

Data Engineering Program Final Project - Background and Study Case Property of Rizky Faiar Aditya

Study Case: Go Hailing Inc. faces growing challenges in Data Access and Quality, hindering decision-making and business agility.

Background and Study Case:



Go-Hailing is a ride-hailing company that provides on-demand transportation services.

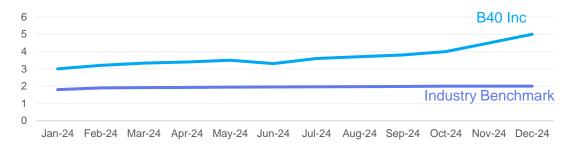
The Sales and Operation Team at Go Hailing. is facing challenges in gathering data for their reports and insights, which is impacting their decision-making processes. To address this, they have enlisted the help of Data Engineering team to design and implement a robust and scalable data pipeline that will streamline data collection and processing, enabling faster reporting and more informed decision-making.

Please note: The case study and data presented in this slide are hypothetical and created for illustrative purposes only

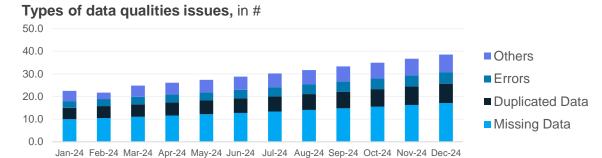
Several challenges that makes Go-Hailing Inc is on the verge to develop a robust data pipeline:

The time required to access critical data increases by 4% each month, as the duration from data collection to reporting grows, which delays their ability to gain insights and make decisions

Time needed teed to collect data to reporting, in hours

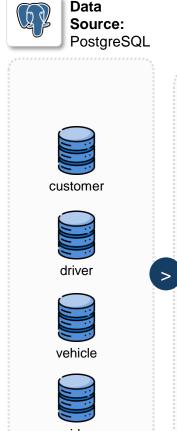


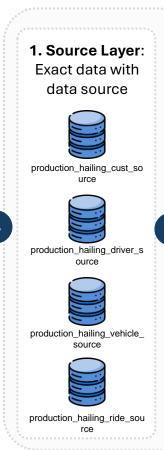
Inconsistent data quality makes it challenging to perform comprehensive data analysis

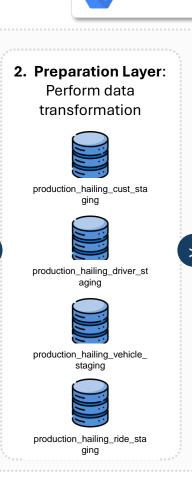


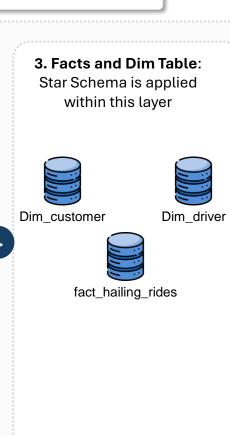
The 1st stage to resolve the issue, I developed the end-to-end data pipeline as a framework reference.

Data Warehouse: BigQuery











Remarks:

- Data pipeline is orchestrated using **Airflow**
- **DBT (Data Build Tool)** is used for data modelling in Data Warehouse
- All apps are containerized using Docker.
- If any tasks failed or retry it will be notified through **Discord** WebHok
- **ELT (Extract** Transform and Load) is used in this data pipeline
- Emphasizing batching data pipeline

Data Engineering Program Final Project – Development Cycle and Process

Property of Rizky Fajar Aditya

The 2nd stage is to develop the codes based on the data pipeline architecture design

