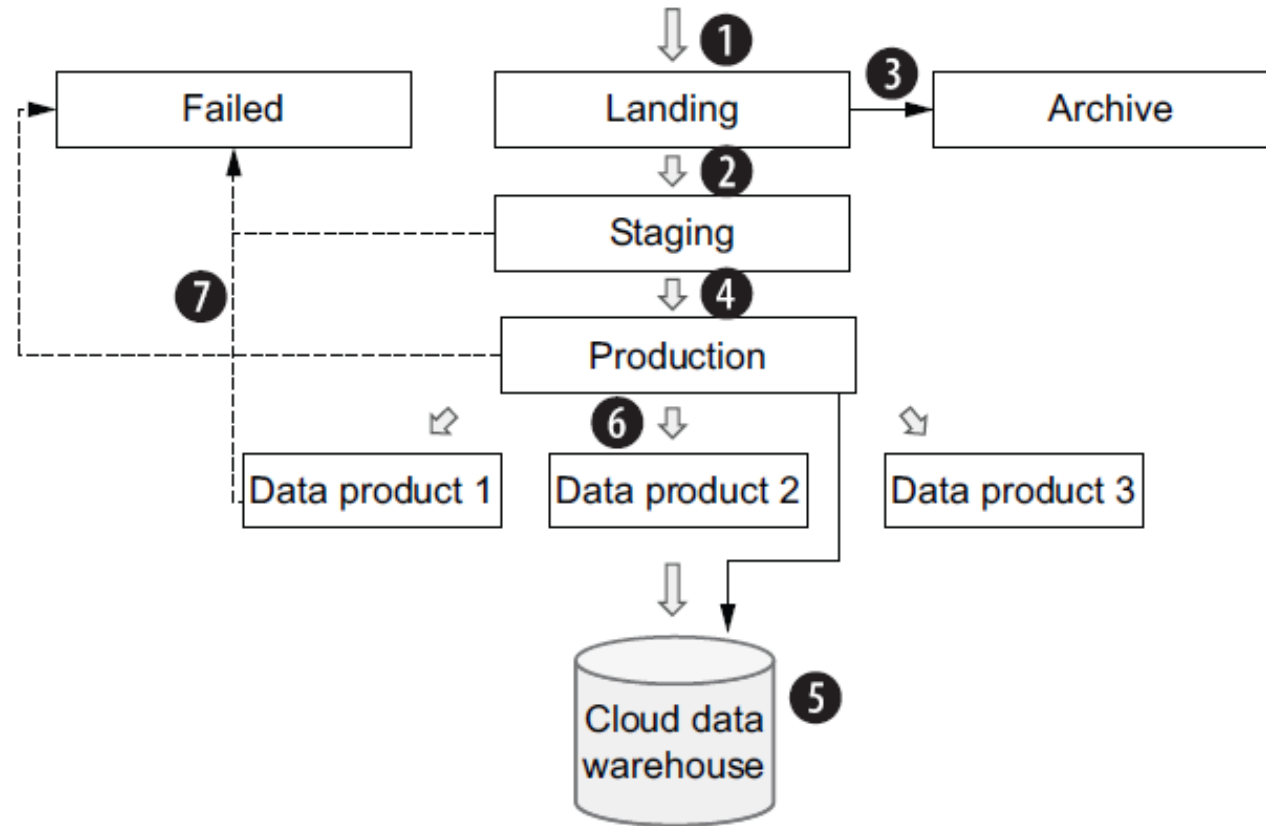


Data Lake Architecture

- CURATED ZONE:
 - Se refina con información externa/extra
 - Se agrega validación y lógica de negocio
 - Se cataloga para hacer gobierno

