# Sephora E-Commerce Review Analysis

# Self Overview



- Hello my name is Nanda Muhammad .I am a final-year diploma student in Electronics Engineering at Politeknik Negeri Malang with a GPA of 3.52/4.00. I have a strong interest in the data field, particularly data engineering, data science, and data analysis. Currently, I am undertaking relevant courses and a bootcamp to acquire foundational skills, with a score of 89/100. I have developing knowledge of data modeling, data quality, and building data pipelines. I am eager to deepen my expertise and learn new things. I am adaptable, a quick learner, and committed to achieving the best results in the data industry. In the past, I worked as an Embedded Engineer Intern, where I tackled engineering problems in micro edge devices.

# **Overview Project**

E-Commerce Data Pipeline GCP

I make a e-commerce data pipeline that using the product and the customer data. The data source is from kaggle, then I make the pipeline using airflow as orchestrator with docker containerization and then for the rest I am using GCP it includes the services like cloud storage, dataproc, bigquery, and looker studio
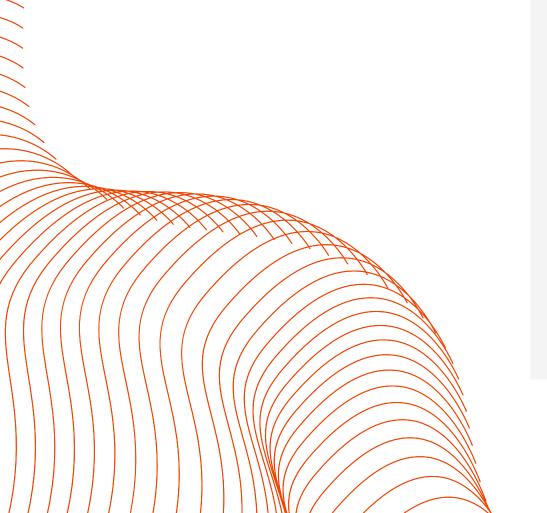
## Background

E-commerce needs to manage their data effectively, so we need to create a robust pipeline to get some insight of the data.

## Purpose

- Building reliable pipeline to support a lot of e-commerce data.
- Ensuring Data Driven Insight for E-Commerce Bussiness

# Objective
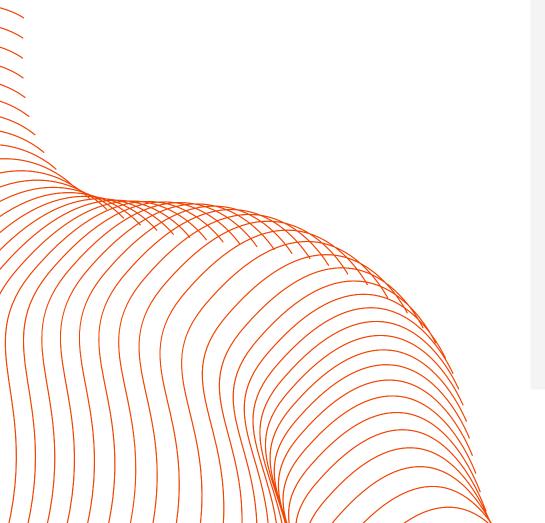
## Project Objective

- Product Performance Analysis, Identify Top Product
- Review Customer Analysis, Customer Segmentation
- Visualize the distribution

## Expected Output

- Robust Data Pipeline
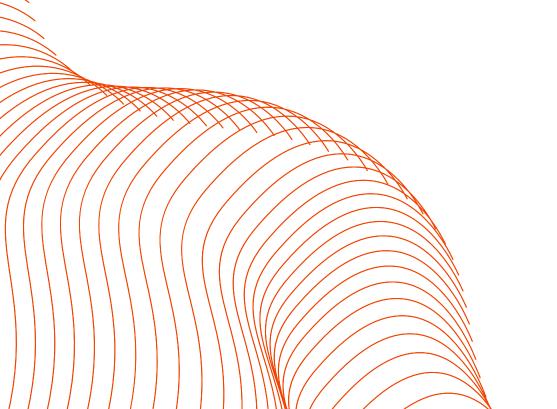- Dashboard from Product and Review analysis
- Getting some actionable insight

