

# Kindle Training – Major Project

By Monish NP

## **Overview**

- 1. Problem Statement
- 2. Architecture
- 3. Solution Approach
- 4. Data Visualization and Reporting
- 5. Technology Stack
- 6. Challenges faced and learnings



### **Problem Statement**



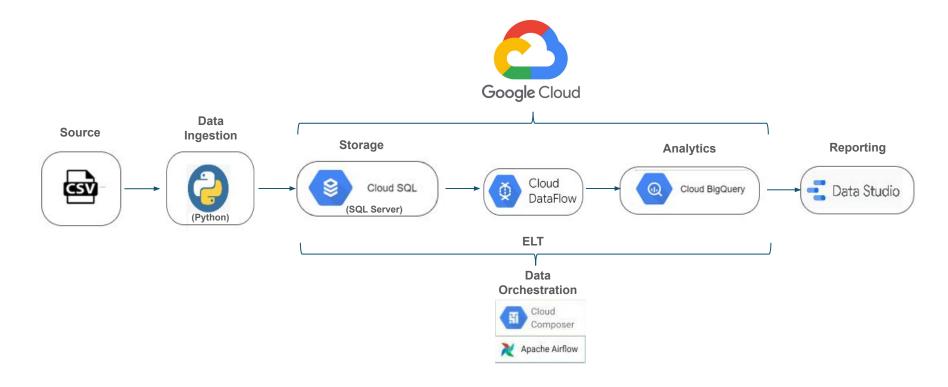
Customer Segmentation to focus on growth

Find 20% of customers who drive 80% of the company revenue

Visualize and report the findings

Analyze the visualized data for insights

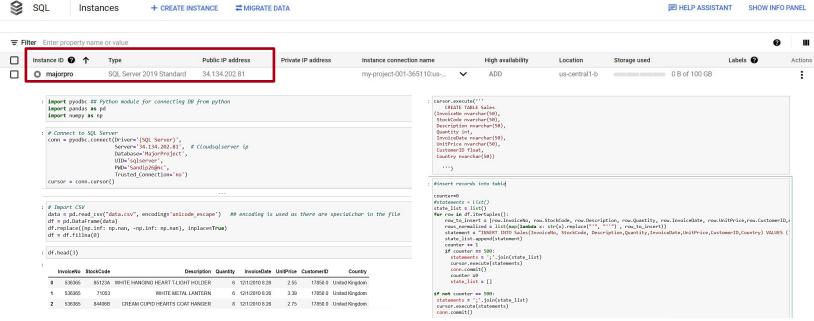
## **Architecture**





#### Step -1

- A Cloud SQL Server instance in GCP is created to import data into Cloud SQL.
- Wrote a python code and configured it with server credential to connect to Server instance.
- The below python code will create database and table in it with respective columns.





#### Step -2

Extract - Designed an Apache beam pipeline and executed it in dataflow to extract data into bigguery table.

```
p = beam.Pipeline(options=pipeline_options)

init = p | 'Begin pipeline with initiator' >> beam.Create(['All tables initializer'])
(init

# Read the file. This is the source of the pipeline. All further

# processing starts with lines read from the file. We use the input

# argument from the command line. We also skip the first line which is a

# header row.

| 'Read from a DB' >> beam.ParDo(ReadSQLTable())

# This stage of the pipeline translates from a CSV file single row

# input as a string, to a dictionary object consumable by BigQuery.

# It refers to a function we have written. This function will

# be run in parallel on different workers using input from the

# previous stage of the pipeline.

| 'White to BigQuery' >> beam.io.WriteToBigQuery('my-project-001-365110.Major_Project.sqlserver_to_bigquery',

schema='InvoiceNo:STRING, StockCode:STRING, Description:STRING, Quantity:INTEGER,'

'InvoiceDate:STRING,UnitPrice:FLOAT,CustomerlD:INTEGER,Country:STRING.'

p.run().wait_until_finish()

if _name_ == '_main_':
    logging.getLogger().setLevel(logging.INFO)
    run()
```

