



# 14 DAYS

## AI CHALLENGE

### DAY 06

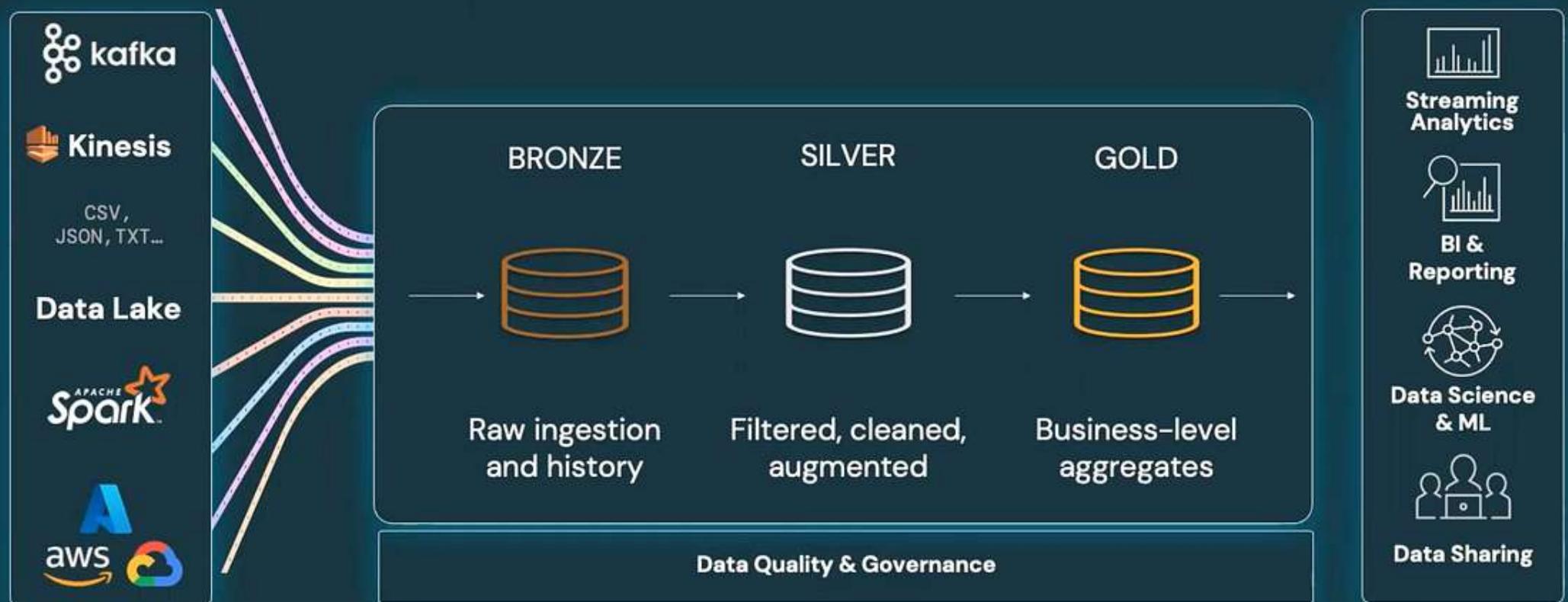
**Topic:**

Medallion Architecture

**Challenge:**

- 1.Design 3-layer architecture
- 2.Build Bronze: raw ingestion
- 3.Build Silver: cleaning & validation
- 4.Build Gold: business aggregates

# Medallion Architecture in the Lakehouse



databricks

Search data, notebooks, recents, and more..

workspace

Day\_6\_DataBrics

File Edit View Run Help Python Tabs: ON Last edit was now

Catalog

Type to search...

For you All

- My organization
  - workspace
  - system
- DeltaShares Received
- samples

## Import required libraries and functions

```
from pyspark.sql import functions as f
from delta.tables import DeltaTable
```

## Ingest Raw Data

```
raw = spark.read.csv(
    "/Volumes/workspace/e-commerce/e-commerce_data/2019-Nov.csv",
    header=True,
    inferSchema=True
)

raw.printSchema()
```

raw: pyspark.sql.connect.DataFrame = [event\_time: timestamp, event\_type: string ... 7 more fields]

```
root
|-- event_time: timestamp (nullable = true)
|-- event_type: string (nullable = true)
|-- product_id: integer (nullable = true)
|-- category_id: long (nullable = true)
|-- category_code: string (nullable = true)
|-- brand: string (nullable = true)
|-- price: double (nullable = true)
|-- user_id: integer (nullable = true)
|-- user_session: string (nullable = true)
```

Searched for "Bronze Layer - Raw Data Ingestion" in the search bar.