# Dataset

**Using 1.110.000 records  data**

- For more details use this link to the <u>kaggle</u>

# Timeline

### Data Extraction

- Extract data from kaggle using kaggle API
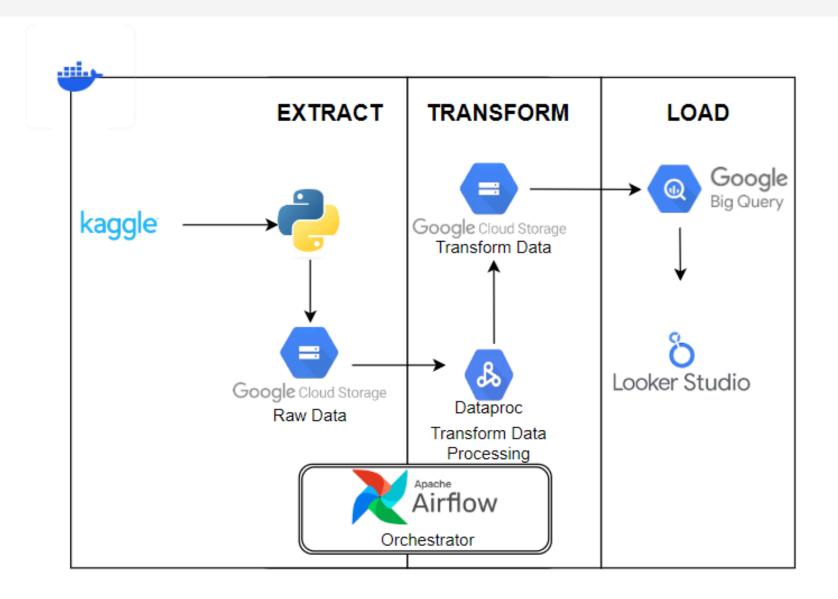
### Data Storage

- Save the data into Data Storage

### Data Preprocessing and Modeling

- Cleaning the data, null values, duplicate values
- Utilise spark to transform the data into star schema

### Data Loading

- Load the data into Data Warehouse Bigquery

### Data Visualization

- Visualize the data using Looker Studio

# Architecture



## Extract

- Extract the data using Python and kaggle API
- Save the raw data to data storage