Figure: Pipeline Code

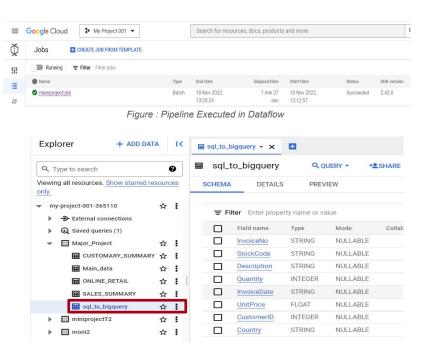


Figure: Table in Bigquery after Pipeline execution



#### Step -3

- ☐ Created an environment in Cloud Composer, uploaded the dag(python file) into the dag folder.
- Transform Queried raw data table and cleaned it using python Dag concept as per the project requirements.
- Load After transformation the data was loaded into respective tables i.e. (online retail. Customer summary, sales summary)
- The above ETL process is orchestrated using Cloud Composer

```
cloudsql_to_bq = BeamRunPythonPipelineOperator(
    task_id="cloudsql_to_bg",
    runner="DataflowRunner",
    py_file="gs://us-central1-mainpro-a06ea0ba-bucket/dags/sql to bq_moni.py",
    ppleline_options={\templocation': 'gs://majorpro/temp/', 'stagingLocation': 'gs://majorpro/temp/'},
    py_options=[],
    py_requirements=['apache-beam[gcp]==2.42.0','cloud-sql-python-connector[pytds]==0.6.1','pyodbc==4.0.34','SQLAlchemy==1.4.41','pymssql==
    py_interpreter='python3',
    py_system_site_packages=False,
    dataflow_config=DataflowConfiguration(job_name='mainproject-job', project_id='my-project-001-365110',
    location="us-central1"),
    dag=dag)

query1 = '''
    SELECT InvoiceNo, StockCode, Description, Quantity, PARSE_DATETIME('%m/%d/%Y %H:%M', InvoiceDate) AS InvoiceDate,
    UnitPrice, CustomerID,Country FROM `my-project-001-365110.Major_Project.sql_to_bigquery`
    '''
```





```
query4 = '''
           SELECT InvoiceNo, StockCode, Description, Quantity, InvoiceDate, UnitPrice, CustomerID, Country,
           Quantity * UnitPrice as ItemTotal
           FROM `my-project-001-365110.Major Project.Main data`;
ONLINE_RETAIL = BigQueryOperator(task_id='ONLINE_RETAIL',
                           destination_dataset_table= "my-project-001-365110.Major_Project.ONLINE_RETAIL",
                           sql=query4,
                           use_legacy_sql=False,
                           create_disposition="CREATE_IF_NEEDED",
                           write_disposition="WRITE_TRUNCATE",
                           dag=dag
query5 = '''
           select CustomerID, sum(ItemTotal) as TotalSales,count(Quantity) as OrderCount, avg(ItemTotal) as AvgOrderValue
           group by CustomerID;
CUSTOMARY_SUMMARY = BigQueryOperator(task_id='CUSTOMARY_SUMMARY',
                           destination dataset table= "my-project-001-365110.Major Project.CUSTOMARY SUMMARY",
                           sql=query5,
                           use_legacy_sql=False,
                           create_disposition="CREATE_IF_NEEDED",
                           write_disposition="WRITE_TRUNCATE",
                           dag=dag)
```

Figure: Query Creating Tables

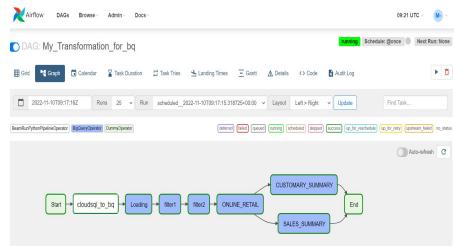
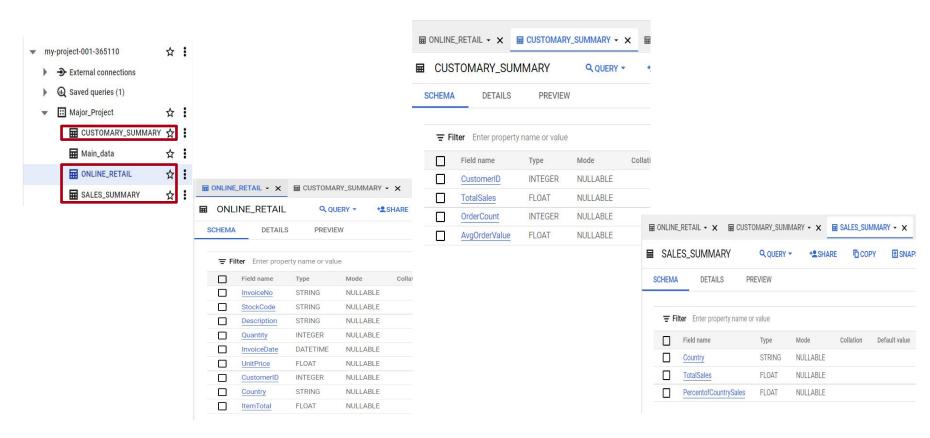


Figure : Air flow UI with Successful Orchestration



### **Tables Created**

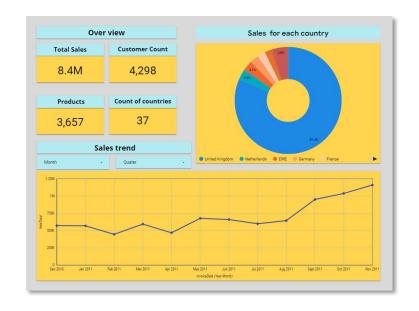




## **Data Visualization and Reporting**

- Using Bigquery as a data source in data studio we can access the tables stored in Bigquery
- The following questions are answered and a visualized which can be accessed from the link below
  - What are the sales figures for each country?
  - What is the overall sales trend?
  - How many new customers are there each month?
  - When do customers make the most purchases?
  - Which is the best selling product in each country?
  - When were the largest orders made?
  - Which customers made the largest orders?

Link - <a href="https://datastudio.google.com/s/o7bBBOcgd\_s">https://datastudio.google.com/s/o7bBBOcgd\_s</a>





## **Technology Stack**

Source

CSV File Cloud SQL Data Ingestion

Python Apache Beam Data Flow Data WareHouse

Big Query

Orchestrion

Cloud Composer Air Flow Reporting

**Data Studio** 



## **Challenges Faced**

- Data ingestion 

  Infinite values, Time consuming during record insertion.
- Running pipeline 

  Startup worker pool in zone Asia failed // replaced zone to us central.
- Run time error error 403 message forbidden □ Enabled data flow api, cloud sql api, granted permission for service account
- Orchestration 

  Name error "google" not defined 

  Created a requirement.txt file with dependencies

## Learnings

- Alternatively we can implement Federated Queries with MySQL option instead of SQL server
- Practical experience working with GCP services like Cloud shell, Cloud Sql, Data flow, Bigquery.
- Building and running Beam Pipeline.
- Creating Dashboard in Data Studio from BigQuery tables.



# Thank you!

