# Azure Databricks Case Study

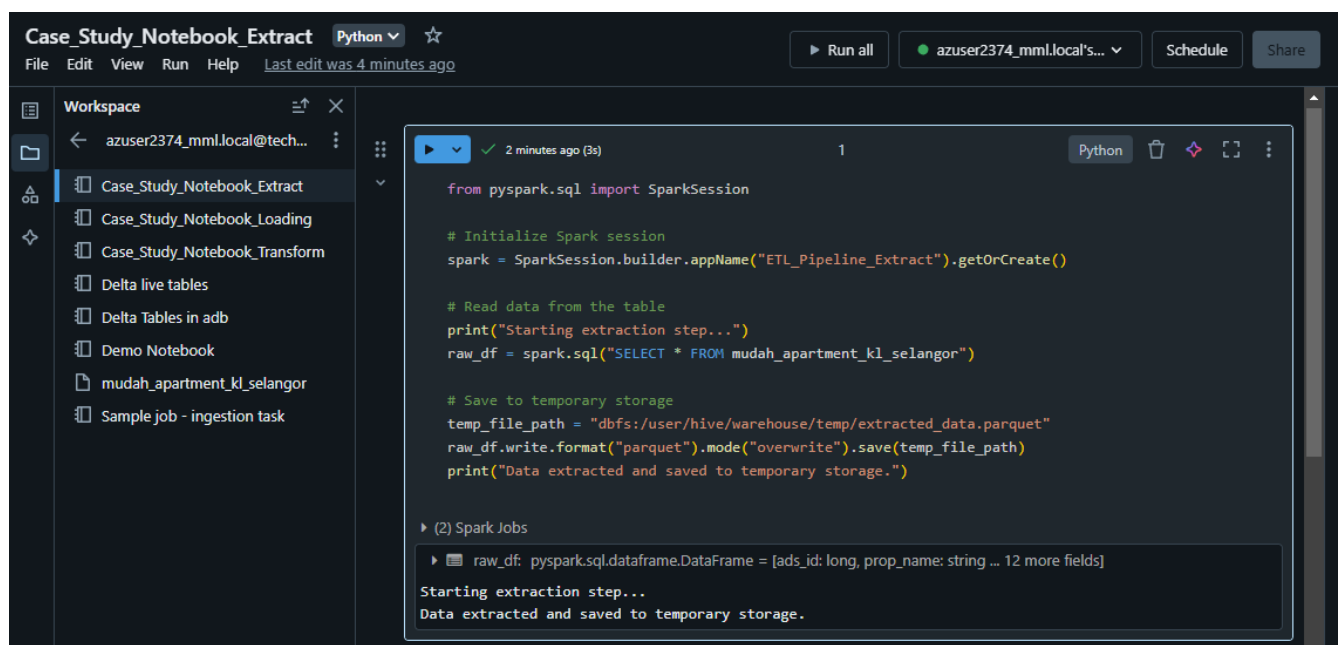**Submitted By-**
Subrat Shukla, DE-1

**Task given: Create an ETL pipeline of ingestion & transform and load queries on any data set and initiate the pipeline from workflow using the notebook.**