Code	Code Objective	File Name	Main Logic/Function
Docker	Build Airflow Build PostgreSQL Build DBT	 docker-compose.yaml (Airflow and dbt) click here >> docker-compose.yaml (postgresQL) click here >> 	 Airflow – build postgre for airflow metadata >> build scheduler for orchestrating tasks >> build webserver to access and trigger DAGS interactively PostgreSQL – build postgre in different container to store dummy data DBT (Data Build Tool) – Build DBT in different container to run models for data modelling within data warehouse. Airflow can send command to DBT to run task using BashOperator (docker exec –it dbt run models). Credentials is also mounted in the volume so it can access BigQuery.
DAG 1	Generate dummy tablesLoad to PostgreSQL	 DAG1_load_to_postgre.py as main DAG << click here >> Generate_data.py as helper << click here Postgres_helper.py as helper << click <p>here >> </p> 	Execute .sql files to create schema in PostgreSQL >> generate and insert data to PostgresSQL
DAG 2	Data Ingestion to BigQuery	 DAG2_load_to_bigquery as main DAG < click here >> Bigquery_helper.py as helper << click here >> 	 Ensure if table exists (creates if not exists) >> fetch data from postgresql >> return to dataframe in pandas >> map data types befor upsert in BigQuery >> perform upsert in bigquery, partitioned by created_at Table that has no created_at/partitioned columed can also be load into bigquery – emphasizing for modularity
DAG	Perform data transformation after ingestion	 DAG3 as main DAG << click here >> Models/staging as models in preparation layer << click here >> Models/facts/ as models in fact and dim table << click here >> Models/marts/ as sql models in marts layer << click here >> 	Perform data transformation (unify phone number) from source layer to preparation layer >> build star schema model in facts and dim layer >> build aggregated data to gain insights from rides activities, and driver idleness.

Data Engineering Program Final Project - Project Results Property of Rizky Fajar Aditya

Project Results

Successfully generate data, load data, and build relationship between tables in PostgresSQL- all orchestrated using Airflow

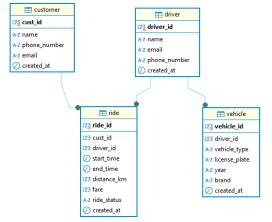


Fig 1. ERD in postgreSQL

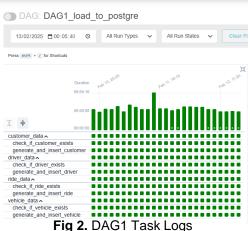


Fig 4. DAG2 Task Logs

- Successfully performed ELT ingesting data from postgre to bigguery incrementally
- Performed data transformation from source table to mart-layer using DBT
- Created STAR SCHEMA in facts tables

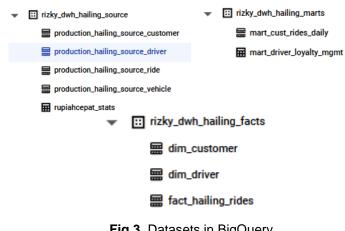
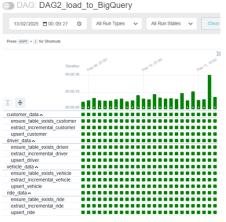


Fig 3. Datasets in BigQuery



Developed a function to send notification through Discord Webhook if a task faile or attempts to retry

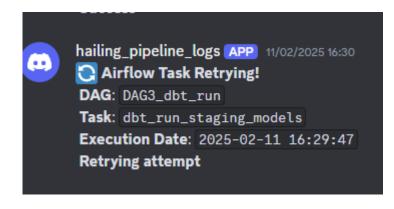


Fig 5. Retry task notification

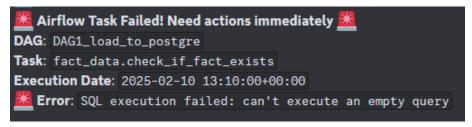
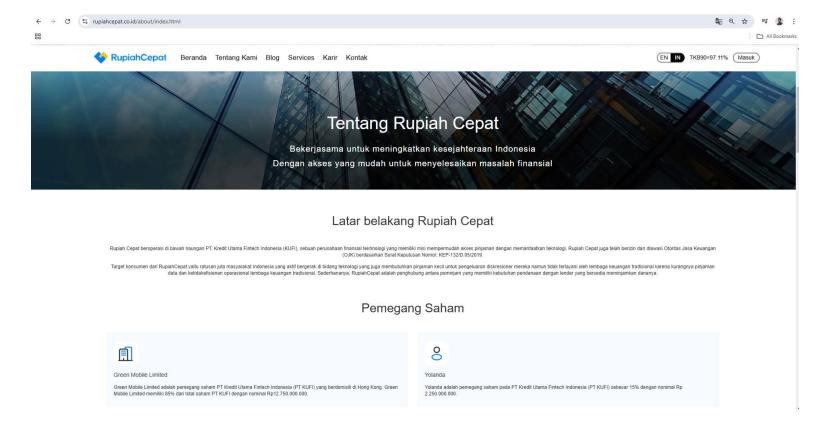


Fig 6. Failed task notification



Final Project - Background and Study Case Data Engineering Program Property of Rizky Fajar Aditya

Study Case: Web Scraping dynamic website of Rupiah Cepat



Project Objective: Perform web scraping from one of P2P lending website (RupiahCepat) and load to Data Warehouse (BigQuery)