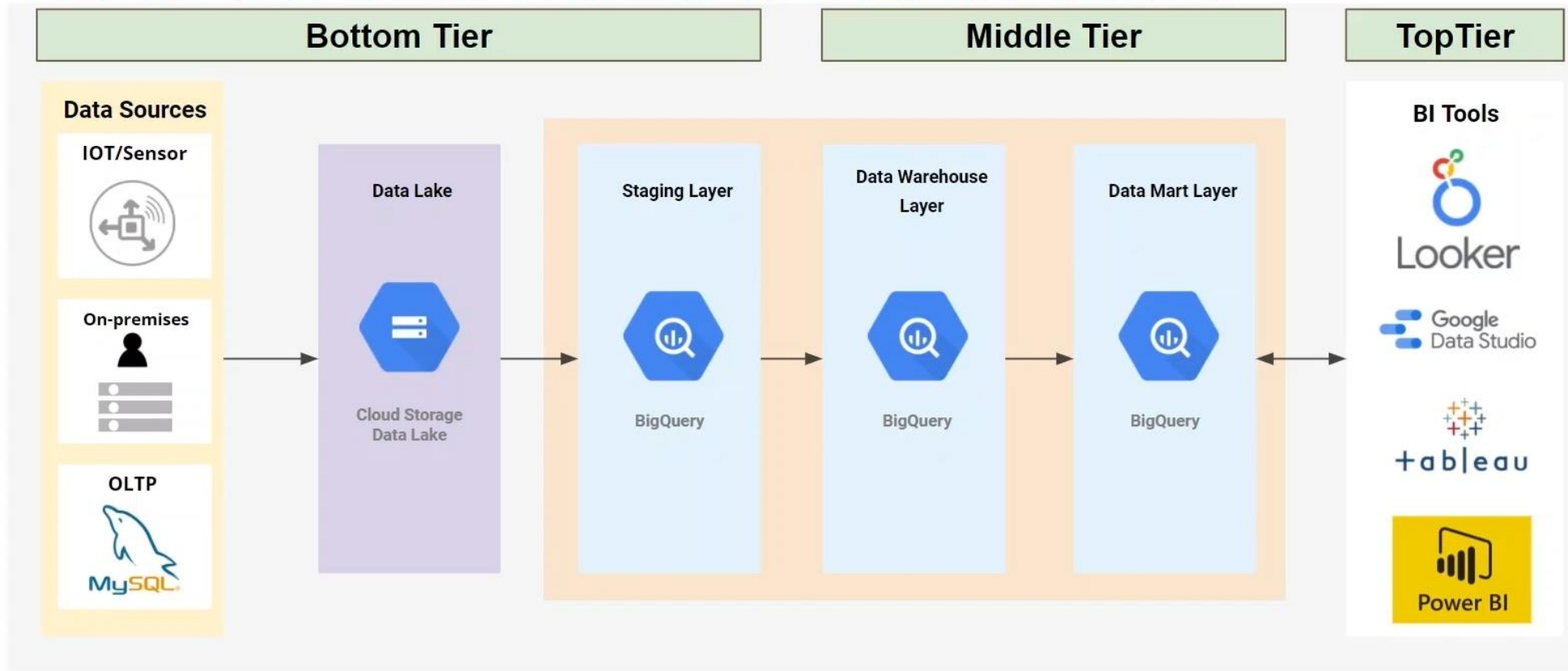


Pasos de ELT al DataLake