## Steps:

- Create a notebook with ETL queries.
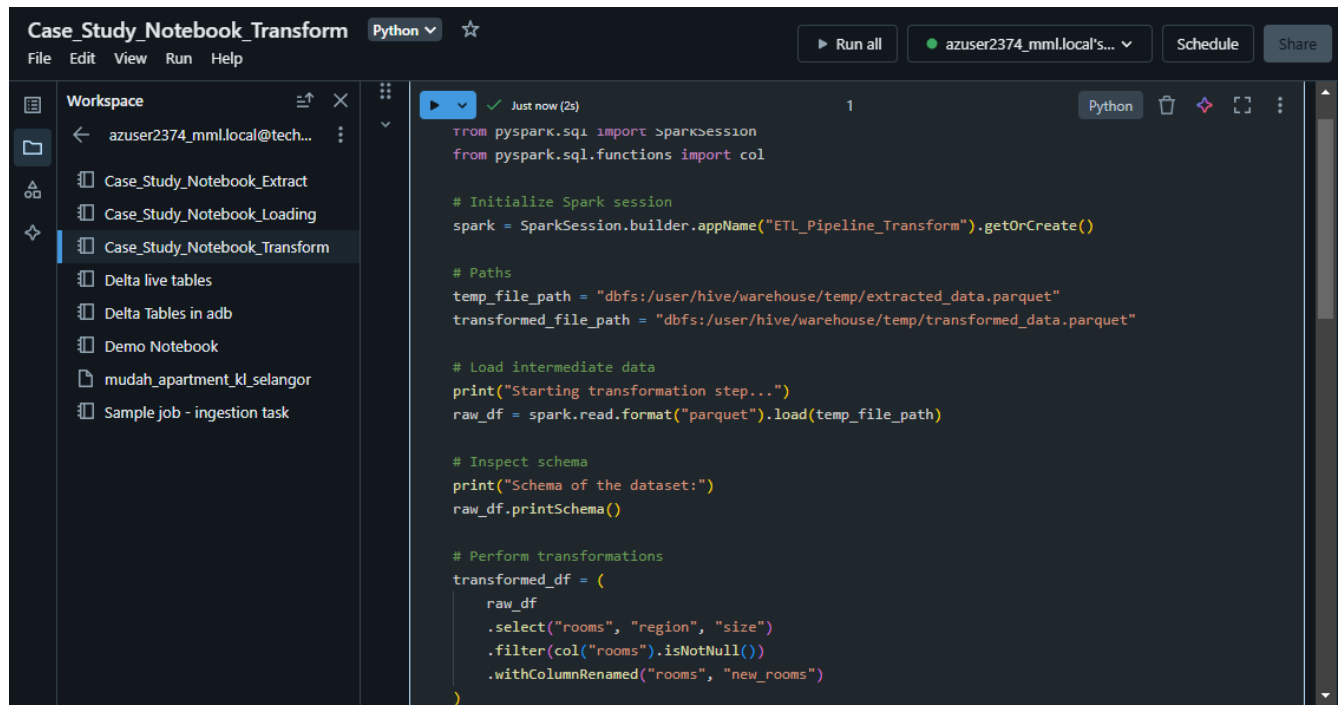- Run the notebook from workflow pipeline in azure databricks workspace.

## 1. Notebook for Ingestion (Extract):

This notebook will load the raw data into the Databricks environment, either from a file or a table.

