# ASSIGNMENT - 4 : DATA WAREHOUSING LAB
Submitted by : Shivie Saksenaa

# 1. Objective

The objective of this lab was to understand how **Google Cloud Storage (GCS)** and **BigQuery** integrate to provide a complete cloud-based data storage and analytics solution.
The goal was to:

- Store data in **GCS**,
- Query it directly using an **external table** in **BigQuery**, and
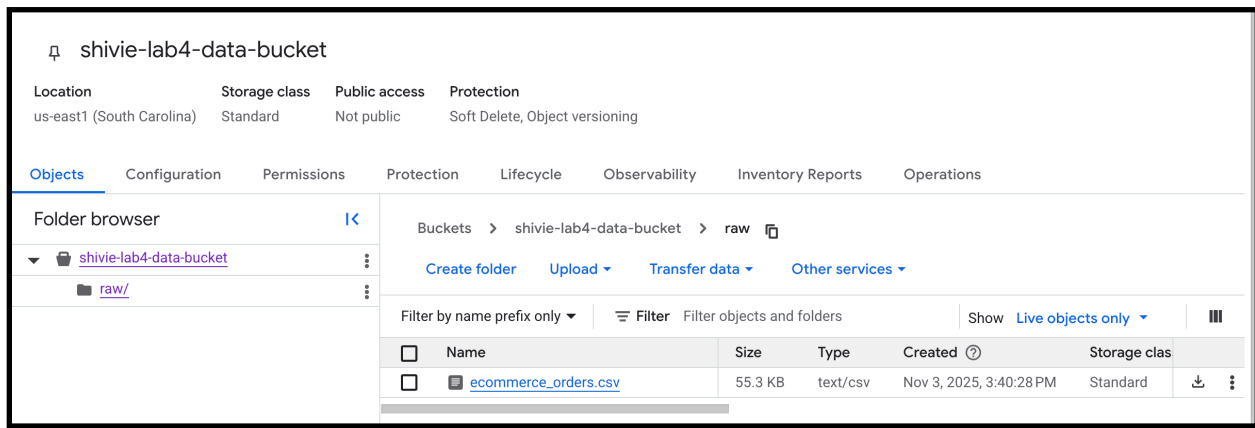- Compare performance with a **managed (partitioned) BigQuery table**.

# 2. Tools and Services Used

| Tool | Purpose |
|------|---------|
| **Google Cloud Storage (GCS)** | Store raw CSV data securely |
| **BigQuery** | Data warehousing and SQL-based analytics |
| **GCP Project ID** | clean-equinox-472523-c4 |
| **Dataset** | ecommerce_dw |
| **Region** | us-east1 |
| **Dataset used** | Ecommerce_Orders.csv |

# 3. Steps Performed

### Step 1: Upload data to Google Cloud Storage

1. Created bucket : shivie-lab4-data-bucket
2. Added folder : raw/
3. Uploaded file : ecommerce_orders.csv
4. Final GCS path:
   gs://shivie-lab4-data-bucket/raw/ecommerce_orders.csv

# Step 2: Created External Table in BigQuery

1. In BigQuery Console , I created the Table **orders.ext**

   - **Source:** Google Cloud Storage
   - **URI:** gs://shivie-lab4-data-bucket/raw/ecommerce_orders.csv
   - **File Format:** CSV
   - **Auto Detect Schema:** Enabled
   - **Destination Dataset:** ecommerce_dw
   - **Table Type:** External Table

**Result : Table was created : clean-equinox-472523-c4.ecommerce_dw.orders_ext**

# Step 3: Created Managed (Partitioned) Table

1. In BigQuery, I created the Table **orders_partitioned**

   - **Source:** Same GCS file
   - **Destination Table:** ecommerce_dw.orders
   - **Table Type:** Native (Managed)
   - **Partition By:** order_date
   - **Cluster By:** state, category

**Result: Table Created: clean-equinox-472523-c4.ecommerce_dw.orders**

**AIM: Understand the difference in performance of External Table and Native table.
Query Tested: on both the tables - External and the Native One**

```
SELECT
  DATE_TRUNC(order_date, MONTH) AS month,
  ROUND(SUM(line_amount_usd), 2) AS total_revenue
FROM `clean-equinox-472523-c4.ecommerce_dw.orders_ext`
GROUP BY month
ORDER BY month;
```

# 4. Result





## Observation and Inference – External Table (orders_ext)

- When I ran the query on the external table, BigQuery read the data directly from my GCS file (gs:// shivie-lab4-data-bucket/raw/ecommerce_orders.csv).
  The job processed 53.91 KB and took around 334 ms. Since the table only references the CSV, BigQuery had to scan the entire file every time I ran the query.
  Caching, partitioning, and clustering weren't available, and most of the time was spent reading data from GCS.

- From this, I inferred that external tables are great for quick, ad-hoc analysis of raw files but not ideal for repeated or large analytical queries. They provide flexibility without needing to load data, but performance is slower compared to a managed table

# Managed Table

## Observation and Inference – Managed Table (orders)

- When I ran the same query on the managed, partitioned table, BigQuery used its internal storage instead of reading from GCS.
  The job showed 0 bytes processed and finished instantly because the data was already stored in BigQuery and cached.
  Since this table is partitioned by order_date and clustered by state and category, BigQuery only scans the required partitions, making queries faster and more efficient.

- From this, I inferred that managed tables are much better for repeated and large-scale analytics. They support partitioning, clustering, and caching, which reduce scan cost and improve performance compared to external tables.