# Dynamic ELT Jobstreet Data using Airflow and GCP

## Muhammad Sidqy D

https://github.com/msidqyd/ELT-to-GBQ-Jobstreetscrape-Detailed

# Education

Physics Engineering, institute Technology Sepuluh Nopember

# Working

- Management Trainee – Korea Tomorrow & Global Indonesia

- Retail Sales Supervisor – Korea Tomorrow & Global Indonesia

- Productivity Specialist - Korea Tomorrow & Global Indonesia

- Field Operation Specialist - Korea Tomorrow & Global Indonesia

# Project Overview

- Modular and scalable ELT pipeline

- Automate job data collection

- Multithreaded Web Scraping with Selenium

- Branching DAGs for dynamic input format handling

- Incremental & full load

- Stored intermediate files in local staging as Parquet

- transformations and load data to Big Query

- Data validation tasks after loading

- Data Modelling : Silver -> Bronze -> Gold

# Background

- Indonesia's unemployment rate in Feb 2025 reached 4.76%, above the healthy range of 3.5%–4.5%.
- The government needs to promote digital upskilling, especially in high-demand roles like Data Engineering, which is projected to be a top-growth job for more than 30%.
- remote job for data engineer in US reach 10% in 2024, but the number has plummet less than 2% in 2025.
- This project aims to collect and analyze Data Engineer job postings to identify the most in-demand skills and tools.

# Problem Formulation

| Gap | Opportunity | Solution | Impact | How to | Project Output |
|---|---|---|---|---|---|
| There's a gap reach 0.26% of unemployee rate need to reduce to reach the bare minimum of healthy economy. | The digital field has thrive, one of them is data engineer | Offering data for learning and coaching tools and skill needed as data engineer | The people who want to learn about data engineer could know the trend skill and demand for data engineer role. | | Show keywords with the frequencies for skill and tools data from job posting in indonesia or overseas platform. |
| | Data Engineer role projected to be a top growth job which reach more than 30% growth | Conduct training with curricullum alligned to data engineer job-market trend. | Private, public, and third sector, could have training with curriculum the most demanding and tren tools to accelerate the student learning process. | 1. Get the trend skills and tools from the job platform.<br>2. Build pipeline to ingest data and load it to database.<br>3. Transform data for collection and analyzing the data. | Dashboard for another information such as total jobs available, total jobs each location, and average wage. Also utilize the data for future research (forecast the job availability) |
| | There are 10% remote jobs in based on job posting from indeed in US for data engineer jobs | Generate indonesian talent to compete in indonesia or global job market. | Actively for contribution of reducing the unemployement rate for current and future condition | | Collect and stay updated with job list data to domestic or global, so the talent couldn't miss the job posting or opportunity |

# Data Platform Understanding



## Architecture

### 1. Scraping
JobStreet search pages and job detail pages (dynamic web content).
Targeting job postings with the keyword: "Data Engineer" in Indonesia and
Singapore with multithreading regarding to scrape job page.

### 2. Ingestion
Scheduled with Apache Airflow to run regularly and support incremental scraping.
Data saved as raw JSON or CSV.

### 3. Storage (Bronze Layer)
Raw data stored locally or in Google Cloud Storage (GCS).
File formats: Parquet/CSV for efficient querying and processing.

### 4. Transformation (Silver Layer)
Parsed and cleaned text (job titles, descriptions, company names, locations).
Extracted structured features: skills, tools, experience levels, job types.
Transformation done using Pandas, SQL

### 5. Gold (Gold Layer)
Final dataset loaded into BigQuery for analytics and dashboarding.
Allows for aggregation, keyword frequency, job listing

### 6. Google Big Query
Word Frequencies : identified skills, tools, traits in-demand at the market.
Dashboard: Further analysis current condition for data engineer job.
Job Market Search: Updated and sorted Data Engineer Job with all information.