## 2. Notebook for Transformation
This notebook will read the extracted data, apply transformations, and prepare it for loading.



```python
from pyspark.sql import SparkSession
from pyspark.sql.functions import col

# Initialize Spark session
spark = SparkSession.builder.appName("ETL_Pipeline_Transform").getOrCreate()

# Paths
temp_file_path = "dbfs:/user/hive/warehouse/temp/extracted_data.parquet"
transformed_file_path = "dbfs:/user/hive/warehouse/temp/transformed_data.parquet"

# Load intermediate data
print("Starting transformation step...")
raw_df = spark.read.format("parquet").load(temp_file_path)

# Inspect schema
print("Schema of the dataset:")
raw_df.printSchema()

# Perform transformations
transformed_df = (
    raw_df
    .select("rooms", "region", "size")
    .filter(col("rooms").isNotNull())
    .withColumnRenamed("rooms", "new_rooms")
)
```
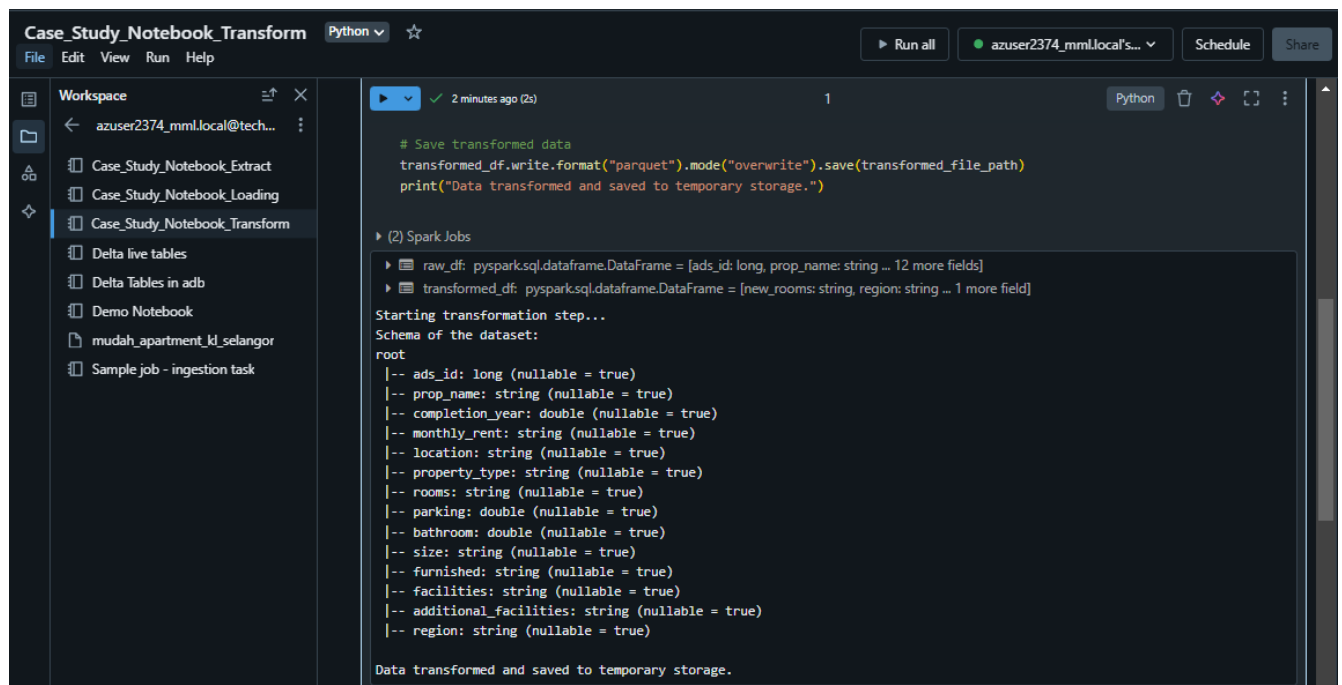


```python
# Save transformed data
transformed_df.write.format("parquet").mode("overwrite").save(transformed_file_path)
print("Data transformed and saved to temporary storage.")
```
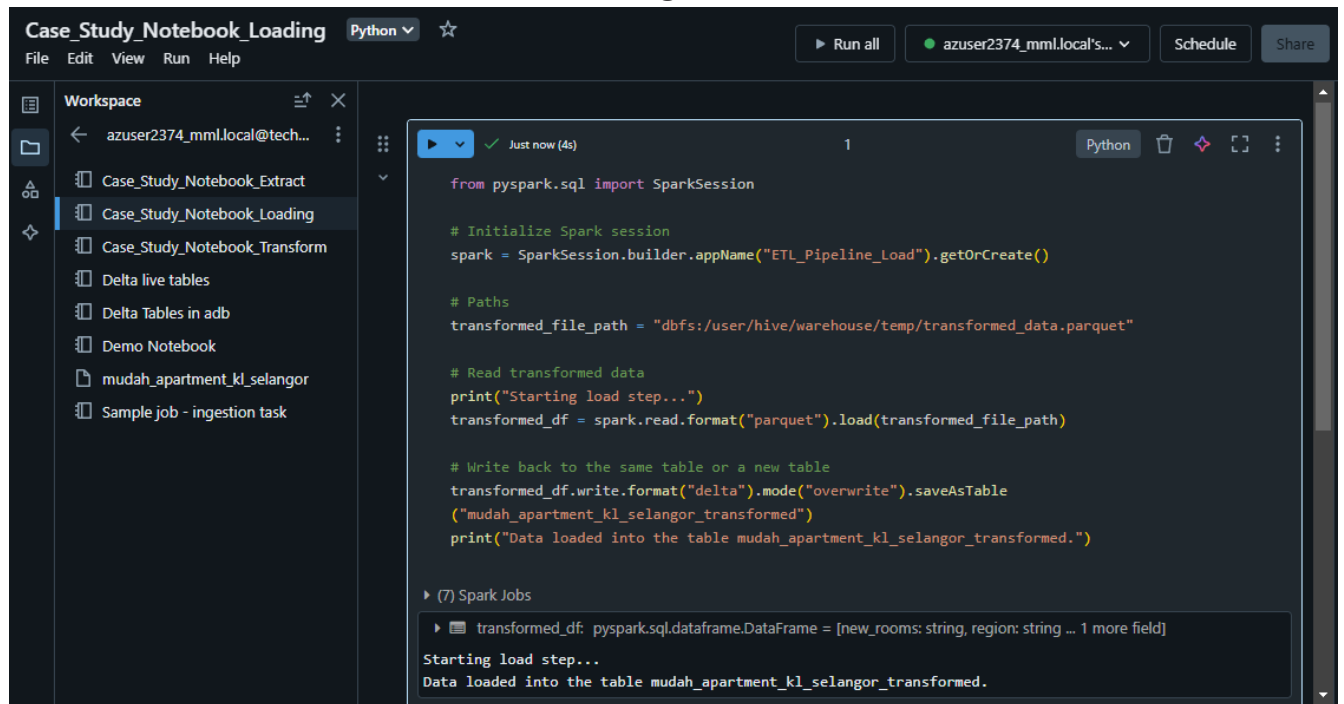
▶ (2) Spark Jobs

