

Automated Data Pipeline for Sentiment Analysis Using Apache Airflow

SIB A-INFOTECH





Team Member



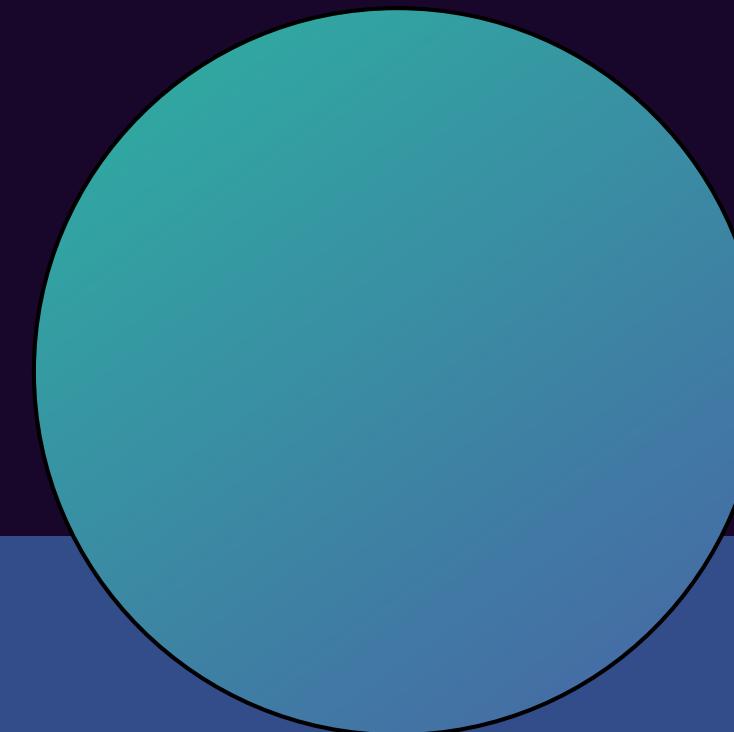
Muhamad Fahrurrozi



Muhammad Mahda



Nareswari Huwaida N



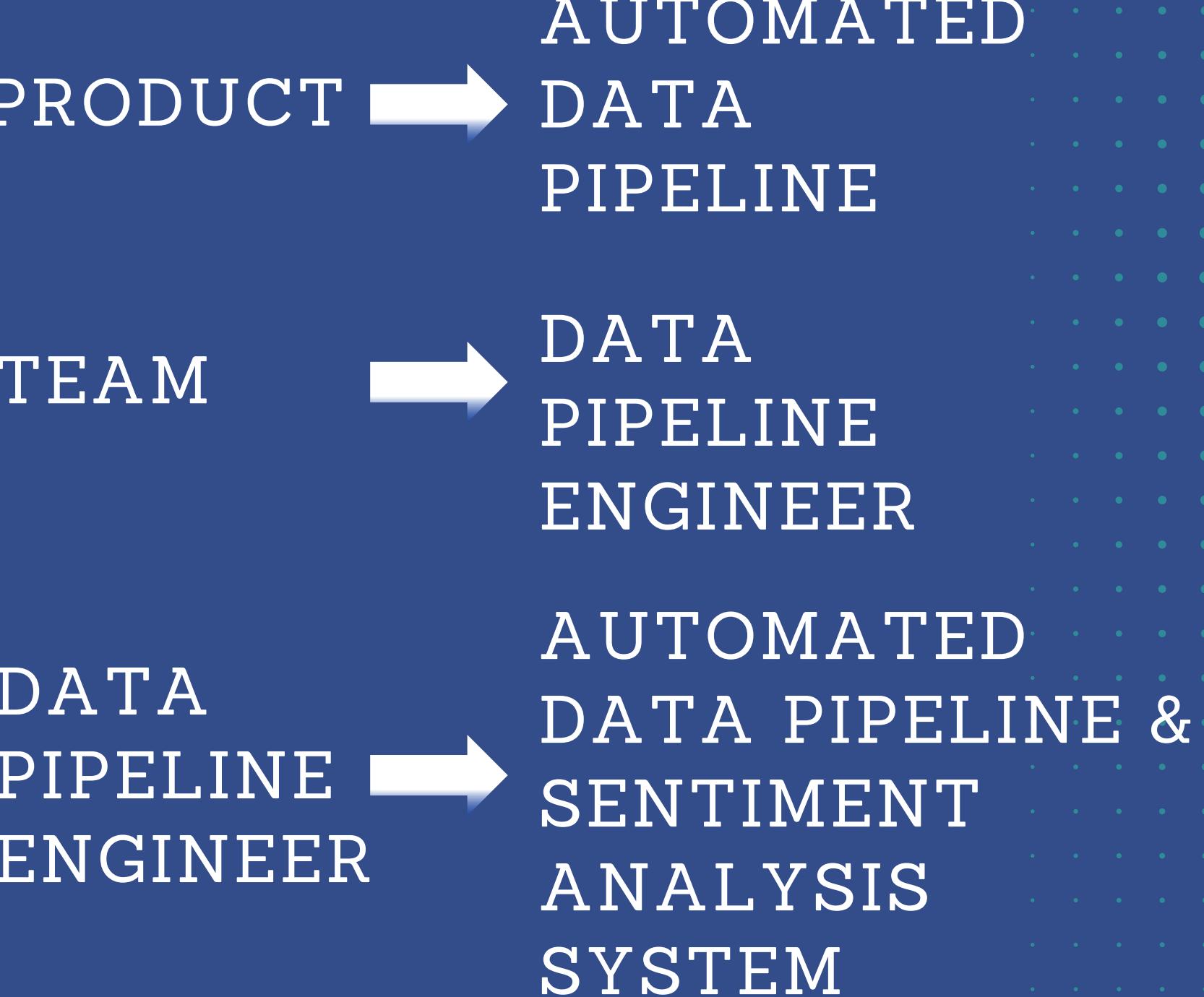
**Muhsyi Fadhilah
Hirani**