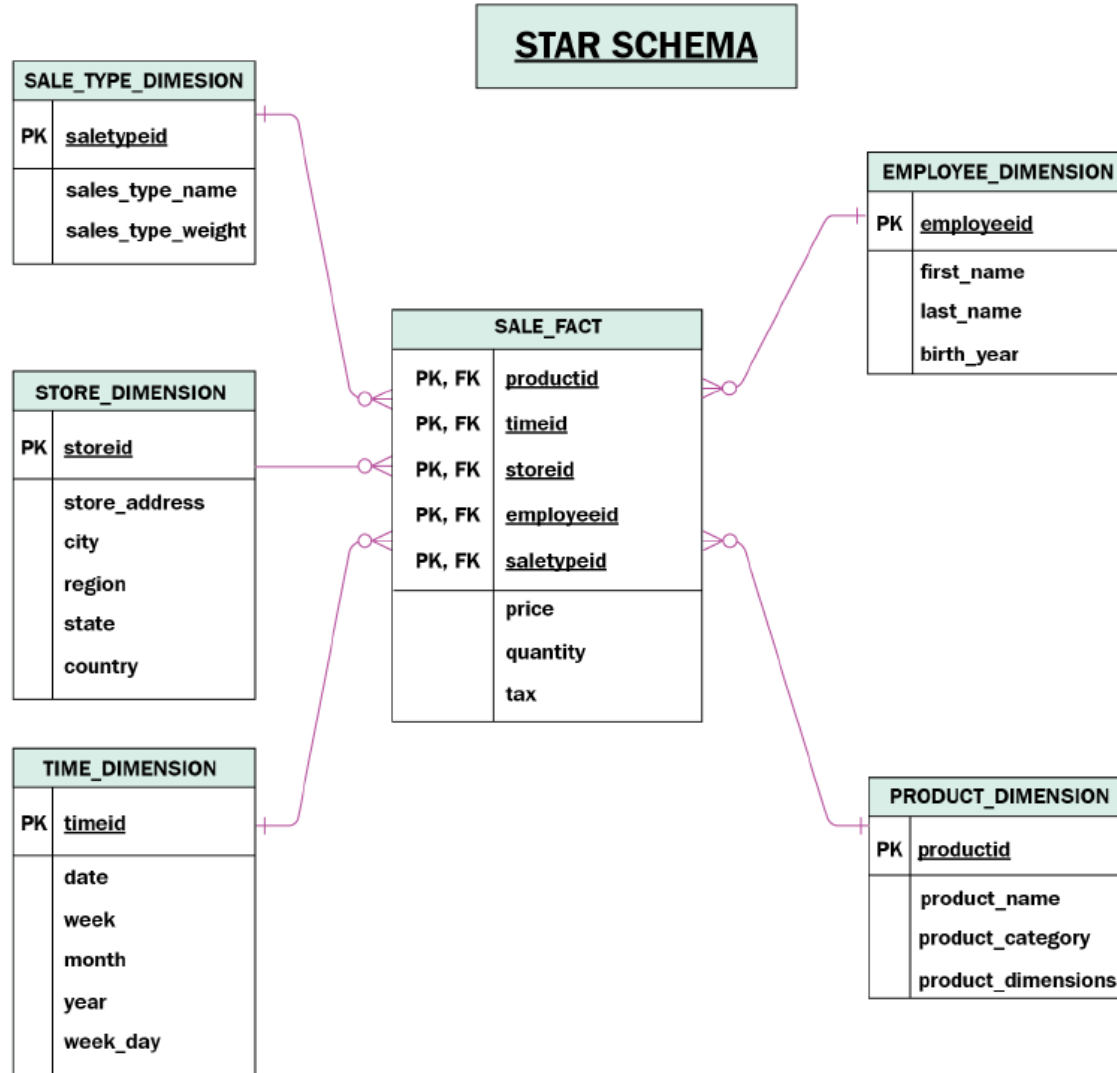
COMO SE VE UN DW?