## Transform

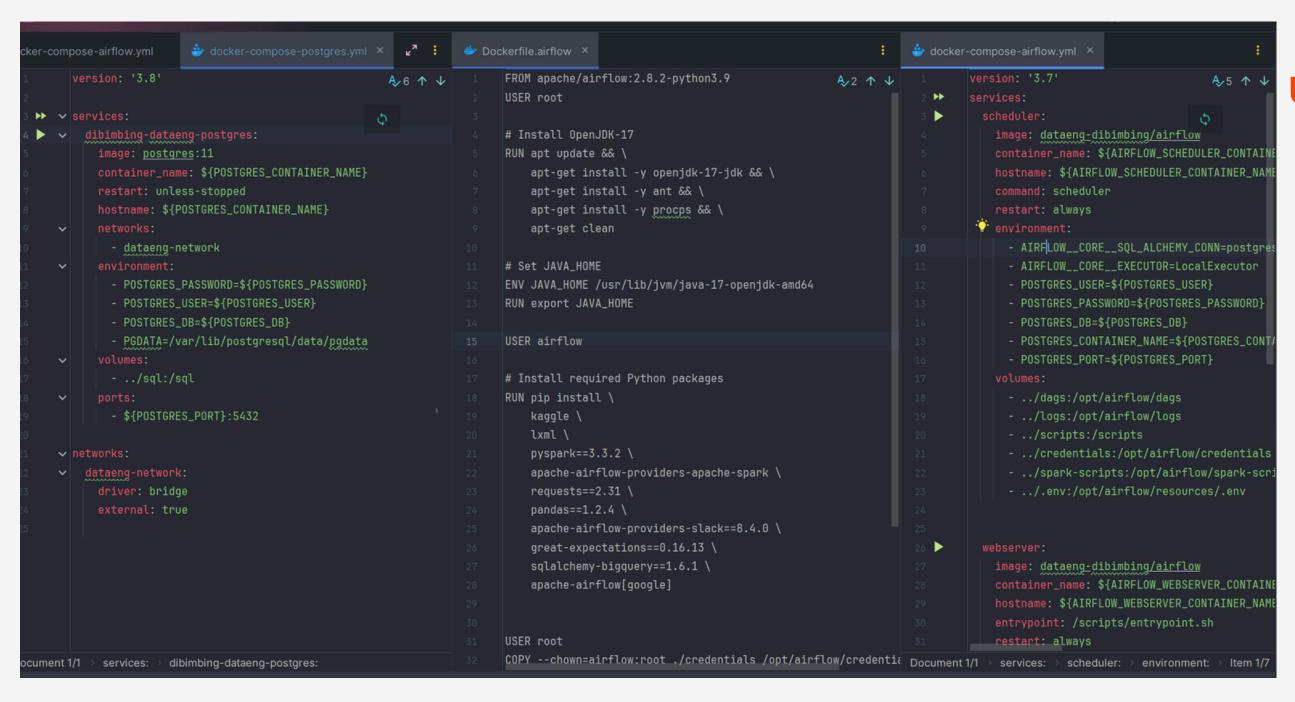- Clean the data using pyspark, duplicate, missing values, etc.

## Load

- Load the data to Data Warehouse, Bigquery
- Get visualization using looker studio

## Data Orchestration

- Using Airflow for scheduling and monitoring data pipeline

## Containerization

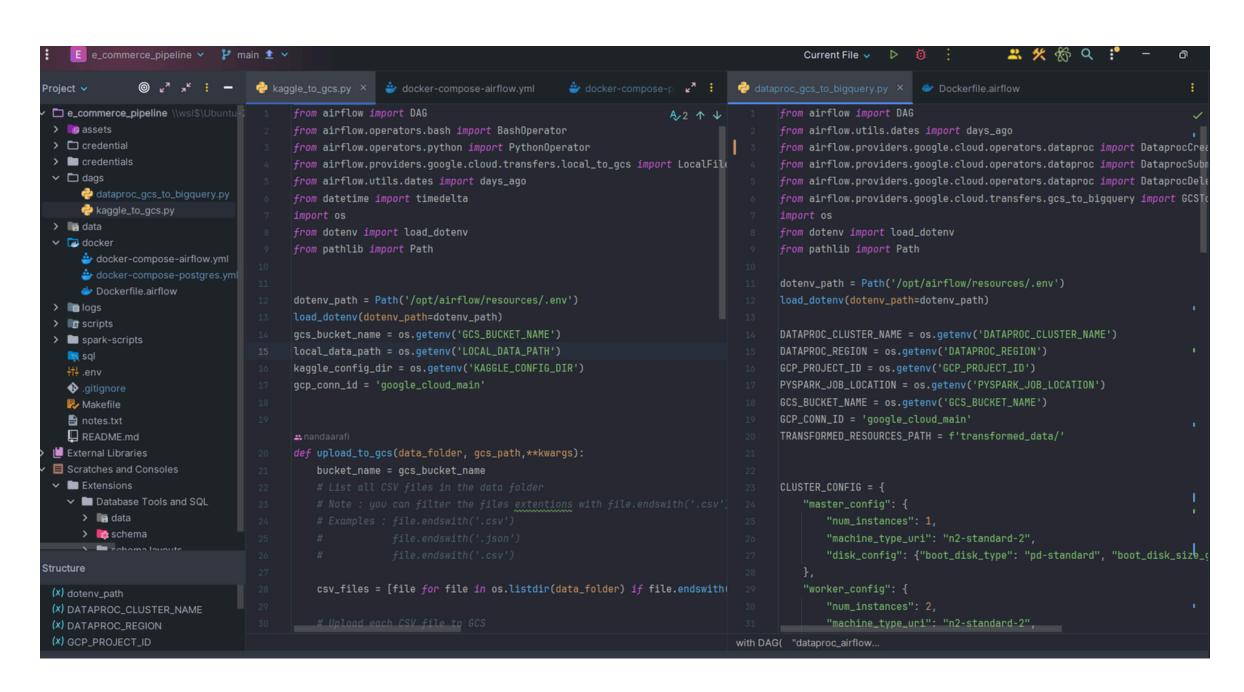- Using Docker for seamless Containerization