# Data Platform Understanding

## Features

1. Performed multithreading using Selenium (headless browser automation).
Scheduled with airflow with config parameter
2. Dynamic configuration and parameterization
3. Modular & Scalable ELT architecture
4. Automate job everyday 21.15 O'clock

```python
# Load source list, based on the country want to scraped.
with open("dags/resources/dynamic-dag/list_source.yaml") as f:
    list_data_source = yaml.safe_load(f)

# Create dag with schedule at 21:15 everyday, also i add the config parameter for reducing manual intervention.
def create_elt_dag(source_data):
    @dag(
        #DAG name
        dag_id=f"Main_ELT_GCP_{source_data}",
        #DAG will be executed everyday at 21.15 PM
        schedule_interval="15 21 * * *",
        #DAG Start Date is 7 June 2025, so after the calendar the DAG coul be run automatically when active.
        start_date=datetime(2025, 6, 7, tzinfo=pytz.timezone("Asia/Jakarta")),
        catchup=False,
        tags=["ELT_Scrape", "To_GBQ"],
        #Config parameter setting
        params={
            "source_type": "CSV",
            "load_type": Param("incremental", description="incremental/full", enum=["full", "incremental"]),
            "table_date": Param(datetime.today().strftime('%Y-%m-%d'), description="Choose Date Table for ELT To GBQ"),
        }
    )
```



Airflow · DAGs · Cluster Activity · Datasets · Security · Browse · Admin · Docs · 21:42 WIB (+07:00)

Trigger DAG: Main_DAG_ELT_GBQ_Source_Data_jobstreetscrape_singapore

Select Recent Configurations
Default parameters

DAG conf Parameters

| source_type: | CSV |
| load_type: | incremental |
| | incremental/full |
| table_date: | 2025-07-19 |
| | Choose Date Table for ELT To GBQ |

Generated Configuration JSON and Dagrun Options

[Trigger] [Cancel]

| trigger_scrape_ingest | trigger_load | trigger_transform | trigger_gold_layer |
| ■ success | ■ success | ■ success | ■ success |
| TriggerDagRunOperator | TriggerDagRunOperator | TriggerDagRunOperator | TriggerDagRunOperator |

# Data Understanding



## 1. Search page

- **job_id**
- **title** — Job will selected if title match with more than 1 keywords
- **url**
- **company_name** — For incremental, filtered by max publish time from GBQ
- **publish_time**
- **location**

From URL use multithreading to open each job page

## 2. Job Page

- **wage**
- **work_type**
- **job_desc**

### Google Big Query (Bronze layer)

### Data Staging

Use Branch Operator CSV or JSON and combine all keywords data

Use Branch Operator CSV or JSON and combine all keywords data

# Data Understanding

## DAG for Scrape & Ingestion

# Transformation & consideration

**Transformation (Silver Layer)**

1. Separate salary to be absolut number, to be salary_min_cleaned and salary_max_cleaned .



2. Fill column for blank value of company.



4. Parsed and cleaned text (title, job_desc, company, location).

4. Fill the work_type with value NULL



5. Drop duplicate regarding to job_id_platform.

6. Null posted time filled by '1900-01-01 00:00:00'

7. Data validation function to ensure the table has no duplicate and null value regarding to job_id_platform.

# Data Modelling

**Bronze Schema**

| Field name | Type |
|---|---|
| Job_ID | INTEGER |
| Role | STRING |
| Company | STRING |
| Location | STRING |
| Publish_Time | STRING |
| URL | STRING |
| job_desc | STRING |
| salary | STRING |
| work_type | STRING |
| country | STRING |

**Silver Schema**

| Field name | Type |
|---|---|
| job_id | INTEGER |
| job_id_platform | STRING |
| role | STRING |
| company | STRING |
| location | STRING |
| posted_time | TIMESTAMP |
| url | STRING |
| job_desc | STRING |
| salary | STRING |
| work_type | STRING |
| country | STRING |
| salary_min_cleaned | NUMERIC |
| salary_max_cleaned | NUMERIC |