Three-tier Data Warehouse Architecture

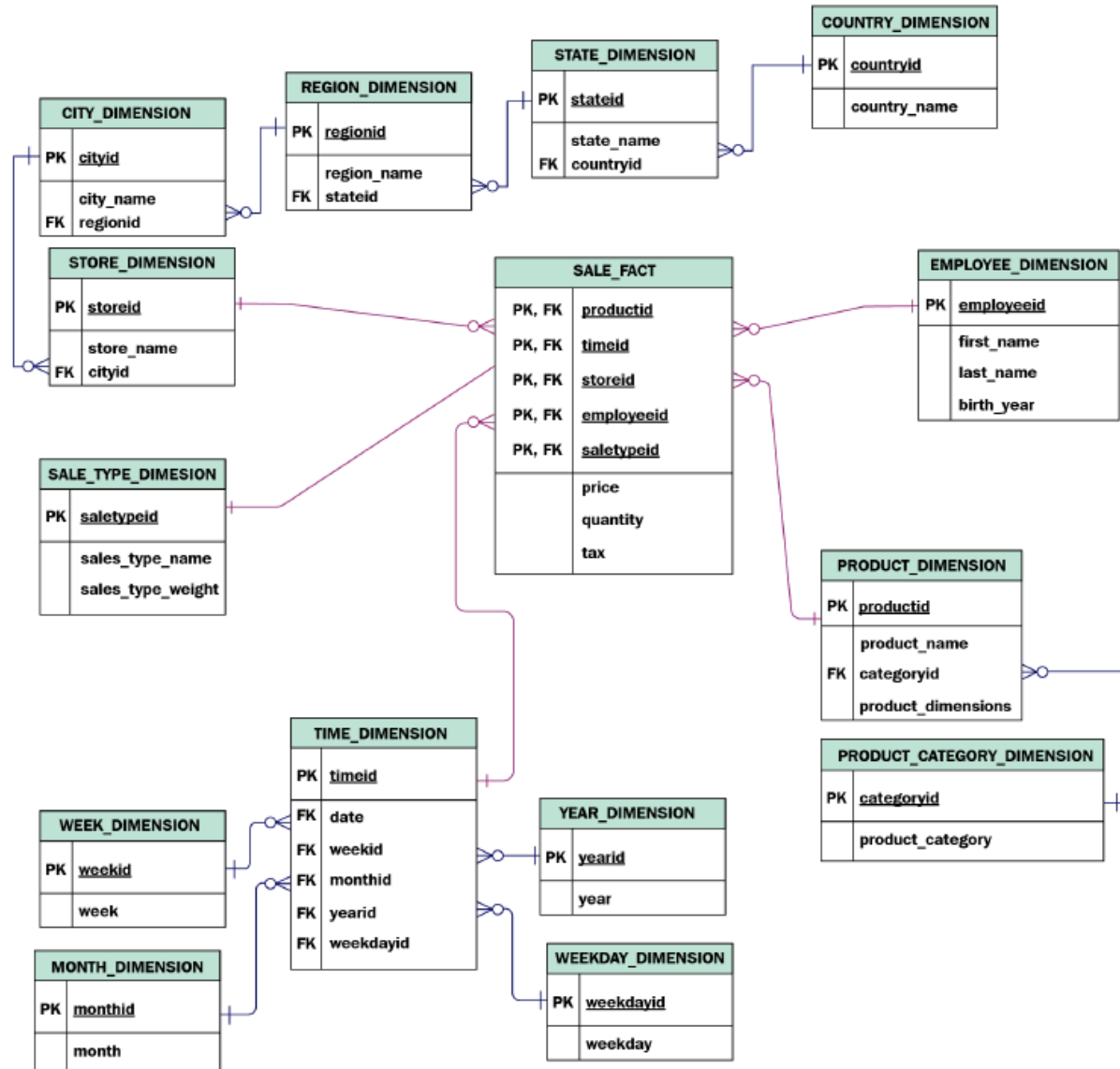


Dimensional Modelling

Star Schema

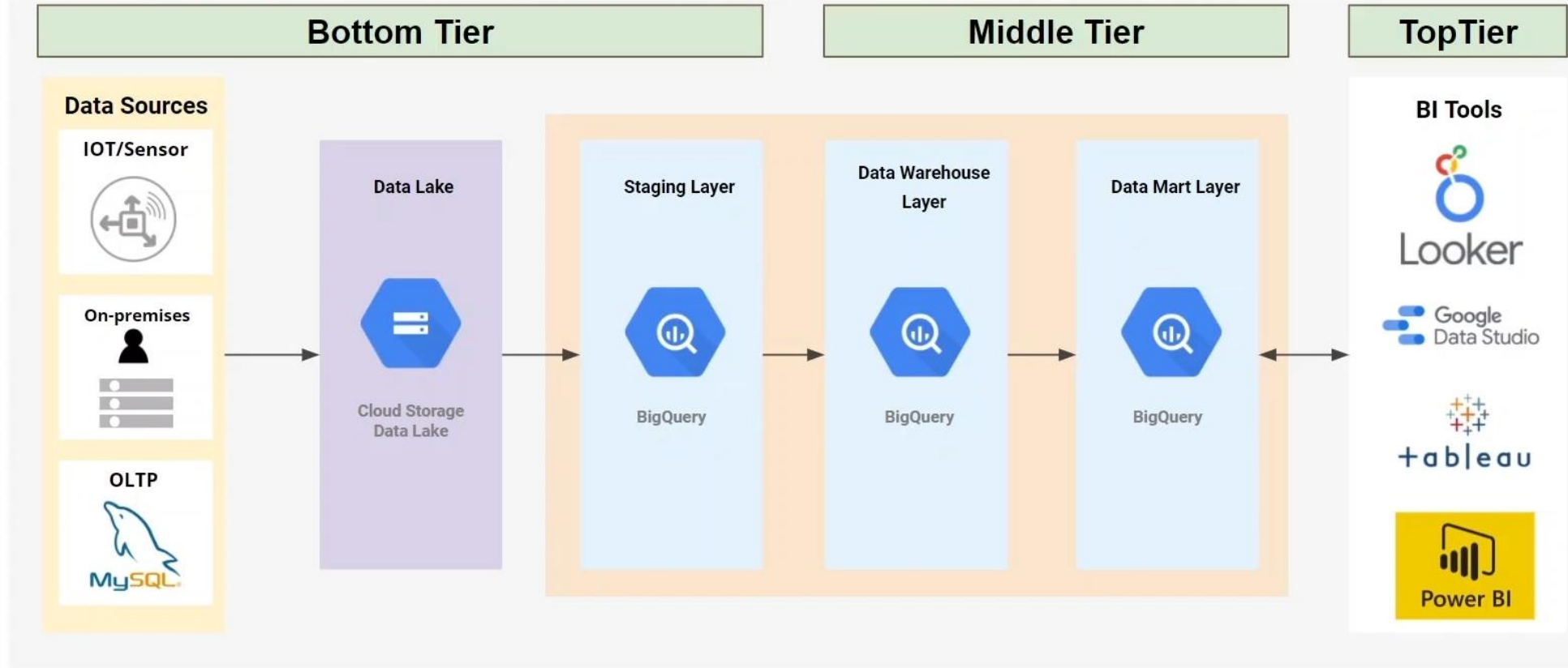


Dimensional Modelling - Snowflake Schema



