**Google Pipeline**

Suppose we are tasked to create an ETL Data Pipeline on Google Cloud with Airflow or Cloud Data Fusion, to extract some employee data for example from various sources, mask sensitive information within the data, and load it into BigQuery. It will require us to build and configure a data engineering infrastructure in google cloud platform.

**Requirements**

- **Data Extraction** from different sources such as databases, CSV files or APIs.
- **Data Masking (Transformation)** is done when sensitive information is identified
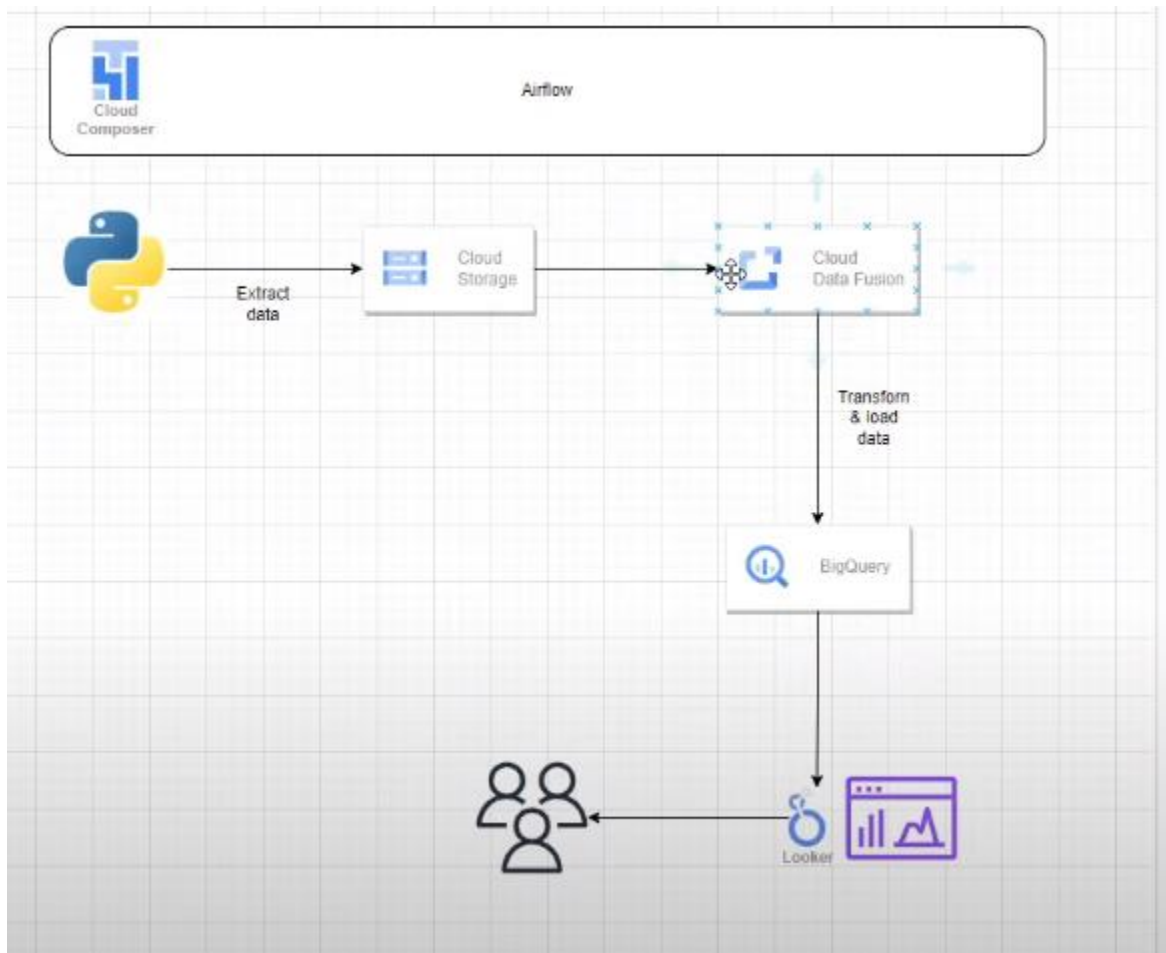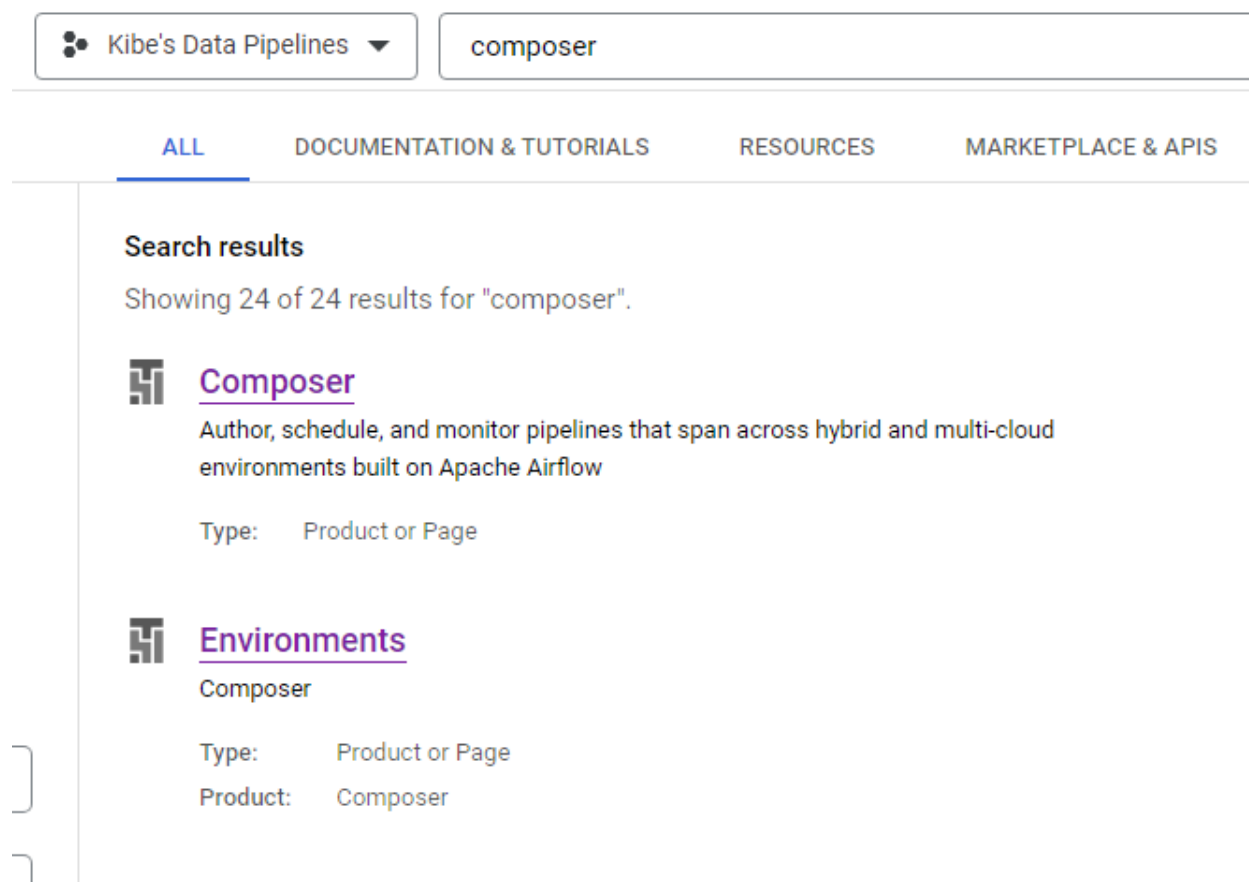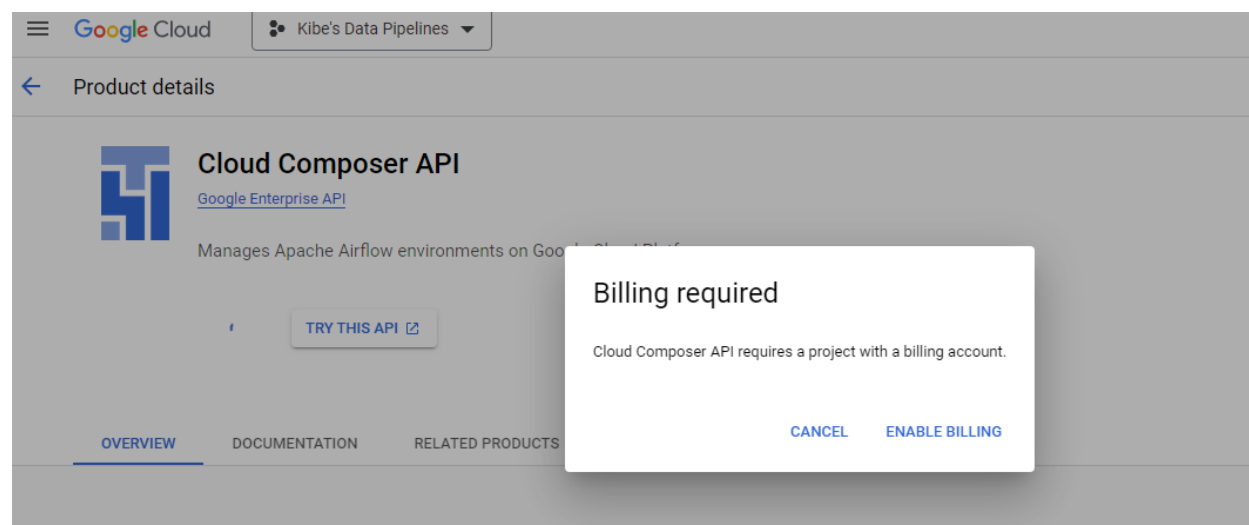- **Data Loading into BigQuery** is after identified data is extracted and masked



