



# 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

+ New



Databricks day 0

Databricks day 2

Databricks Day 4

Databricks Day 5

Databricks Day 6 ×



Home

Workspace

Recents

Catalog

Jobs &amp; Pipelines

Compute

Marketplace

SQL

SQL Editor

Queries

Dashboards

Genie

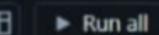
Alerts

Query History

SQL Warehouses

Data Engineering

File Edit View Run Help Python ▾ Tabs: ON ▾ Last edit was 1 hour ago



Serverless ▾

Schedule

Share

# Databricks Day 6

## Medallion Architecture

Markdown



Designed 3-layer architecture

Created a path in Delta Volume

05:42 PM (26s)

5

dbutils.fs.mkdirs("/Volumes/workspace/ecommerce/delta/bronze")

+ New

Home

Workspace

Recents

Catalog

Jobs &amp; Pipelines

Compute

Marketplace

SQL

SQL Editor

Queries

Dashboards

Genie

Alerts

Query History

SQL Warehouses

Data Engineering



Databricks day 0

Databricks day 2

Databricks Day 4

Databricks Day 5

Databricks Day 6 ×



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 1 hour ago



Run all

Serverless ▾

Schedule

Share



05:42 PM (26s)

5

0000115.ts.mquilts@volumes/workspace/ecommerce/delta/gold

&gt; lib See performance (3)

True

Loaded the data



05:44 PM (1m)

7

Python



```
1. # BRONZE: Raw ingestion
2. from pyspark.sql import functions as F
3. raw = spark.read.csv("/Volumes/workspace/ecommerce/ecommerce_data/2019-Nov.csv",
4. header=True, inferSchema=True)
5. raw.withColumn("ingestion_ts", F.current_timestamp())\
6. .write.format("delta").mode("append").save("/Volumes/workspace/ecommerce/delta/bronze/events")
```

&gt; lib See performance (1)

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

+ New

Home

Workspace

Recents

Catalog

Jobs &amp; Pipelines

Compute

Marketplace

SQL

SQL Editor

Queries

Dashboards

Genie

Alerts

Query History

SQL Warehouses



Databricks day 0

Databricks day 2

Databricks Day 4

Databricks Day 5

Databricks Day 6 ×



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 1 hour ago



Run all

Serverless ▾

Schedule

Share

## Builded Silver : cleaning &amp; validation of Data

```
▶   ✓ 05:52 PM (15s) 9
1 # SILVER: Cleaned data
2 bronze = spark.read.format("delta").load("/Volumes/workspace/ecommerce/delta/bronze/events")
3 silver = bronze.filter(F.col("price") > 0) \
4     .filter(F.col("price") < 10000) \
5     .withColumn("price_tier",
6         F.when(F.col("price") < 10, "budget")
7             .when(F.col("price") < 50, "mid")
8             .otherwise("premium"))
9 silver.write.format("delta").mode("overwrite").save("/Volumes/workspace/ecommerce/delta/silver/events")
```

&gt; See performance (1)

```
> bronze: pyspark.sql.connect.DataFrame = [event_time: timestamp, event_type: string ... 8 more fields]
> silver: pyspark.sql.connect.DataFrame = [event_time: timestamp, event_type: string ... 9 more fields]
```

+ New

Home

Workspace

Recents

Catalog

Jobs &amp; Pipelines

Compute

Marketplace

SQL

SQL Editor

Queries

Dashboards

Genie

Alerts

Query History

SQL Warehouses

Data Engineering



Databricks day 0

Databricks day 2

Databricks Day 4

Databricks Day 5

Databricks Day 6 ×



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 1 hour ago

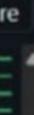


Run all

Serverless ▾

Schedule

Share



## Build Gold: business aggregates

```
▶   ✓ 06:01 PM (7s)           11: Cell 4
1 # GOLD: Aggregates
2 silver = spark.read.format("delta").load("/Volumes/workspace/e-commerce/delta/silver/events")
3 product_perf = silver.groupBy("product_id")\
4     .agg( F.countDistinct(F.when(F.col("event_type")=="view", "user_id")).alias("views"),
5           F.countDistinct(F.when(F.col("event_type")=="purchase", "user_id")).alias("purchases"),
6           F.sum(F.when(F.col("event_type")=="purchase", F.col("price").cast("double"))).alias("revenue"))
7     .withColumn(
8         "conversion_rate",
9         F.when(F.col("views") != 0, F.col("purchases")/F.col("views")*100)
10        .otherwise(None)
11    )
12 product_perf.write.format("delta").mode("overwrite").save("/Volumes/workspace/e-commerce/delta/gold/products")
```

&gt; lib See performance (1)

&gt; └─ silver: pyspark.sql.connect.DataFrame = [event\_time: timestamp, event\_type: string ... 9 more fields]

&gt; └─ product\_perf: pyspark.sql.connect.DataFrame = [product\_id: integer, views: long ... 3 more fields]