CHAPTER 1

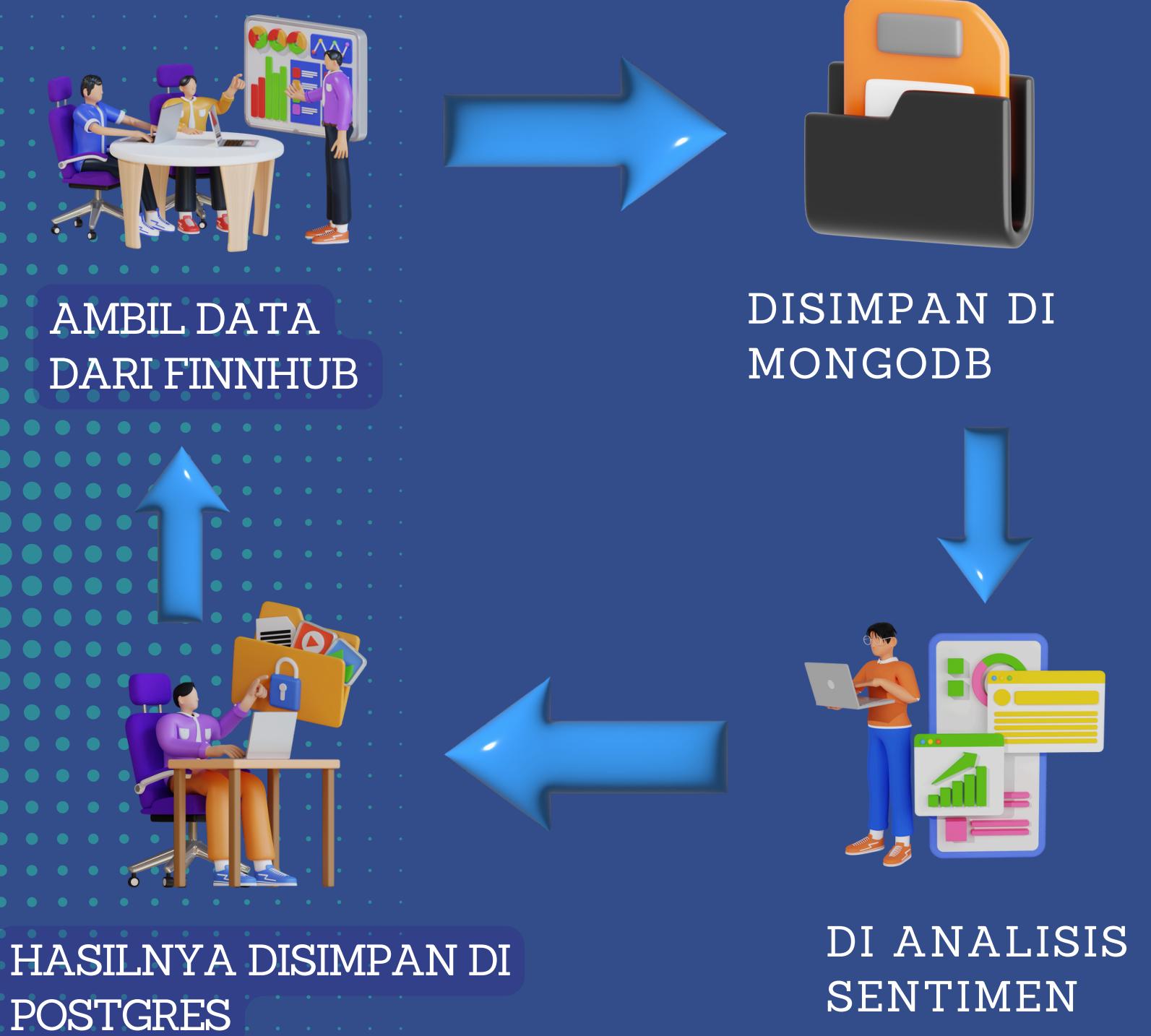
Project Overview



Project Overview



BUSSINES PROCESS



DATASET TRANSFER FINNHUB KE ATLAS

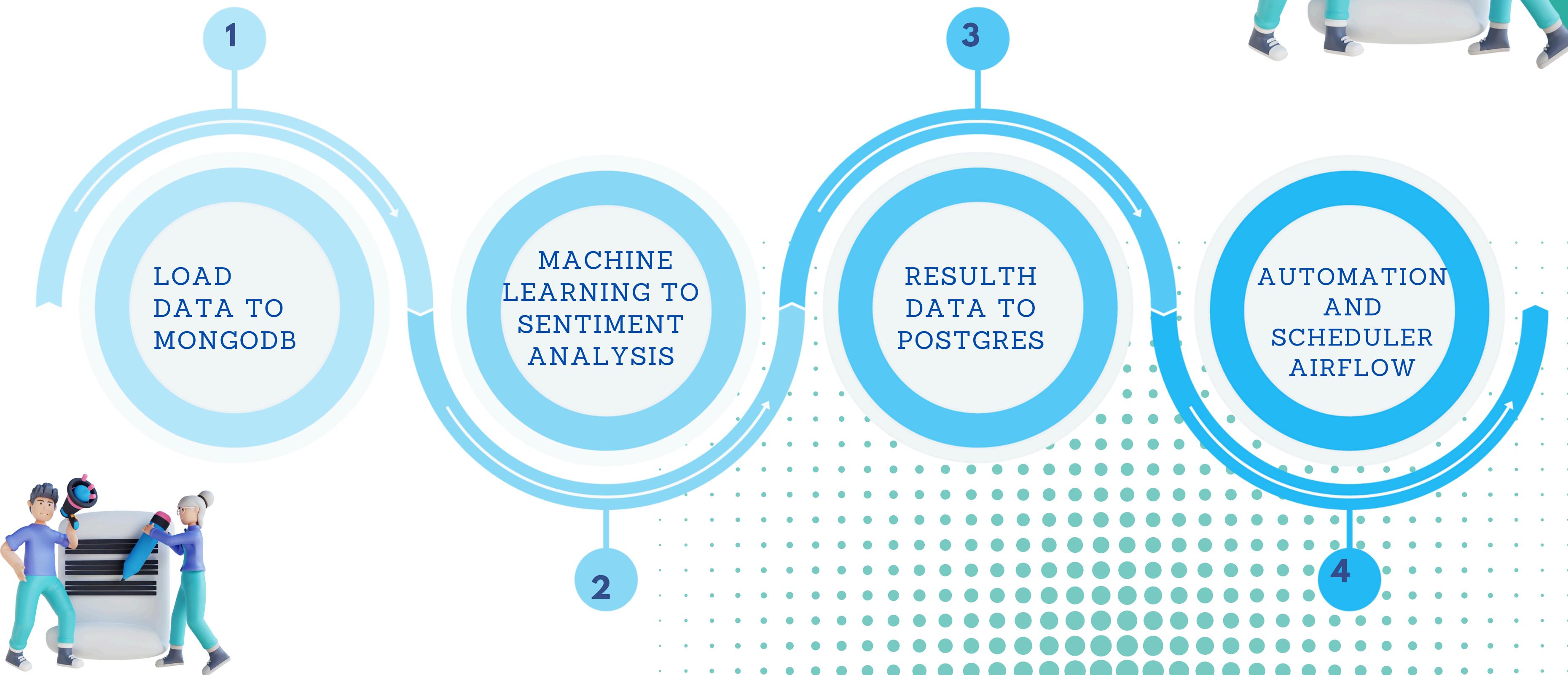


AUTOMATED DATA PIPELINE MENGGUNAKAN AIRFLOW



SENTIMENT ANALYSIS + TRANSFER ATLAS KE POSTGRES

WORKFLOW





PROBLEM

Proses manual dalam pengumpulan dan penyimpanan data saat ini tidak hanya kurang efisien tetapi juga memakan waktu yang lama, menyebabkan keterlambatan dalam analisis dan respon terhadap informasi pasar yang cepat berubah.