Figure 1: Data engineering infrastructure

This work flow is supposed to inform our Google pipeline project, following below steps

**Step 1: Create composer environment**



**Error**



We would have been able to create a composer environment which normally takes some 20-30 minutes.

**Step 2: Create a data fusion environment**



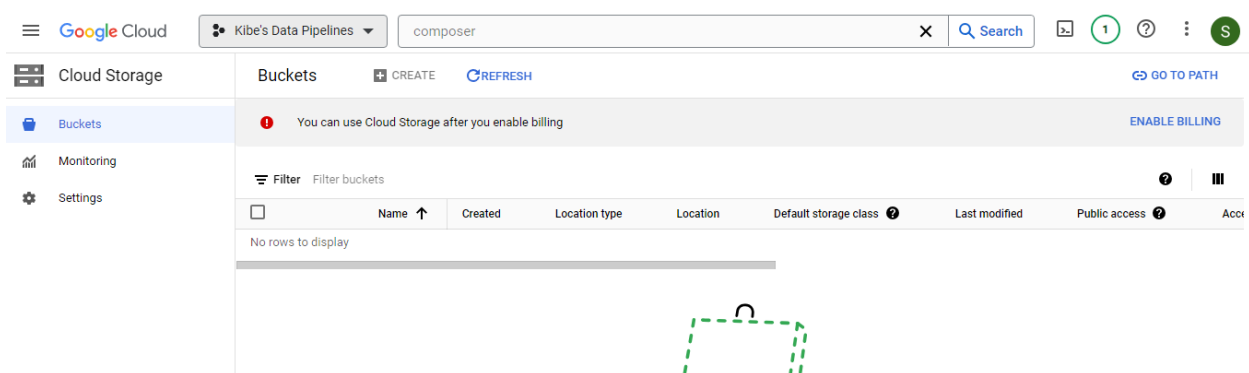Option 1

Option 2 recommended

**Error**



We would have been able to create a Data fusion Instance. Usually takes some time to create

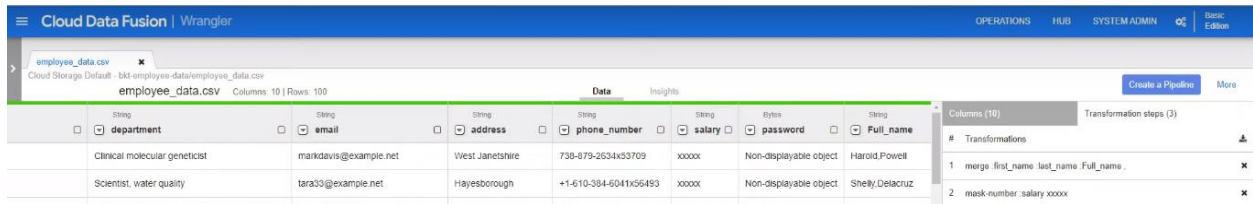## Step 3: Creating dummy data then load it to cloud storage

Basically, a python script with faker library is used to create dummy data, which afterwards is then loaded into the cloud storage bucket, saved as a CSV file. This is easily created and configured through any IDE for example Visual code, which connects direct to a project in google platform using a Gmail account.

**Error**

**Stem 4: Data pipeline**

Our dummy data is now supposed to be transferred into BigQuery using a data pipeline customized in the data fusion environment. A lot of data transformation (masking) can be done in Data fusion without any coding.



Figure 2: Sample Cloud Data fusion environment



Figure 3: sample ETL-pipeline on Cloud Data Fusion

**Step 4: Setup BigQuery Environment**





**Description**

With these packages set up and running, we would now be able to perform ETL, the load the processed data into BigQuery using google pipelining as illustrated by our data flow on the first slide.