The notebook title is "Day\_6\_Databricks X".

The notebook content is as follows:

```
bronze = raw.withColumn("product_name", f.split(f.col("category_code"), r"\.").getItem(1) # second part after dot).filter(f.col("product_name").isNotNull())

bronze = bronze \
    .withColumn("ingestion_ts", f.current_timestamp()) \
    .withColumn("source_file", f.col("_metadata.file_path")) \
    .select(
        "event_time", "event_type", "product_id", "product_name",
        "category_id", "category_code", "brand", "price",
        "user_id", "user_session", "ingestion_ts", "source_file"
    )

bronze.write.format("delta") \
    .mode("append") \
    .save("/Volumes/workspace/e-commerce/e-commerce_data/delta/bronze/events")
```

A note at the bottom of the code cell says "See performance (?)".

The sidebar shows the catalog and workspace sections.

 databricks  
Free Edition

Search data, notebooks, recents, and more... CTRL + P

workspace ▾

Workspace ▾ Day\_6.Databricks +

Catalog File Edit View Run Help Python ▾ Tabs: ON ▾ Last edit was 3 minutes ago

Type to search...

For you All

My organization workspace system

Delta Shares Received samples

# Verify Bronze table.  
# Display first 3 rows  
display(bronze.limit(3))

# Verify schema, including ingestion\_ts & source\_file  
bronze.printSchema()

See performance (1)

Table +

rent_type	product_id	product_name	category_id	category_code	brand	price	user_id	user_session
1	1003461	smartphone	20530135556318826...	electronics.smartphone	xiaomi	489.07	520088904	4d3b30da-a5e4-49df-b
2	5000088	sewing_machine	20530135661008660...	appliances.sewing_machi...	janome	293.65	530496790	8c5f4f83-366c-4f70-86
3	3601530	kitchen	20530135638107759...	appliances.kitchen.washer	lg	712.87	518085591	3bfb58cd-7892-48cc-8

3 rows | 1.15s runtime Refreshed 12 hours ago

root

```
|-- event_time: timestamp (nullable = true)
|-- event_type: string (nullable = true)
|-- product_id: integer (nullable = true)
|-- product_name: string (nullable = true)
|-- category_id: long (nullable = true)
|-- category_code: string (nullable = true)
|-- brand: string (nullable = true)
|-- price: double (nullable = true)
|-- user_id: integer (nullable = true)
|-- user_session: string (nullable = true)
|-- ingestion_ts: timestamp (nullable = false)
|-- source_file: string (nullable = false)
```

Databricks Free Edition

Search data, notebooks, recents, and more... CTRL + P

workspace ▾

Day\_6\_Databricks

File Edit View Run Help Python Tabs: ON Last edit was 3 minutes ago

Silver Layer - Cleaning & Validation

For you All

My organization workspace system

Delta Shares Received samples

# 1. Read Bronze Delta

```
bronze = spark.read.format("delta") \
    .load("/Volumes/workspace/ecommerce/ecommerce_data/delta/bronze/events")
```

bronze: pyspark.sql.connect.DataFrame = [event\_time: timestamp, event\_type: string ... 10 more fields]

# 2. Clean and validate

```
silver = bronze \
    .filter(f.col("price") > 0) \
    .filter(f.col("price") < 10000) \
    .dropDuplicates(["user_session", "event_time"]) \
    .withColumn("event_date", f.to_date(f.col("event_time"))) \
    .withColumn("price_tier", \
        f.when(f.col("price") < 10, "budget") \
        .when(f.col("price") < 50, "mid") \
        .otherwise("premium"))
```

silver: pyspark.sql.connect.DataFrame = [event\_time: timestamp, event\_type: string ... 12 more fields]

# 3. Write to Delta Silver

```
silver.write.format("delta") \
    .mode("overwrite") \
    .save("/Volumes/workspace/ecommerce/ecommerce_data/delta/silver/events")
```

See performance (1)