GOALS

Menghemat waktu dan
Meminimalisir human error
karena proses pengumpulan
data analisis data dilakukan
secara otomatis



OBJECTIVE

Membangun sistem
otomatis yang dapat
mengumpulkan,
menganalisis sentimen, dan
menyimpan data berita
untuk memberikan
wawasan pasar



BUSINESS METRICS

- Ketepatan waktu
- Responsivitas Data

CHAPTER 2

Implementation



Step-by-Step Process



Finnhub Data Extraction

Menggunakan Finnhub API untuk mengambil data berita



Data Loading

Menyimpan data mentah ke dalam MongoDB



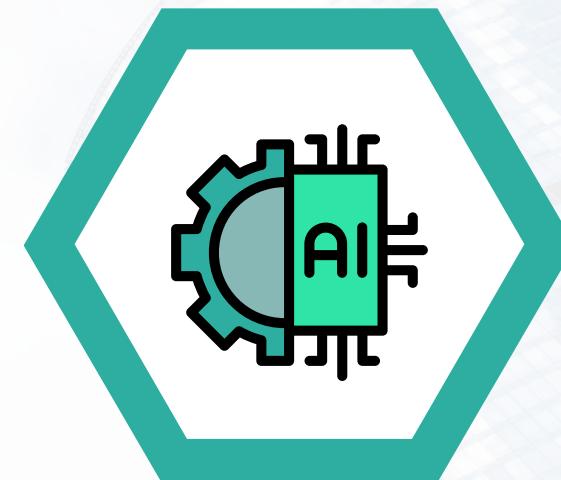
Sentiment Analysis

Menganalisis sentimen dari data tersebut



Result Storage

Menyimpan hasil analisis sentimen ke dalam PostgreSQL



Automation

Menjadwalkan tugas dengan Apache Airflow

Technologies Used

- 1 Docker Compose
- 2 Python Scripts
- 3 Apache Airflow
- 4 MongoDB & PostgreSQL



KONFIGURASI DOCKER COMPOSE

Postgres

Menggunakan postgres:12,
pengaturan pengguna dan database

airflow-init

Inisialisasi DB dan pembuatan
pengguna admin

airflow-webserver & airflow-scheduler

Instal requirements dan
menjalankan layanan terkait

Volumes & Networks

Direktori mounting dan pengaturan
jaringan

CODE SNIPPETS

Highlight key parts of the code



```
1 &airflow-common
2   image: apache/airflow:2.8.4-python3.10
3   environment:
4     - AIRFLOW__CORE__EXECUTOR=LocalExecutor
5     - AIRFLOW__CORE__SQLALCHEMY_CONN=postgresql+psycopg2://airflow:airflow@postgres:5432/airflow
6     - AIRFLOW__CORE__FERNET_KEY=FB0o_zt4e3Ziq3LdUU07F2Z95cvFFx16hU8jTeR1ASM=
7     - AIRFLOW__CORE__LOAD_EXAMPLES=False
8     - AIRFLOW__CORE__LOGGING_LEVEL=INFO
9   env_file:
10    - ./.env
11   user: "${AIRFLOW_UID:-50000}:0"
12   volumes:
13     - ./dags:/opt/airflow/dags
14     - ./source:/opt/airflow/source
15     - ./output:/opt/airflow/output
16     - ./logs:/opt/airflow/logs
17     - ./plugins:/opt/airflow/plugins
18     - ./requirements.txt:/requirements.txt

1 services:
2   postgres:
3     image: postgres:12
4     environment:
5       - POSTGRES_USER=airflow
6       - POSTGRES_PASSWORD=airflow
7       - POSTGRES_DB=airflow
8     volumes:
9       - postgres-db-volume:/var/lib/postgresql/data
10    healthcheck:
11      test: [ "CMD", "pg_isready", "-U", "airflow" ]
12      interval: 5s
13      retries: 5
14    ports:
15      - "5432:5432"
16    networks:
17      - airflow_network

1 airflow-init:
2   <<: *airflow-common
3   container_name: airflow_init
4   entrypoint: /bin/bash
5   command:
6     - -c
7     - >
8       airflow db init &&
9       airflow users create --role Admin
10      --username airflow
11      --password airflow
12      --email airflow@airflow.com
13      --firstname airflow
14      --lastname airflow
15   restart: on-failure
16   depends_on:
17     postgres:
18       condition: service_healthy
19   networks:
20     - airflow_network
```