# VISUALIZATION

# Code Overview



## Usage

- Using Docker Airflow and Postgres for containerization
- Dockerfile airflow, for jdk installation, python installation and python package installation.
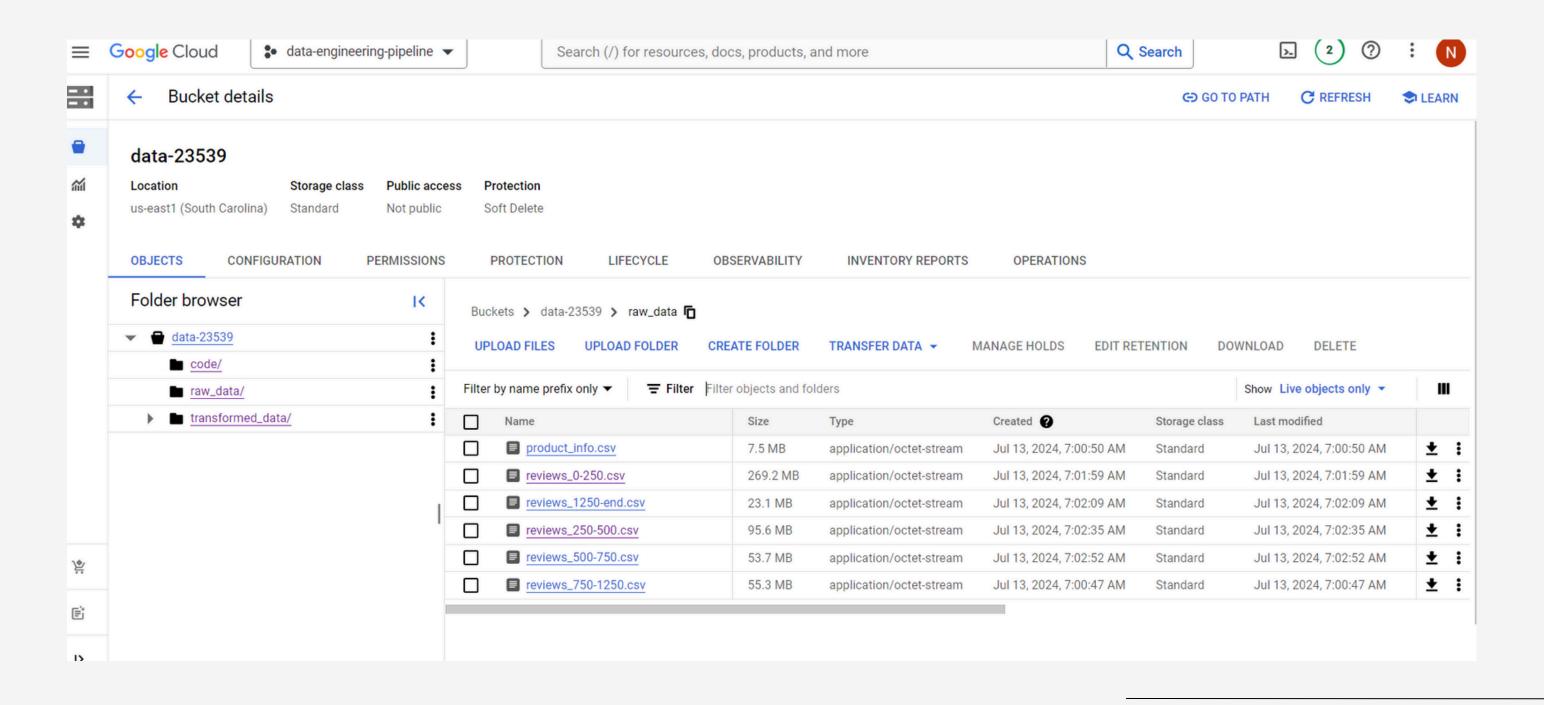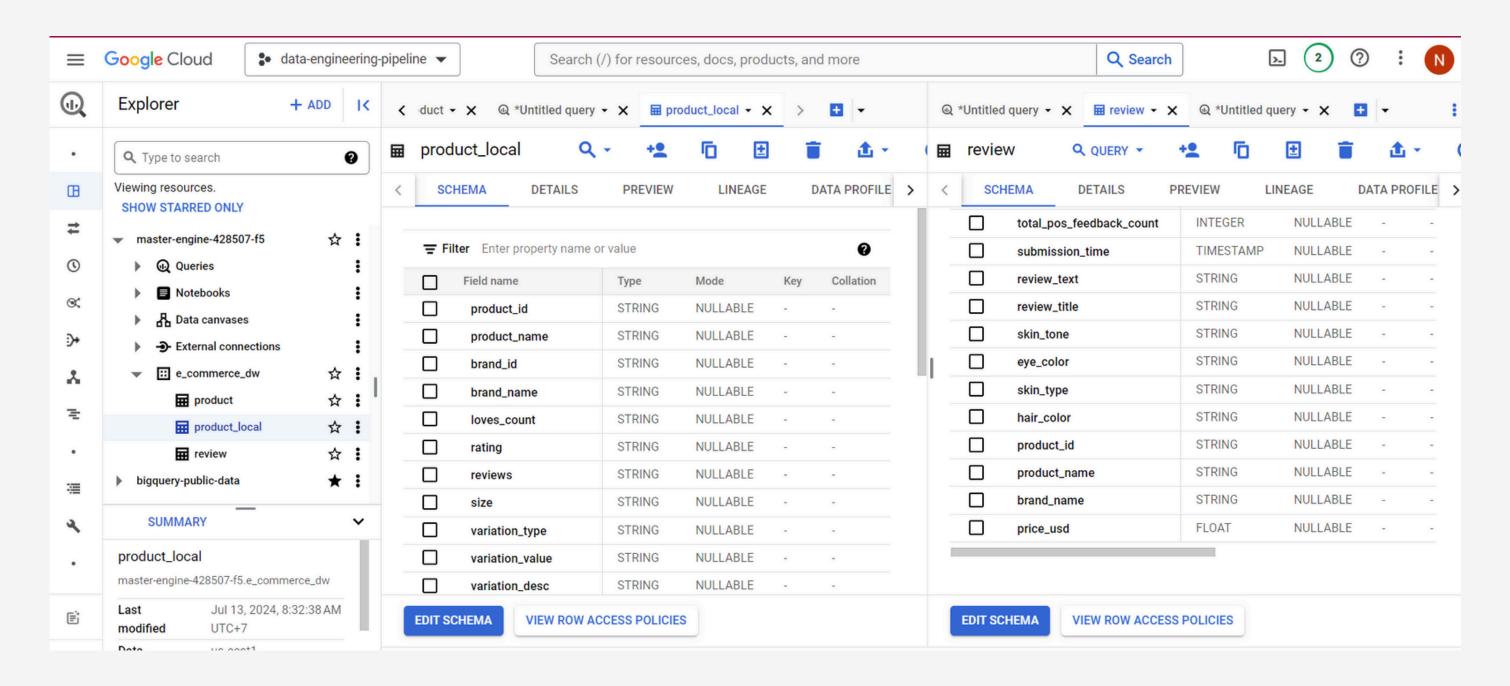
# Code Overview