DATA MART

Three-tier Data Warehouse Architecture



Data Mart Cubes

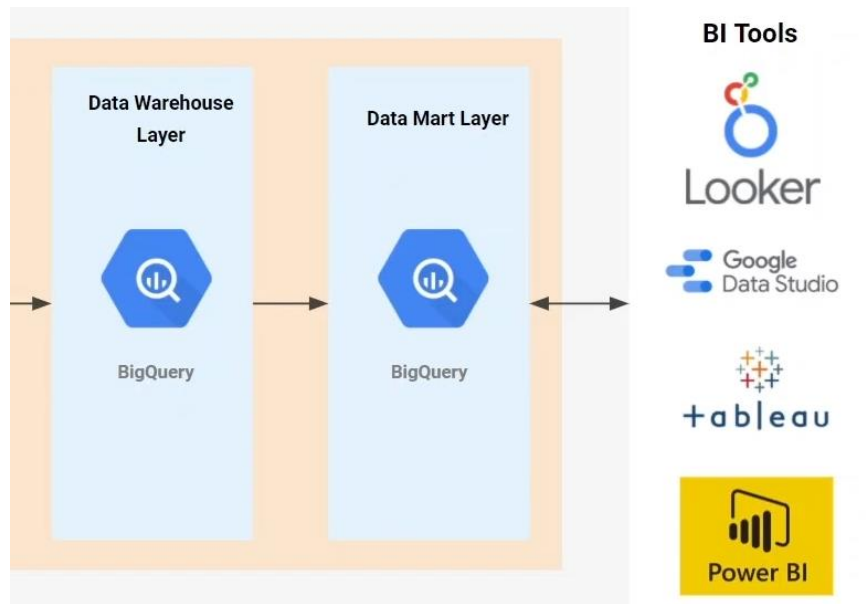


Diagram illustrating a Data Mart Cube structure with a fact table and dimensions.