**Gold Schema**

| Field name | Type |
|---|---|
| word | STRING |
| frequency | INTEGER |

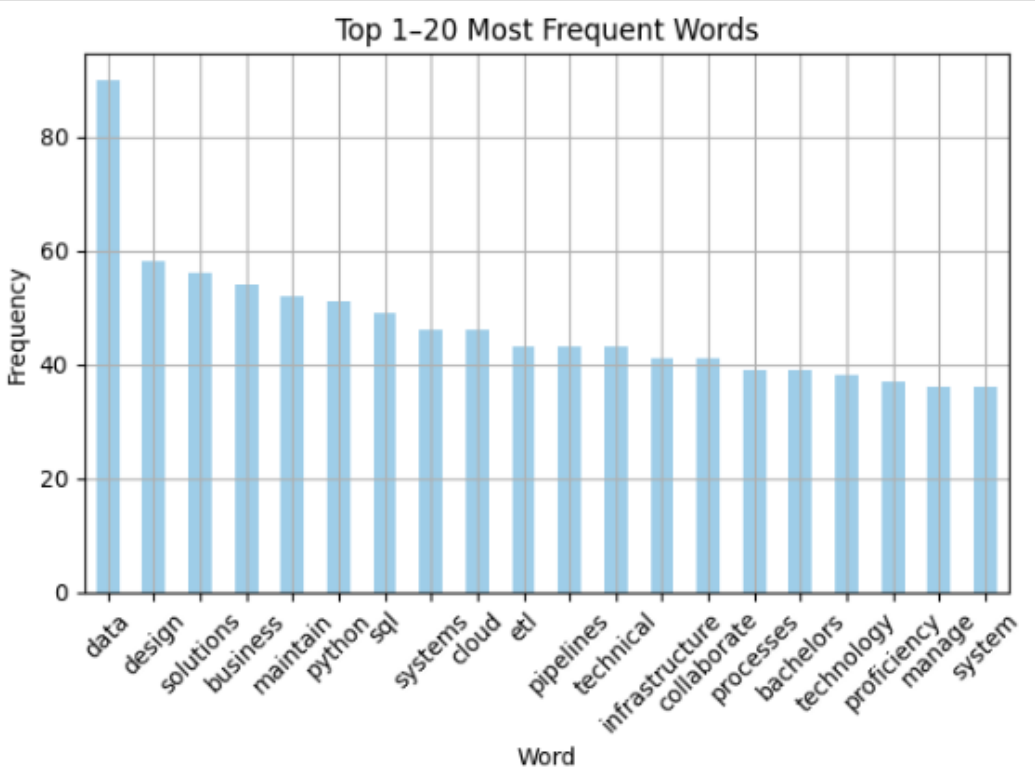| Field name | Type |
|---|---|
| job_id | INTEGER |
| job_id_platform | STRING |
| role | STRING |
| company | STRING |
| location | STRING |
| posted_time | TIMESTAMP |
| url | STRING |
| job_desc | STRING |
| salary | STRING |
| work_type | STRING |
| country | STRING |
| salary_min_cleaned | NUMERIC |
| salary_max_cleaned | NUMERIC |

word_frequency table utilized for identify skill and tools in demand

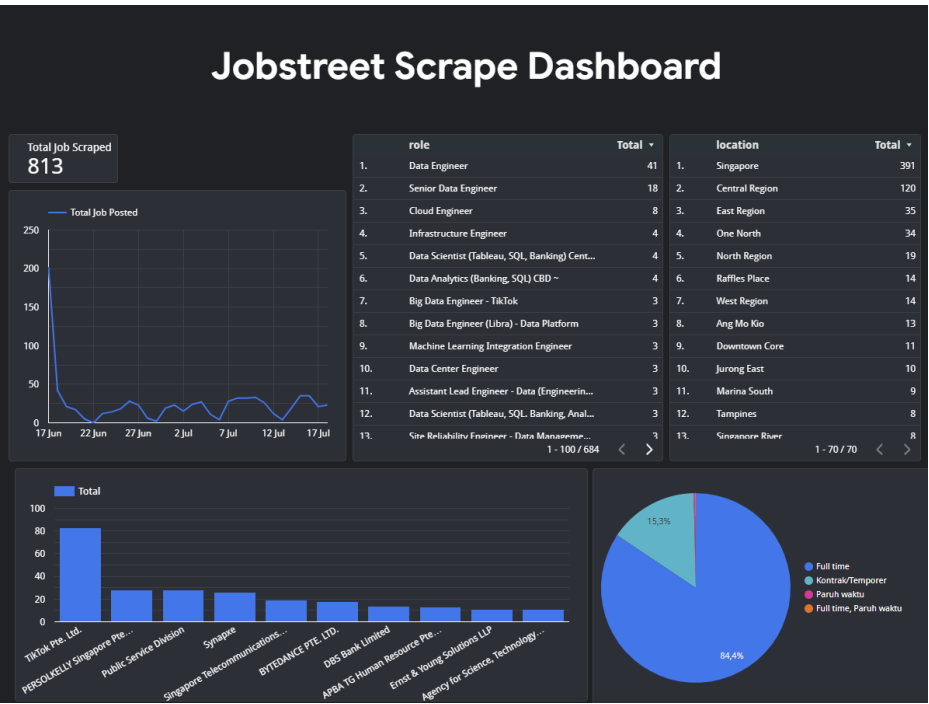jobstreetscrape_singapore_detailed table utilized for :
1. Daily dashboard such as count job posting each day, count location, etc
2. Job updated, thus the work searcher could prepare more for the job