Data Engineering Program **Capstone 3 – Development Cycle and Process** Property of Rizky Fajar Aditya

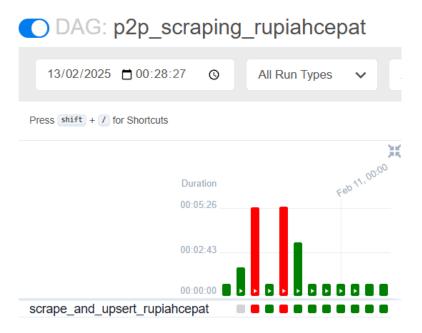
Project;s Result:

No	Activity	Code Development	Main Logic/Function
1	Build Airflow in Docker	• << click here >>	Build Airflow using requirements.txt that includes Selenium and WebDriver
2	Build Selenium	• << click here >>	I built the selenium and chrome driver in same docker-compose.yml file as Airflow. This makes it simpler to running the apps
3	Create DAG4 to automate data scraping and load to bigquery	• << click here >>	Check if table exists in BigQuery >> Scrape data using remote server (this allows airflow to command selenium container to scrape the data >> find element (in my case in <div div[@data-v-35084c97])[1]"="">> map to column/headers >> return df >> load to BigQuery using upsert (any changes within the data will be inserted as new row and set the old data inactive)</div>

Data Engineering Program Final Project - Project Results Property of Rizky Fajar Aditya

Project Results

- Successfully scraped the data from Rupiah cepat
 - Load to BigQuery using merge update
 - Successfully automated using Airflow



Mark state as... ▼ Filter DAG by task + p2p_scraping_rupiahcepat / @2025-02-12, 00:00:00 WIB / scrape_and_upsert_rupiahcepat <> Code E Event Log All Levels All File Sources Wrap See More Download afa42d46628a *** Found local files: *** * /opt/airflow/logs/dag_id=p2p_scraping_rupiahcepat/run_id=scheduled__2025-02-11T17:00:00+00:00/task_id=scrape_and_upsert_rupiahcepat/attempt=1.log [2025-02-13, 00:05:07 WIB] {local_task_job_runner.py:123} ▶ Pre task execution logs [2025-02-13, 00:05:07 WIB] {scraping_helper.py:54} INFO - 🗹 Table purwadika.rizky_dwh_hailing_source.rupiahcepat_stats already exists. [2025-02-13, 00:05:09 WIB] {scraping helper.py:83} INFO - Loading URL: https://www.rupiahcepat.co.id/about/index.html [2025-02-13, 00:05:19 WIB] {scraping_helper.py:86} INFO - Page title: Rupiah Cepat | Pinjaman Online [2025-02-13, 00:05:19 WIB] {scraping helper.py:166} INFO - DataFrame to be inserted: [2025-02-13, 00:05:19 WIB] {scraping_helper.py:167} INFO - organization_name ... is_active RupiahCepat ... [1 rows x 13 columns] [2025-02-13, 00:05:22 WIB] {scraping_helper.py:174} INFO - 💟 Data successfully batch-loaded into staging table purwadika.rizky_dwh_hailing_source.rupiahcepat_staging [2025-02-13, 00:05:24 WIB] {scraping_helper.py:208} INFO - V Upsert completed successfully in purwadika.rizky_dwh_hailing_source.rupiahcepat_stats [2025-02-13, 00:05:25 WIB] {scraping_helper.py:214} INFO - 💹 Staging table purwadika.rizky_dwh_hailing_source.rupiahcepat_staging successfully dropped. [2025-02-13, 00:05:25 WIB] {scraping_helper.py:225} INFO - Data scraping and upsert complete! [2025-02-13, 00:05:25 WIB] {python.py:240} INFO - Done. Returned value was: None [2025-02-13, 00:05:25 WIB] {taskinstance.py:340} ▶ Post task execution logs

Fig 1. DAG4 Historical Task Logs

Fig 2. DAG4 Detailed Task Logs

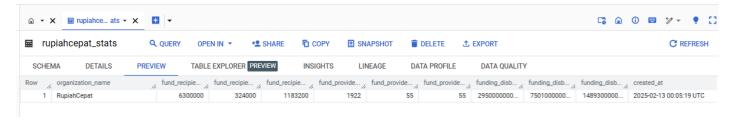


Fig 3. Data scraped and ingested in BigQuery