▶ 🔲 raw_df: pyspark.sql.dataframe.DataFrame = [ads_id: long, prop_name: string ... 12 more fields]
▶ 🔲 transformed_df: pyspark.sql.dataframe.DataFrame = [new_rooms: string, region: string ... 1 more field]

```
Starting transformation step...
Schema of the dataset:
root
 |-- ads_id: long (nullable = true)
 |-- prop_name: string (nullable = true)
 |-- completion_year: double (nullable = true)
 |-- monthly_rent: string (nullable = true)
 |-- location: string (nullable = true)
 |-- property_type: string (nullable = true)
 |-- rooms: string (nullable = true)
 |-- parking: double (nullable = true)
 |-- bathroom: double (nullable = true)
 |-- size: string (nullable = true)
 |-- furnished: string (nullable = true)
 |-- facilities: string (nullable = true)
 |-- additional_facilities: string (nullable = true)
 |-- region: string (nullable = true)

Data transformed and saved to temporary storage.
```
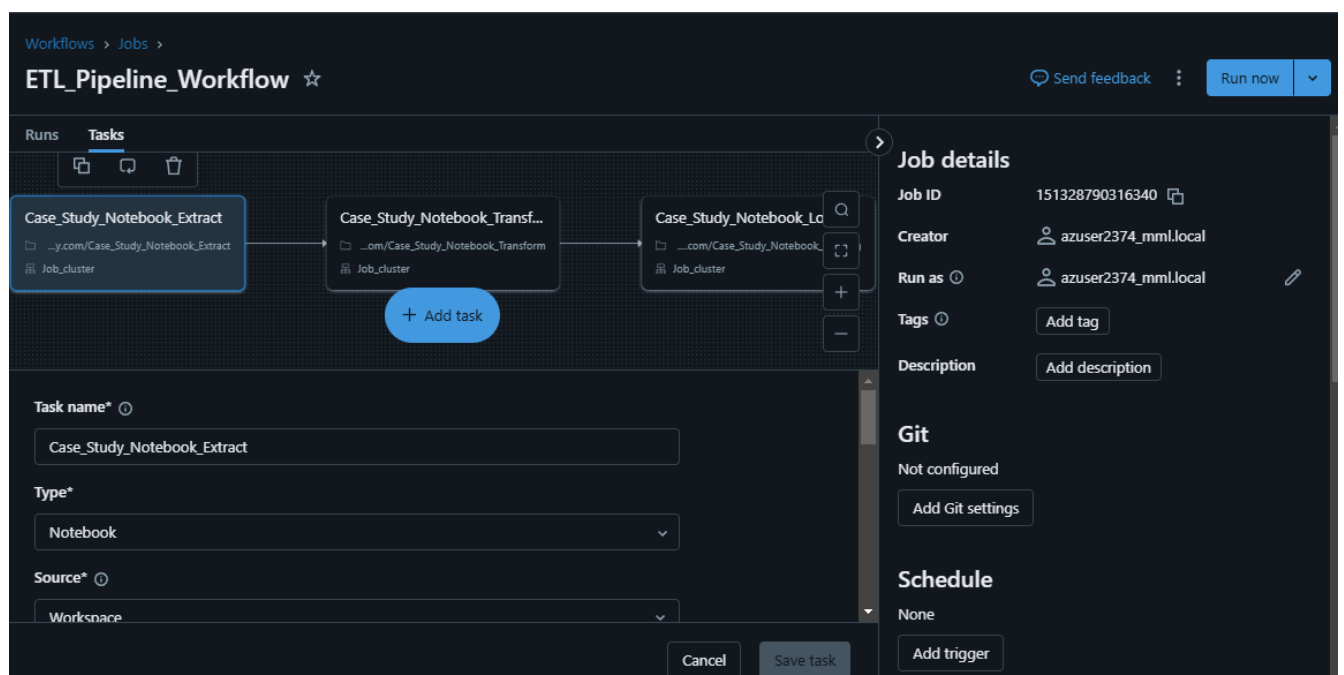
## 3. Notebook for Loading

