# End-to-End ELT Pipeline with Airflow, GCP, BigQuery & Looker Studio

## 1. Introduction

This project demonstrates the implementation of an automated ELT (Extract, Load, Transform) pipeline using Google Cloud Platform (GCP) services and Apache Airflow. The pipeline processes a global health dataset, loads it into BigQuery, and produces visual insights using Looker Studio.

## 2. Architecture Overview

Data Flow: GCS → Apache Airflow → BigQuery → Looker Studio

The pipeline extracts a CSV file from a GCS bucket, loads it into BigQuery via Airflow, performs transformations to create country-level tables, and visualizes the data in Looker Studio.

## 3. Environment Setup & Requirements

The medical research team receives a global health statistics file with disease data by country.

### Setting up Airflow Environment

1. Launch a VM instance in GCP: Go to Compute Engine > VM Instances > Create Instance.

2. Name it 'airflow1'. Leave all other configurations at default.

3. Enable HTTPS traffic under Networking.

4. Under Security, allow full access to all Cloud APIs (Airflow needs access to BigQuery).

### Installing Airflow on VM

1. SSH into the VM. If blocked, configure the firewall appropriately.

2. Install Python using:

> *sudo apt install python3 python3-pip python3-venv*

3. Create and activate a virtual environment:

```
python3 -m venv ~/airflow-env

source airflow-env/bin/activate
```

4. Install Airflow with GCP support:

```
pip3 install apache-airflow[gcp]
```
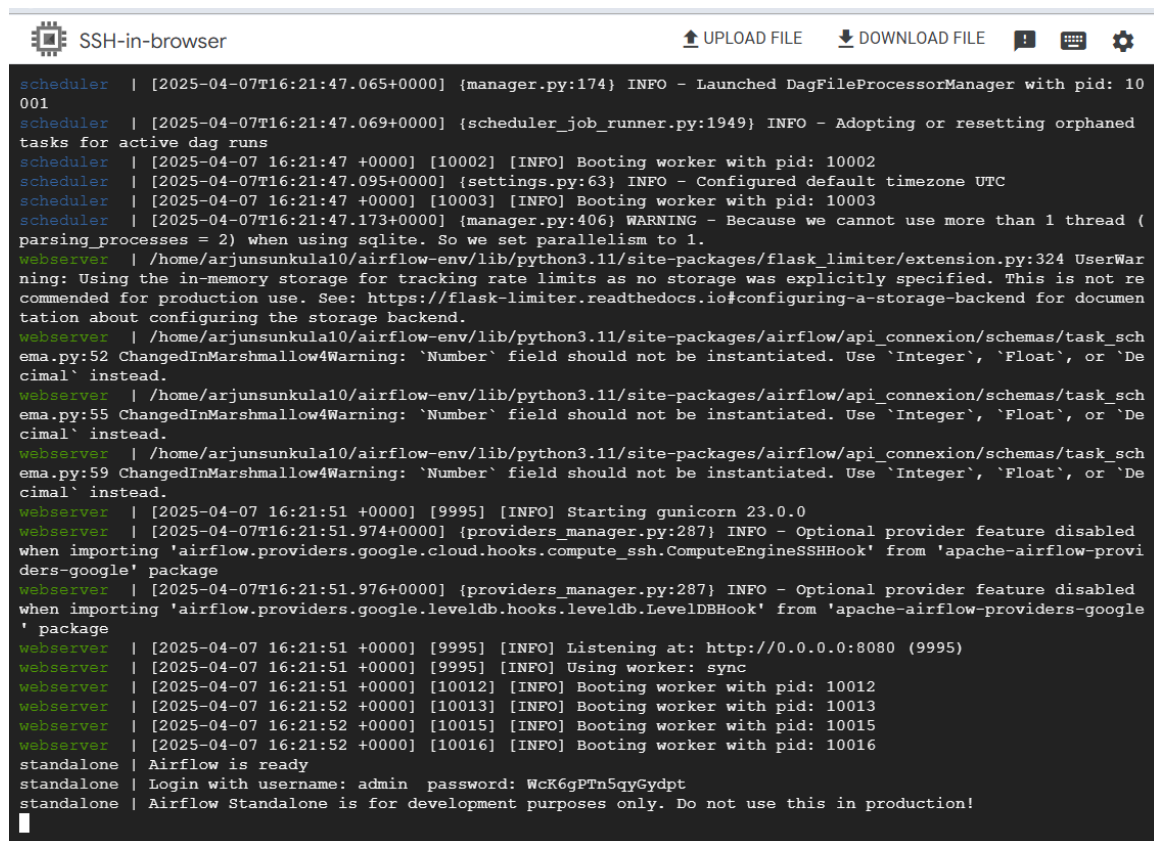
5. Run Airflow in standalone mode:

```
airflow standalone
```



6. Note the generated username and password for logging into the Airflow UI.

7. To run Airflow in the background:

```
nohup airflow standalone > airflow.log &
```

8. Ensure port 8080 is open in the firewall to access the Airflow UI.



## Optimizing Airflow Interface

Disable example DAGs in the Airflow config:

Set `load_examples = False` in airflow.cfg

Note the `dags_folder` path (e.g., /home/arjunsunkula10/airflow/dags)