## Airflow DAG

- kaggle_to_gcs, for dag kaggle api to gcs storage
- dataproc_gcs_to_bigquery
  - creating dataproc cluster
  - submit spark job to transform data to store it on gcs
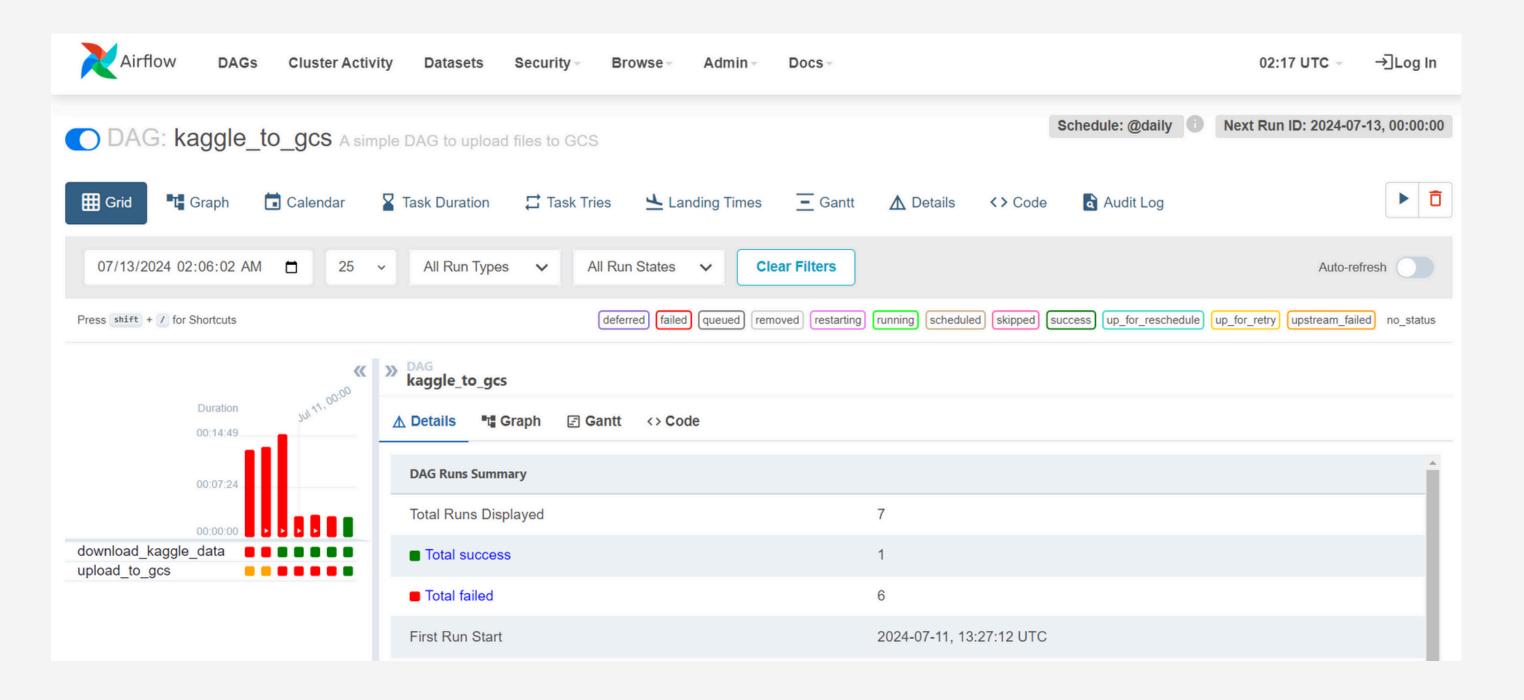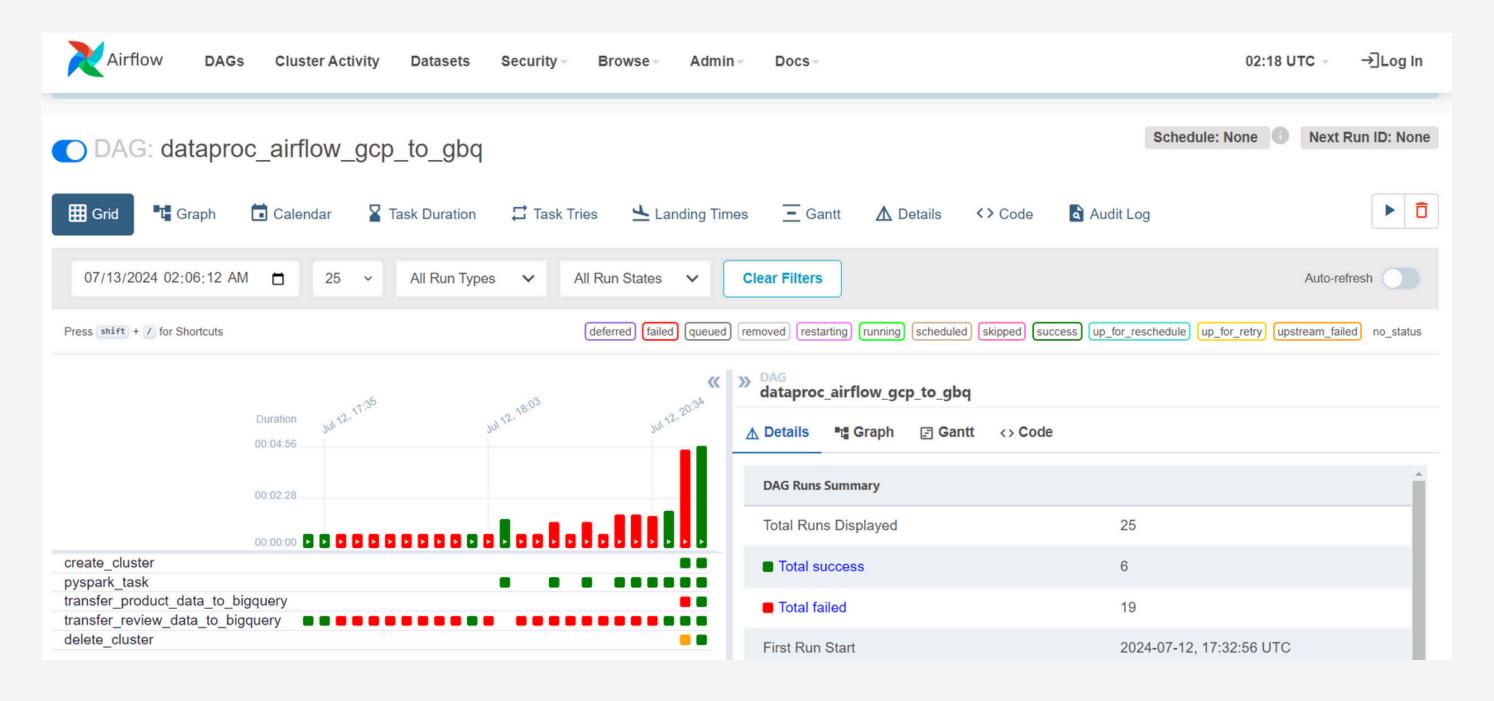  - store it to google bigquery

# GCP Overview

# GCP Overview

# Airflow Overview

# Airflow Overview

# Limitation and Improvement

## Limitation

- Using Denormalized schema that have low cost query but have more cost on storage and data redundancy
- Bad visualization and performance metrics
- No machine learning implementation.

## Improvement

- Data Modeling, For improvement we can use normalized schema (star-schema or snowflake schema),  or for more advanced or more big data use One Big Table or nested schema. Please see this article Here
- Data Quality, check some quality of the data using tools like Dataplex from GCP or using open source tool like Great Expectation
- More advanced visualization, including some more advanced performance metrics
- Some sentiment analysis or some machine learning model for predicting the product

# Thank You

# Problem Overview

- The problem found out when I try to upload the data from GCS to BigQuery, result in bad record because of csv format data
- Solved with change the transformed data to GCS with a parquet then upload it to BigQuery using parquet format which is parquet include the data schema.