Fact Table Query:

```
SELECT SUM(net_price)
FROM fact_sales
WHERE date_key = 140101
AND product_sk = 32
```

Dimensions:

- date_key** (Rows)
- product_sk** (Columns)

	32	33	34	35	total
140101	149.60	31.01	84.58	28.18	40710.53
140102	132.18	19.78	82.91	10.96	73091.28
140103	196.75	0.00	12.52	64.67	54688.10
140104	178.36	9.98	88.75	56.16	95121.09
.....
total	14967.09	5910.43	7328.85	6885.39	lots

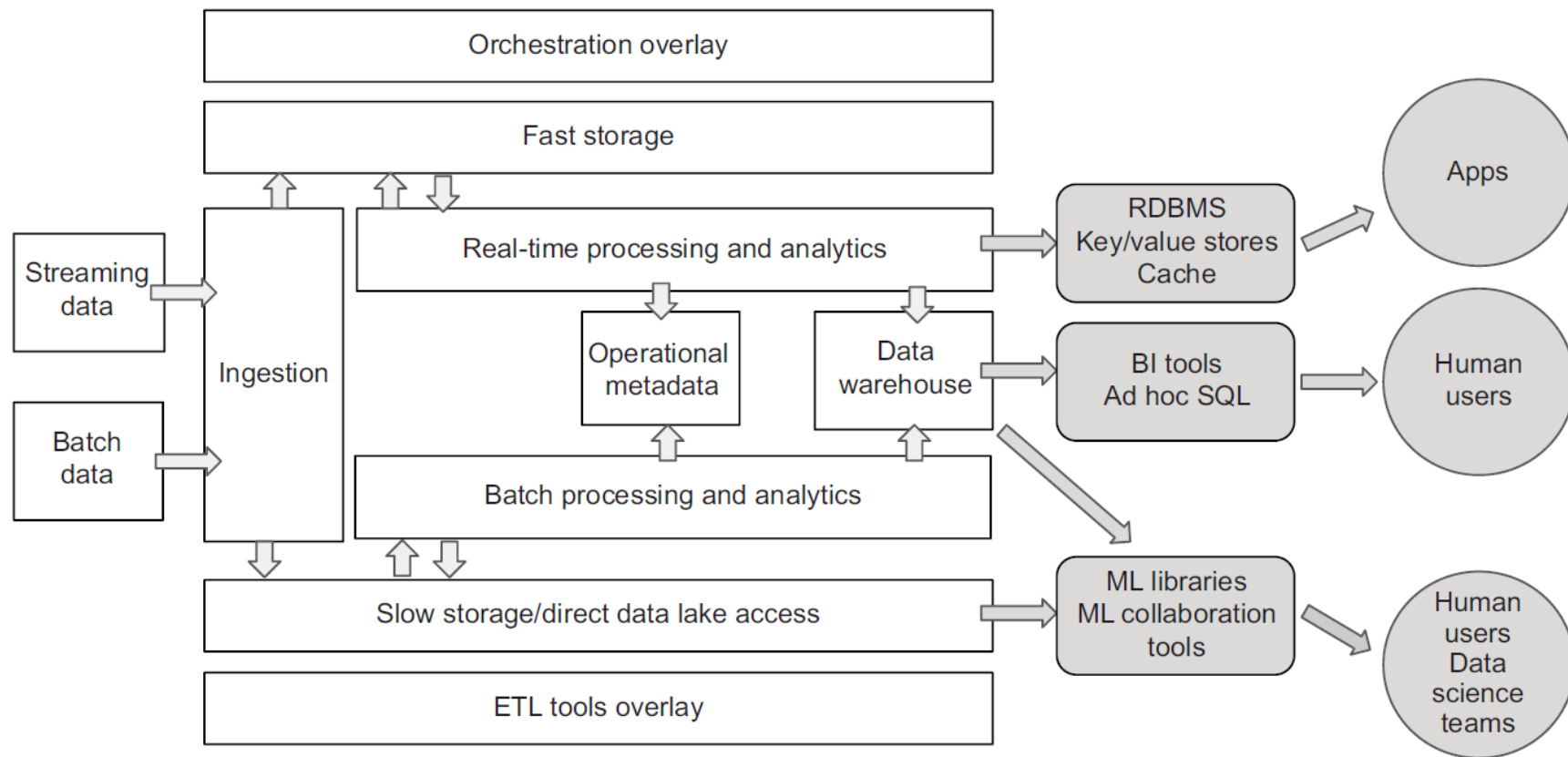
Queries for specific dimensions:

- Query for date_key = 140101:**

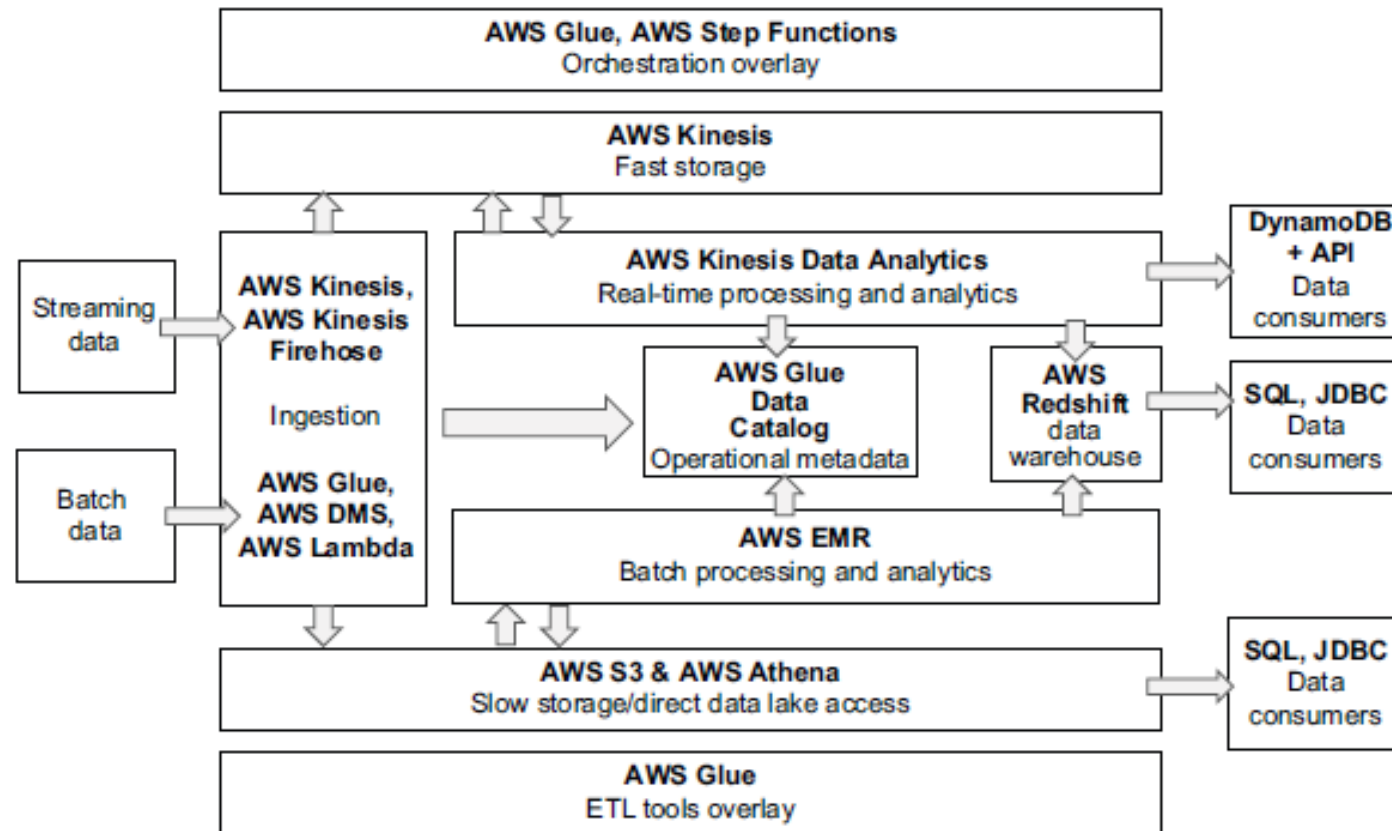
```
SELECT SUM(net_price)
FROM fact_sales
WHERE date_key = 140101
```
- Query for product_sk = 32:**

```
SELECT SUM(net_price)
FROM fact_sales
WHERE product_sk = 32
```