# Dashboard

## 1. Word Frequency Count
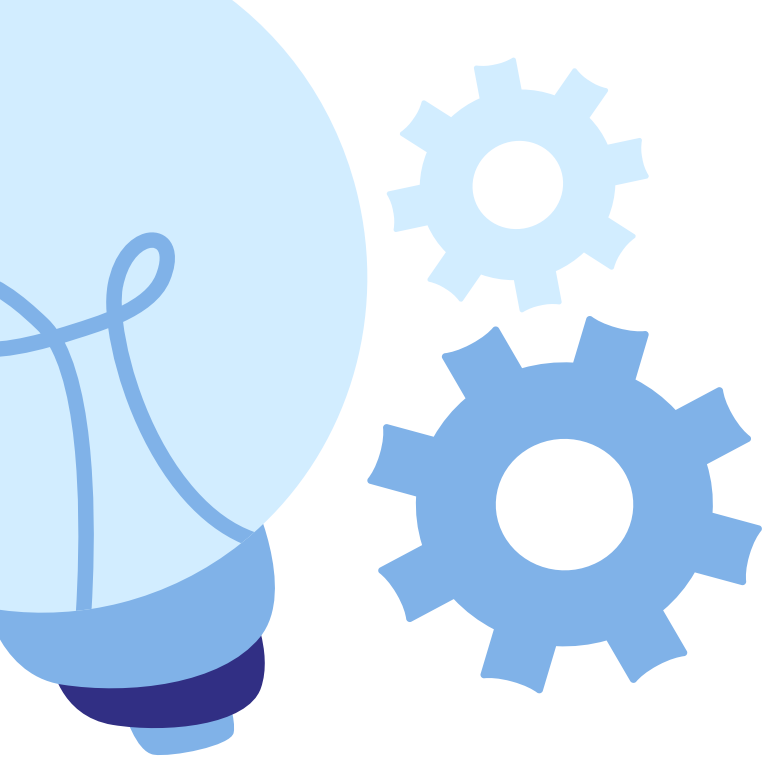


## 2. Dashboard



## 3. Job List

# Conclusion

This project could fulfill the background problem to enhance the public sector, private sector, or third sector to have curriculum aligned to the market demand, moreover the dashboard result or the data could be utilize by job seeker to prepare of their best for every job posted.

The showed pipeline run successfully which consist:

1. Built a scalable, modular ELT pipeline using Airflow with main DAG triggering child DAGs.

2. Enabled dynamic configuration and support for full & incremental loads via parameters.

3. Performed data transformation, validation, and structured into bronze, silver, and gold layers.

4. Implemented multithreaded web scraping from JobStreet for efficient data extraction.
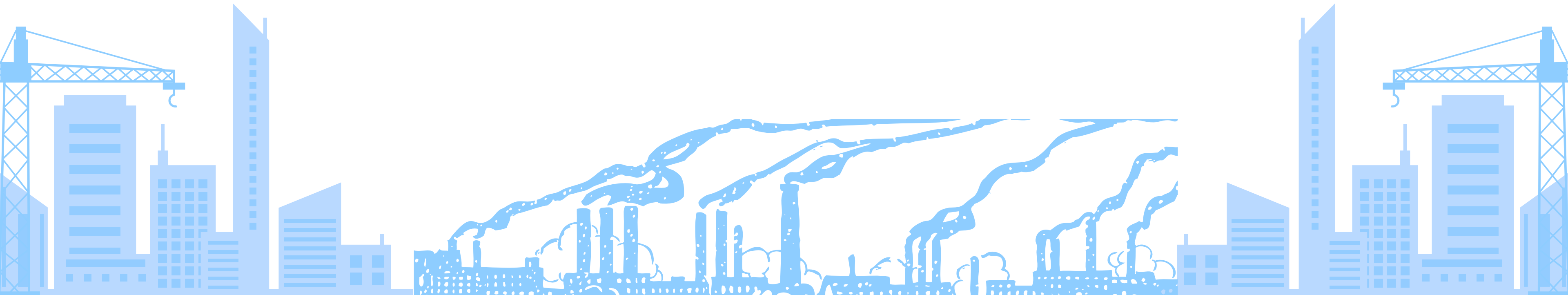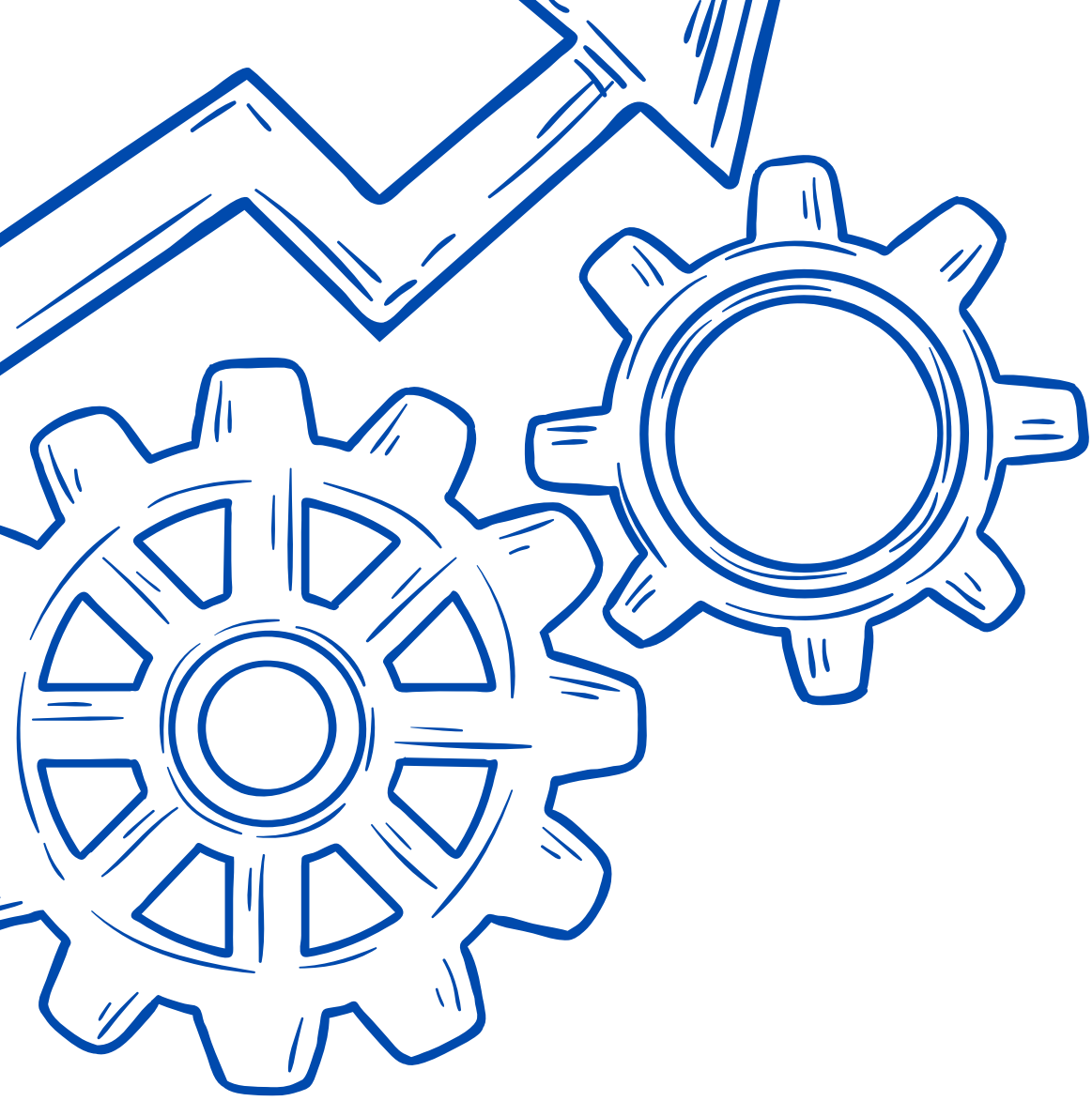
# Recommendation