This notebook will load the transformed data into the final destination, such as a Delta table or another storage format.



## 4. Creating Databricks Workflow:
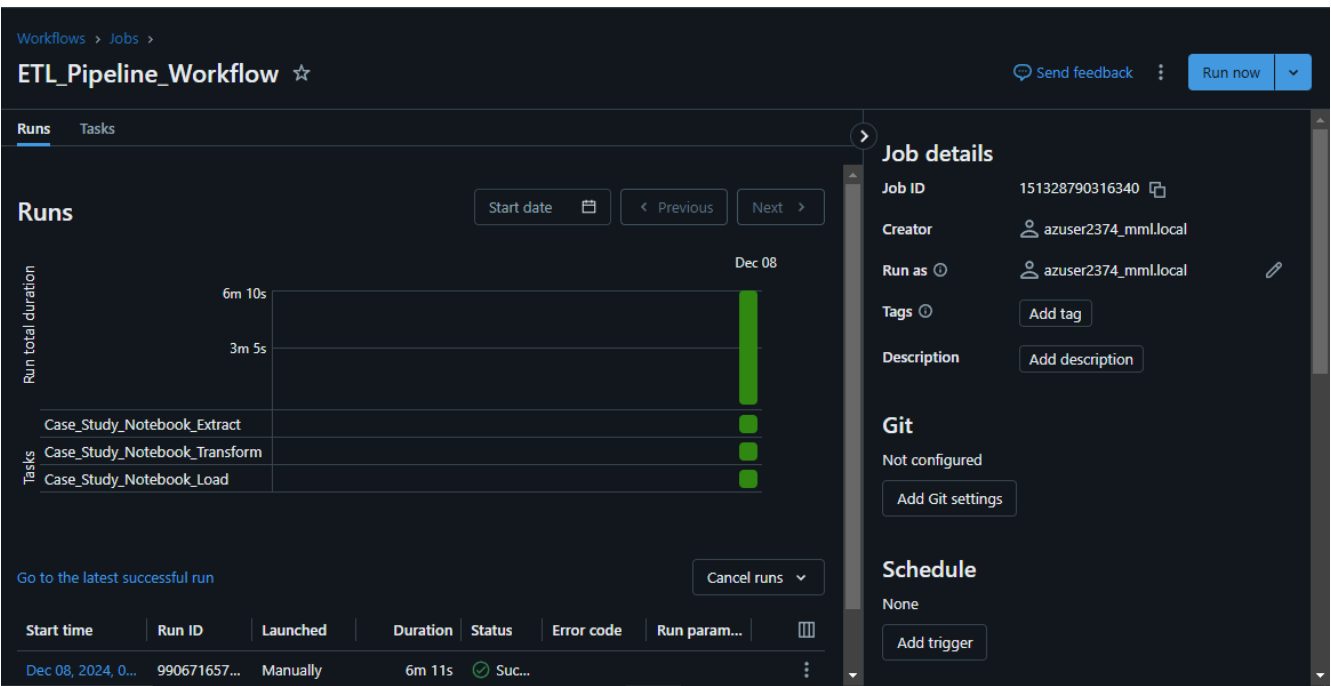
- For Extract process-

- For Transformation process-



- For Loading process-

## 5. Running the Workflows:



- Hence, The ETL Pipeline consisting processes of Extraction, Transformation, Loading has successfully created and running.

**--Thank You!**