# Neobank Data Pipeline

Consultants: Alex, Enrico, Ikechi, Luiggi, Marlin

# Business Problem

Neobank is a online business delivering banking services all around the globe

The CEO wants use analytics to improve the company's performance. She asks the DE team to support the Finance and Marketing departments with self serving data insights to help their decision making:
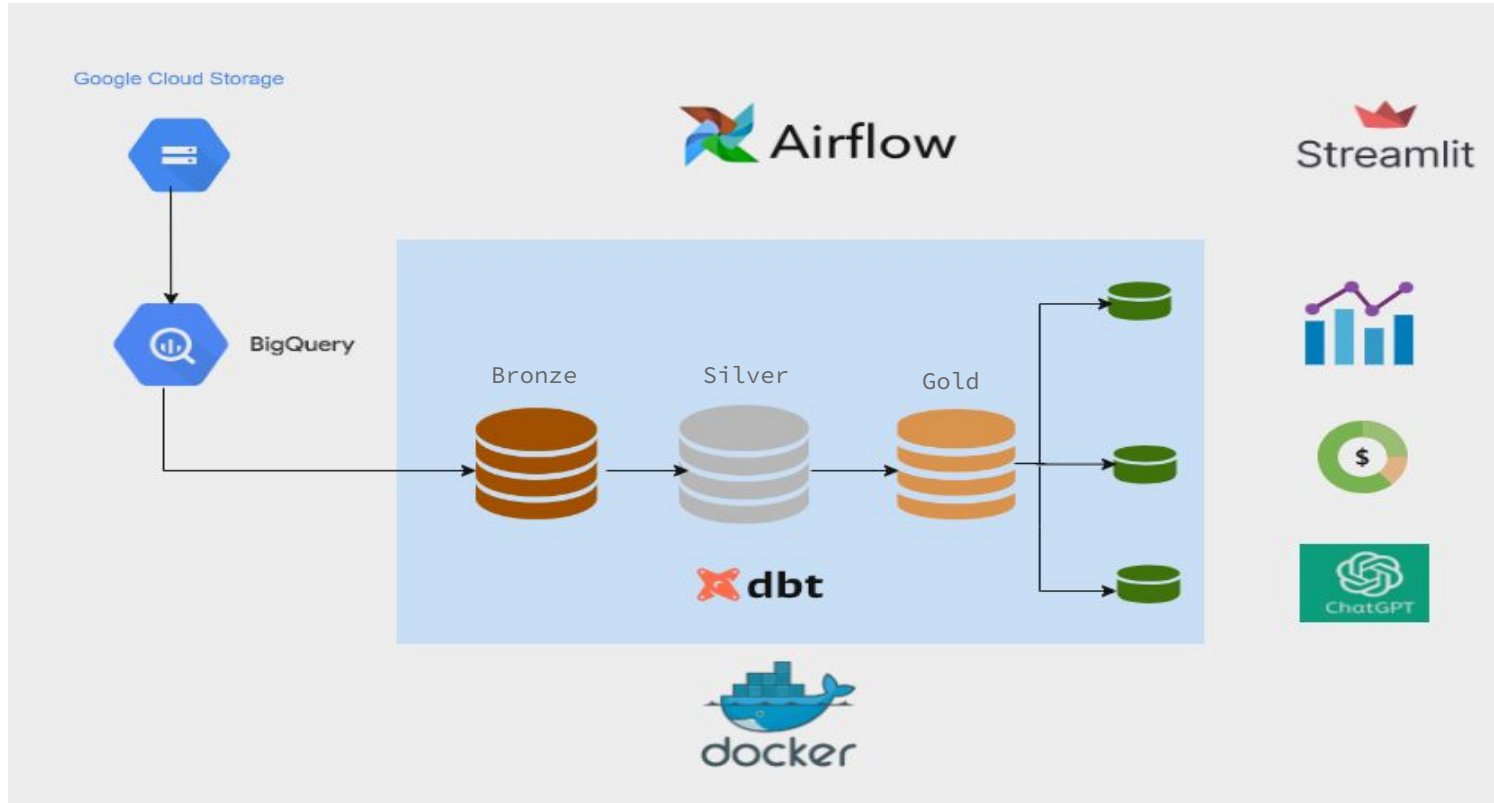
- What is the total amount of transactions by month/year?
- How much was charged in fees? Is it going up?
- Who are our users/customers?
- What are are the best communication channels?
- Are there any transactional patterns that could be of interest?