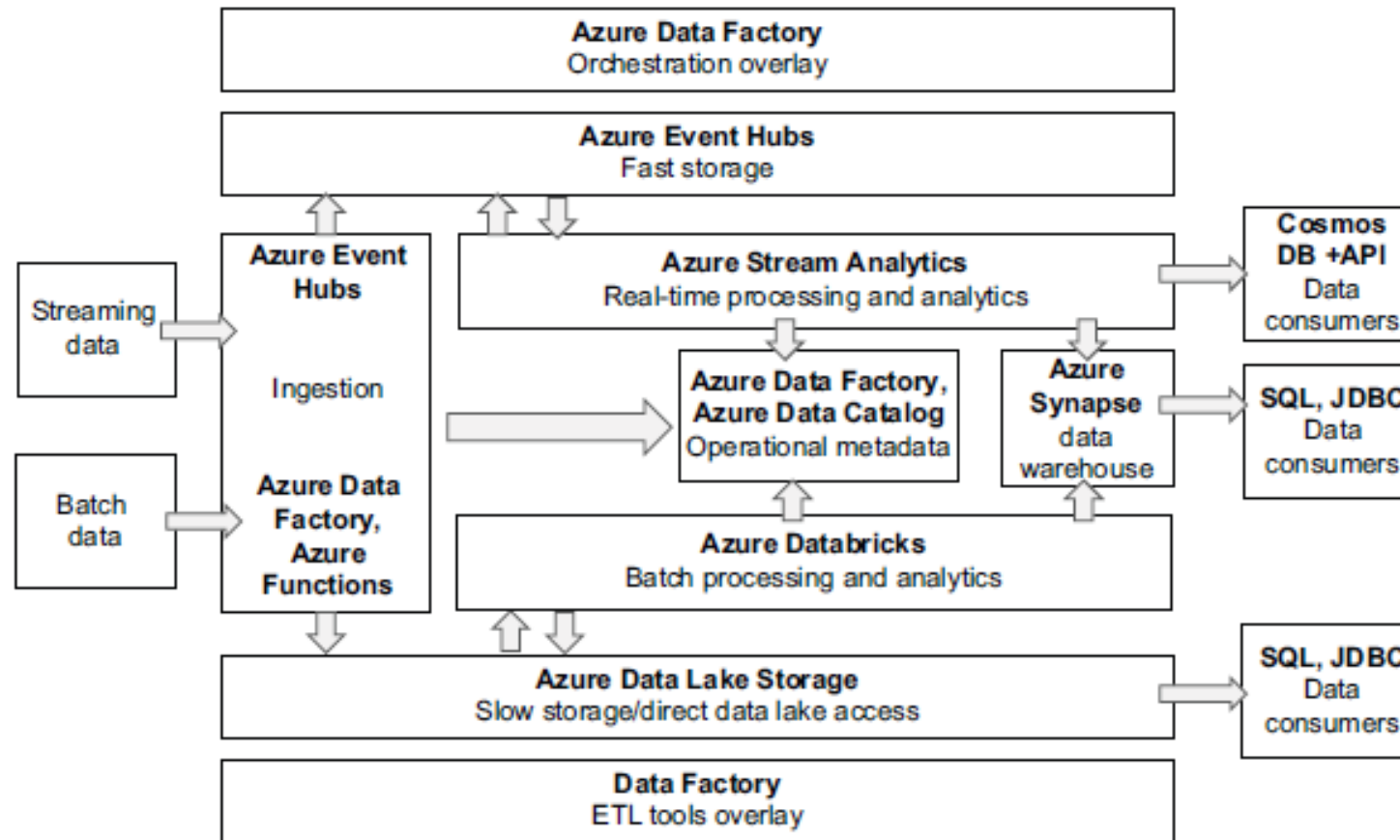
Arquitectura general de un DataWarehouse



Arquitectura general de un DataWarehouse - AWS



Arquitectura general de un DataWarehouse - AZURE



Arquitectura general de un DataWarehouse - GCP