The future improvement for this project could be add more the scrape data source not only jobstreet, but could be glints, dealls, etc.
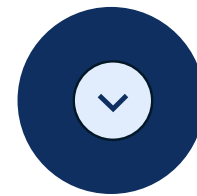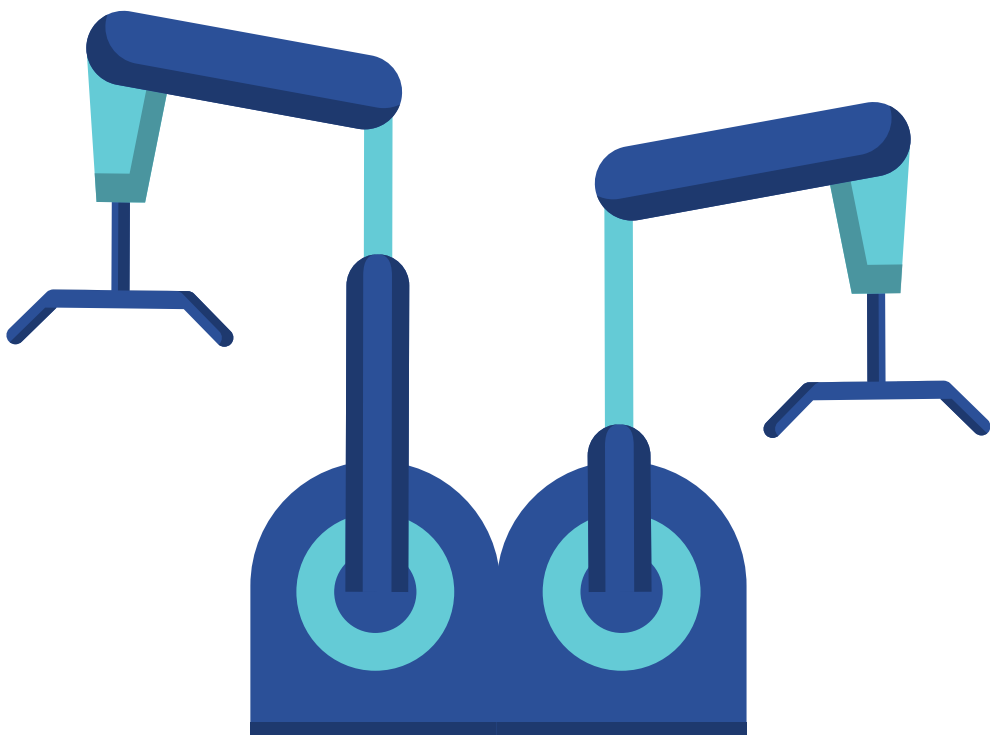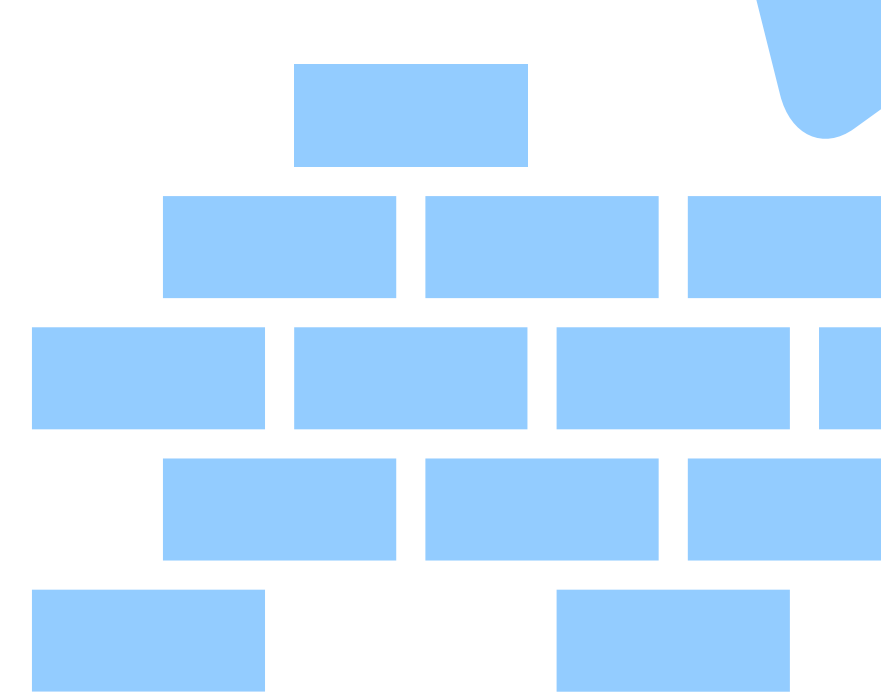
For error handling, could append retry at the DAG or task also notification if the dag is successful run or not.

At report section, could append forecast job based on posted time, moreover at the word frequency could be added a word filter more.

**Thank You**
**谢谢**
**Terima Kasih**
**감사합니다**