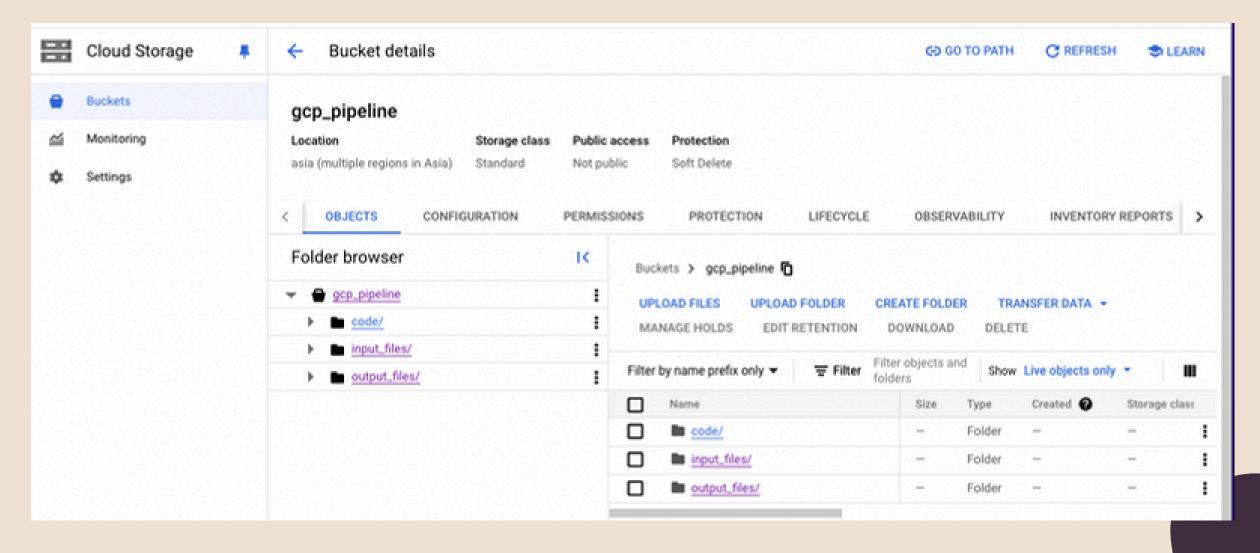# ETL GCP-BIGQUERY

Then Tsze Yen
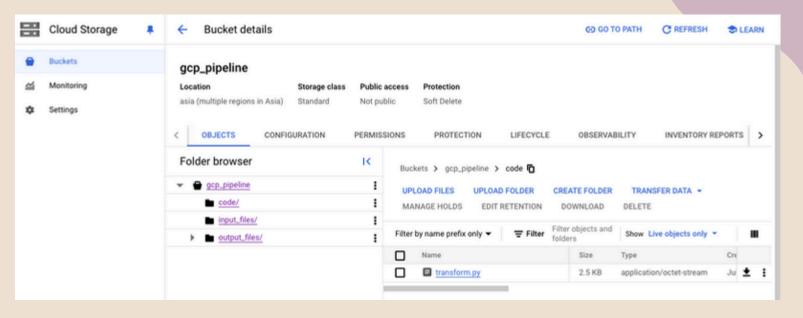
# STEP 1

-Create GCP Cloud Storage Bucket: gcp_pipeline
-Load data into GCP Cloud Storage: input_files/

# STEP 2

-Load the transform.py
to code file (created
your own)



```
 1 transform(used).py > ...
 7
 8   S3_DATA_SOURCE_PATH = 'gs://gcp_pipelines/input_files/superstore_dataset.csv'
 9   S3_DATA_OUTPUT_PATH = 'gs://gcp_pipelines/output_files/'
10
11   # # Define the schema
12   schema = StructType([
13       StructField("order_id", StringType(), False),
14       StructField("order_date", StringType(), False),
15       StructField("ship_date", StringType(), False),
16       StructField("customer", StringType(), True),
17       StructField("manufactory", StringType(), True),
18       StructField("product_name", StringType(), True),
19       StructField("segment", StringType(), True),
20       StructField("category", StringType(), True),
21       StructField("subcategory", StringType(), True),
22       StructField("region", StringType(), True),
23       StructField("zip", IntegerType(), True),
24       StructField("city", StringType(), True),
25       StructField("state", StringType(), True),
26       StructField("country", StringType(), True),
27       StructField("discount", FloatType(), True),
28       StructField("profit", FloatType(), True),
29       StructField("quantity", IntegerType(), True),
30       StructField("sales", FloatType(), True),
31       StructField("profit_margin", FloatType(), True)
32   ])
33
34   def func_run():
35       spark = SparkSession.builder.appName('airflow_with_emr').getOrCreate()
36
37       # all_data = spark.read.option("header", "true").schema(schema).csv(S3_DATA_SOURCE_PATH)
38       all_data = spark.read \
39           .option("header", "true") \
40           .option("dateFormat", "M/d/yyyy") \
41           .option("quote", "\"") \
42           .option("escape", "\"") \
43           .option("multiLine", "true") \
44           .option("inferSchema", "false") \
45           .schema(schema) \
46           .csv(S3_DATA_SOURCE_PATH)
47
48       # Convert order_date and ship_date to proper date format
49       all_data = all_data.withColumn("order_date", to_date("order_date", "M/d/yyyy")) \
50           .withColumn("ship_date", to_date("ship_date", "M/d/yyyy"))
51
52       selected_data = all_data.select(
53           'order_id', 'order_date', 'ship_date', 'customer', 'manufactory', 'product_name',
54           'segment', 'category', 'subcategory', 'state', 'country', 'discount',
55           'profit', 'quantity', 'sales', 'profit_margin'
```
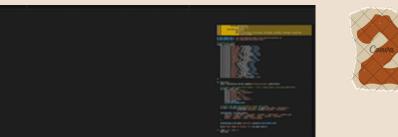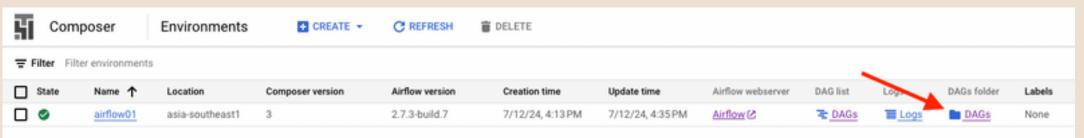
# STEP 3

-Load the load.py to DAGs from GCP Composer

# AIRFLOW OUTPUT

# BIGQUERY OUTPUT

# THANK YOU