## 4. GCS Bucket Setup

Create a bucket named 'airflow-bucket01' and upload 'global_health_data.csv'.
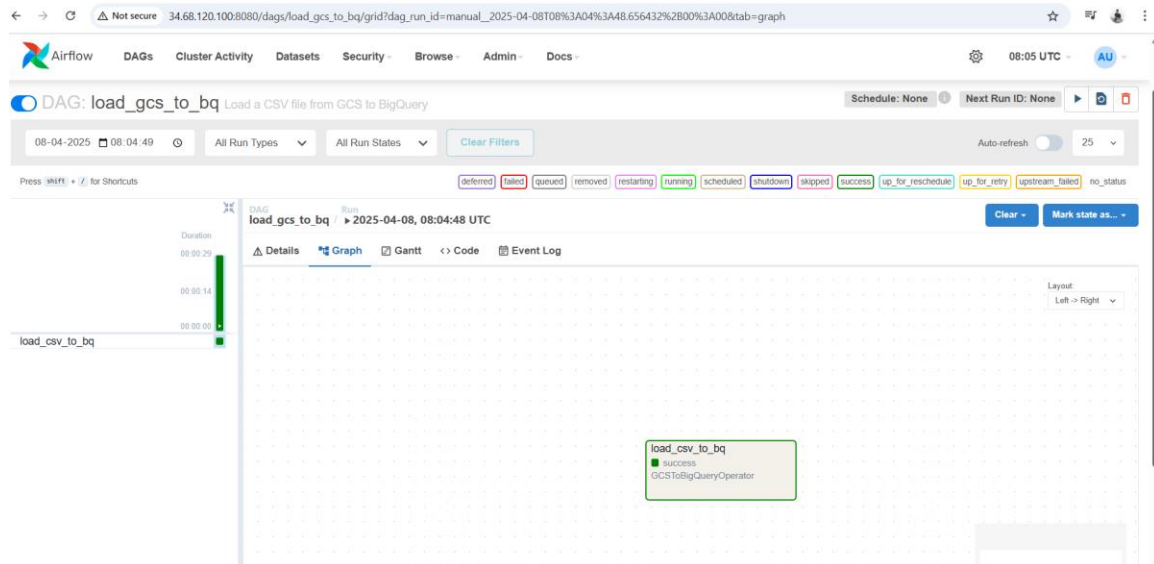


## 5. BigQuery Loading

1. In BigQuery, create a dataset in your project.

2. Load the CSV file from GCS into a table at the path:

*lateral-rider-456109-s6.initial_dataset.global_data*

3. Upload the DAG to the Airflow VM to automate loading.

4. Confirm data is visible in BigQuery.



# 6. File Detection and Data Transformation

Implement a sensor DAG to verify file availability in the bucket before proceeding with the load.

Upon confirmation, the data is processed to generate country-specific tables.



# 7. Creating Views for Optimization

Build views on top of transformed tables to support analytical queries.

These views reduce query size and enhance performance by focusing on filtered fields.



## 8. Visualization and Reporting

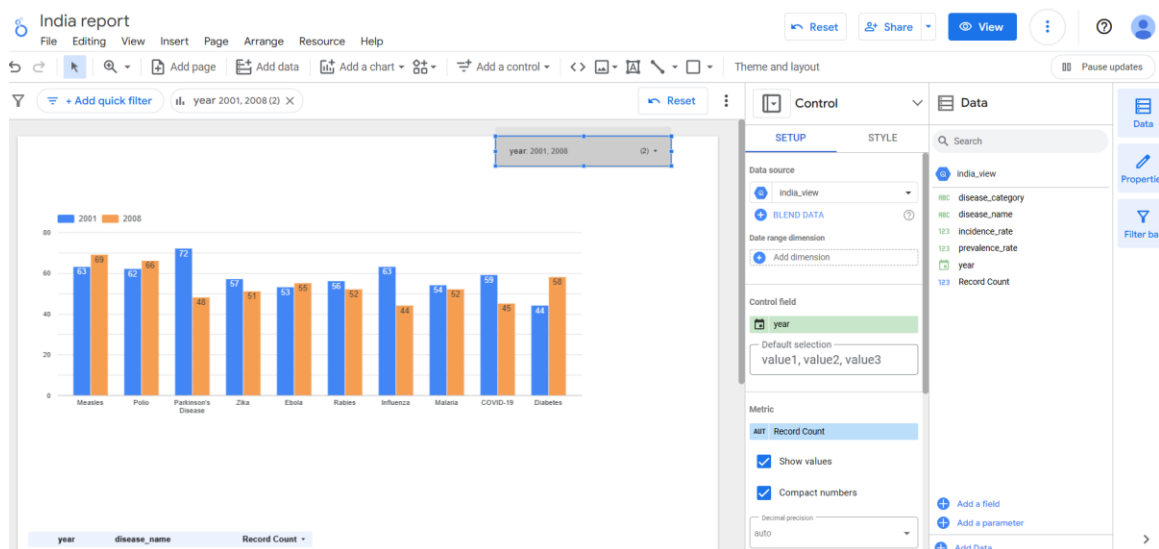Connect Looker Studio to BigQuery and create reports based on view tables.

Sample visualizations include disease trends by year and country.

Automated email alerts can be configured to share report summaries.



## 9. Conclusion

The pipeline provides a complete ELT solution from data ingestion to visualization. It supports reliable data validation, scalable storage, and cost-efficient querying.