**Project Title:** ETL Project with Data Fusion and BigQuery And Create Dashboard

**Project Overview:** This project is about building a full ETL (Extract, Transform, Load) process using Google Cloud tools. The goal is to clean employee data, move it through a cloud pipeline, and create a dashboard to see useful insights. All steps are done using simple and scalable tools provided by Google Cloud Platform (GCP).

**Tools Used:** - Python (for data cleaning) - Google Cloud Storage (to store data) - Google Cloud Data Fusion (to create ETL pipeline) - BigQuery (to store and query final data) - Looker Studio (to create dashboard)
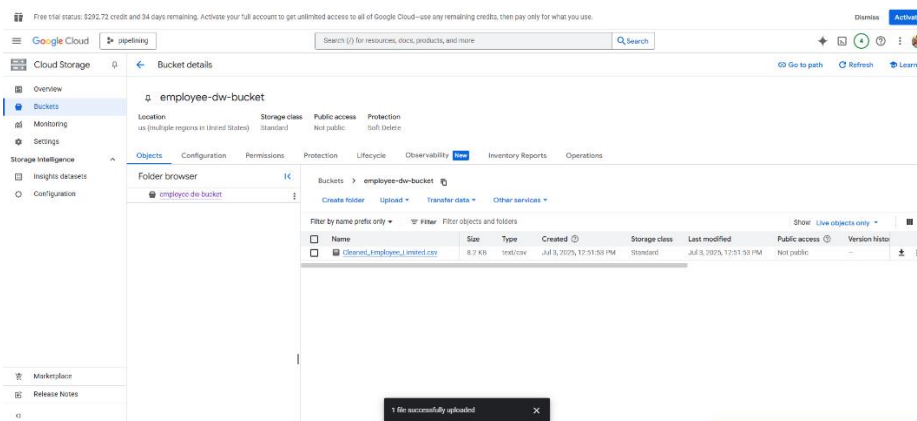
**Steps:**

1. **Data Cleaning with Python:**
   - Loaded the dataset (CSV file) using pandas
   - Removed duplicate records
   - Handled missing values
   - Limited the dataset to 200 clean rows
   - Saved the cleaned file

2. **Upload to Google Cloud Storage:**
   - Created a bucket in GCS
   - Uploaded the cleaned CSV file to the bucket



3. **Create ETL Pipeline with Data Fusion:**
   - Created a new Data Fusion instance
   - Apply data masking techniques to sensitive information in Cloud Data Fusion before loading it into BigQuery.
   - Designed a pipeline to read from Cloud Storage
   - Set the pipeline to load data into BigQuery
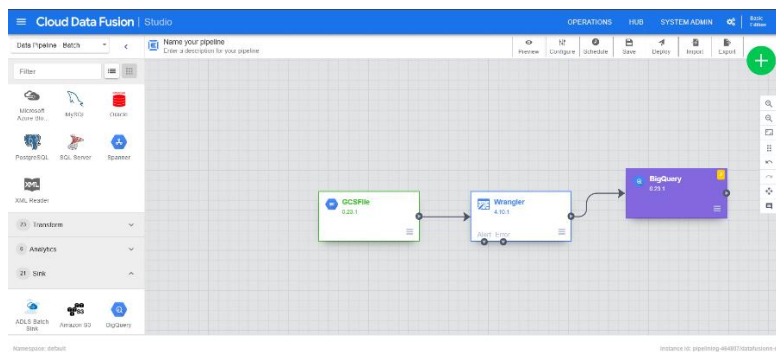   - Deployed the pipeline and confirmed data reached BigQuery

4. **Visualize in Looker Studio:**
    o Connected Looker Studio to BigQuery table
    o Created graphs and charts to show employee insights like age, city, etc.
    o Shared the dashboard for viewing

**Result:** The final project is a working cloud data pipeline that starts from raw data and ends in a clean, visual dashboard. This setup can be used for HR analytics or any company working with employee records.

**Conclusion:** This project shows how to build a complete ETL and BI (Business Intelligence) system using only GCP tools and Python, which is good for real-world use.