Databricks Free Edition

Search data, notebooks, recents, and more... CTRL + P

workspace ▾

Day\_6\_DataBrics X +

File Edit View Run Help Python Tabs: ON Last edit was 4 minutes ago

Catalog

Type to search... 🔍

For you All

- My organization
- workspace
- system
- Delta Shares Received
- samples

# 4. Verify Silver data  
display(silver.limit(3))  
silver.printSchema()

See performance (1)

Table +

	event_time	event_type	product_id	product_name	category_id	category_code	brand	price
1	2019-11-17T08:43:01.000+00:00	view	1005105	smartphone	20530135556318826...	electronics.smartphone	apple	1363.95
2	2019-11-17T08:43:08.000+00:00	view	1480279	desktop	20530135610928667...	computers.desktop	hp	967.82
3	2019-11-17T08:43:30.000+00:00	view	28722200	shoes	20530135652284507...	apparel.shoes	respect	57.4

3 rows | 8.79s runtime Refreshed 12 hours ago

```
root
|-- event_time: timestamp (nullable = true)
|-- event_type: string (nullable = true)
|-- product_id: integer (nullable = true)
|-- product_name: string (nullable = true)
|-- category_id: long (nullable = true)
|-- category_code: string (nullable = true)
|-- brand: string (nullable = true)
|-- price: double (nullable = true)
|-- user_id: integer (nullable = true)
|-- user_session: string (nullable = true)
|-- ingestion_ts: timestamp (nullable = true)
|-- source_file: string (nullable = true)
|-- event_date: date (nullable = true)
|-- price_tier: string (nullable = false)
```

 **databricks**  
Free Edition

Search data, notebooks, recents, and more... CTRL + P

workspace ▾

**Day\_6.Databricks** +

File Edit View Run Help Python Tabs: ON Last edit was 5 minutes ago

Run all Serverless Schedule Share

**Catalog**

Type to search... 

For you All

My organization

- > workspace
- > system

Delta Shares Received

- > samples

# Gold Layer - Business Aggregates

```
# 1. Read Silver Delta
silver = spark.read.format("delta") \
    .load("/Volumes/workspace/e-commerce/e-commerce_data/delta/silver/events")
```

silver: pyspark.sql.connect.DataFrame = [event\_time: timestamp, event\_type: string ... 12 more fields]

```
# 2. Aggregate metrics per product
product_perf = silver.groupBy("product_id", "product_name") \
    .agg(
        f.countDistinct(f.when(f.col("event_type") == "view", f.col("user_id"))).alias("views"),
        f.countDistinct(f.when(f.col("event_type") == "purchase", f.col("user_id"))).alias("purchases"),
        f.sum(f.when(f.col("event_type") == "purchase", f.col("price"))).alias("revenue")
    ) \
    .withColumn("conversion_rate",
                f.round(f.expr("try_divide(purchases, views) * 100"), 3)
)
```

product\_perf: pyspark.sql.connect.DataFrame = [product\_id: integer, product\_name: string ... 4 more fields]

```
# 3. Write to Delta Gold
product_perf.write.format("delta") \
    .mode("overwrite") \
    .save("/Volumes/workspace/e-commerce/e-commerce_data/delta/gold/products")
```

Python     

15  
16  
17

 databricks  
free edition

Search data, notebooks, recents, and more... CTRL + P workspace ▾ J

Workspace Catalog + Day\_6\_DataBrics

File Edit View Run Help Python Tabs: ON Last edit was 6 minutes ago

For you All

- My organization
  - workspace
  - system
- Delta Shares Received
- samples

12 hours ago (6s) # 4. Verify Gold layer display(product\_perf.limit(5)) product\_perf.printSchema()

18 > See performance (1)

Table +

product_id	product_name	views	purchases	revenue	conversion_rate
1	kitchen	489	15	700.1400000000001	3.067
2	accessories	401	7	1283.95	1.746
3	tools	1557	15	634.2399999999999	0.963
4	tools	1155	40	8280.23	3.463
5	clocks	2141	43	2256.85	2.008

5 rows | 5.52s runtime Refreshed 12 hours ago

root

```
|-- product_id: integer (nullable = true)
|-- product_name: string (nullable = true)
|-- views: long (nullable = false)
|-- purchases: long (nullable = false)
|-- revenue: double (nullable = true)
|-- conversion_rate: double (nullable = true)
```

[Shift+Enter] to run and move to next cell [Ctrl+Shift+P] to open the command palette [Esc H] to see all keyboard shortcuts