**NeoBank**

# Engineering Challenge

Neobank gathered unorganized data of 2.7 million transactions across the globe

- Extracting the data from the different data sources (csv files)
- Storing the data in data warehouse
- Data modelling according to the business requirements
- Connecting to an interactive Dashboard
- Connecting a LLM to the database
- Making the process scalable and auto scheduled

**NeoBank**

# Architecture



**NeoBank**

# Extraction and Load

- Extraction: Created data lake buckets. Raw data was uploaded to these GCP buckets


Google Cloud Storage

- Load (to Google Bigquery): Created data tables, schema definition and uploaded data to Google Bigquery bronze layer
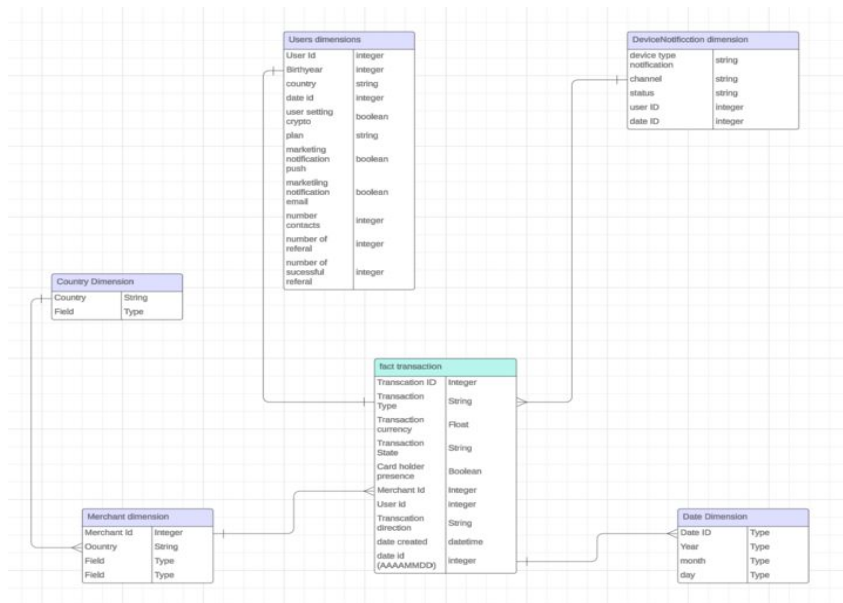

Google Big Query

**NeoBank**

# ⚙️ dbt Transformation

Bronze Layer: Raw Data uploaded from Data Lake (GCP bucket)

Silver Layer: Data Cleansing - Deleting unnecessary columns, formatting, Separating columns

Gold Layer: Building fact table and all dimension tables linking back to the fact table, creating a final star schema

## Final Star Schema



**NeoBank**

**NeoBank**

Streamlit **+ Demo**

https://dedemo-neobank.streamlit.app/

NeoBank

# QR Code



**NeoBank**

# Challenges

- Github collaboration
- Package dependencies
- Connecting our Silver layer to a gpt3.5 model
- Securely deploying Streamlit (python versions, Streamlit secrets)

**NeoBank**

# Future Improvements

- Auto orchestration of the extract and load process with Airflow
- Auto creation of table schema when uploading into BigQuery
- Create a CI/CD pipeline for seamless integration
- For own workflow:
  - Adding test files and documentation
  - Makefile to build, compile and organise source code
- Improving the LLM. We have a working variant, but we need to finetune it for more complex questions

**NeoBank**

# Thank You

NeoBank

# Tech Stack

Google Big Query

Apache Airflow

ChatGPT

docker

GitHub

dbt

Streamlit

NeoBank