CODE SNIPPETS

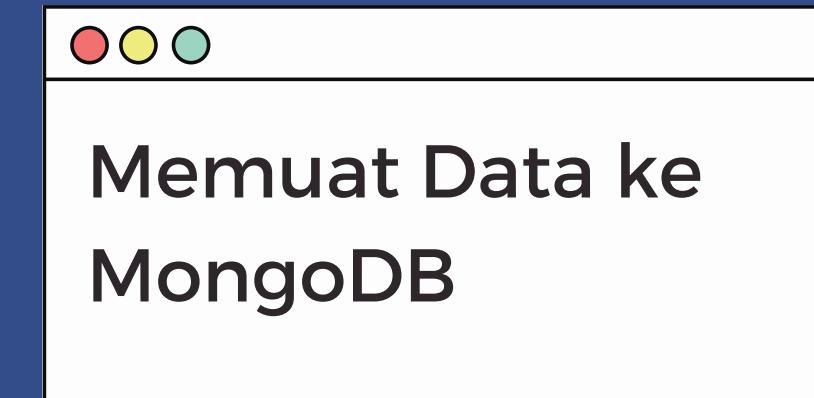


Highlight key parts of the code

```
● ● ●  
1 airflow-webserver:  
2   <<: *airflow-common  
3   command: >  
4     bash -c "python -m pip install --upgrade pip && pip install --no-cache-dir -r /requirements.txt && airflow webserver"  
5   ports:  
6     - 8080:8080  
7   container_name: airflow_webserver  
8   depends_on:  
9     airflow-init:  
10       condition: service_completed_successfully  
11   restart: always  
12   networks:  
13     - airflow_network  
14  
15 airflow-scheduler:  
16   <<: *airflow-common  
17   command: >  
18     bash -c "python -m pip install --upgrade pip && pip install --no-cache-dir -r /requirements.txt && airflow scheduler"  
19   container_name: airflow_scheduler  
20   depends_on:  
21     airflow-init:  
22       condition: service_completed_successfully  
23   restart: always  
24   networks:  
25     - airflow_network
```

Ringkasan Proses

Mengumpulkan dan Menyimpan Data ke MongoDB



Analisis Sentimen dan load ke PostgreSQL

Inisialisasi Koneksi

Analisis Sentimen

Mengonversi ke DataFrame

Load DF ke PostgresSQL

Code Snippets

```
● ● ●  
1 def get_mongo_client(mongo_uri):  
2     try:  
3         client = pymongo.MongoClient(mongo_uri)  
4         client.admin.command('ping')  
5  
6         print("Successfully connected to MongoDB!")  
7         return client  
8  
9     except pymongo.errors.ConnectionFailure as e:  
10        print(f"Connection failed: {e}")  
11  
12    return None
```

Koneksi
MongoDB

● ● ●

Pengambilan Data

```
1 def scrape_news():  
2     api_key = os.getenv("FINNHUB_API_KEY")  
3     finnhub_client = finnhub.Client(api_key=api_key)  
4     news = finnhub_client.general_news('general', min_id=0)  
5     return news
```

● ● ●

Memuat Data ke MongoDB

```
1 def extract_load():  
2     news = finnhub_loader.scrape_news()  
3     collection = mongodb_loader.load('news', 'finnhub_news')  
4     collection.insert_many(news)  
5     print("Successfully load news to MongoDB")
```

Code Snippets

● ● ●

Analisis Sentimen

```
1 class SentimentAnalysis:  
2     def execute(self):  
3         analysis = TextBlob(self.text)  
4         if analysis.sentiment.polarity > 0:  
5             data = {'text': self.text, 'sentiment': 'positive'}  
6         elif analysis.sentiment.polarity == 0:  
7             data = {'text': self.text, 'sentiment': 'neutral'}  
8         else:  
9             data = {'text': self.text, 'sentiment': 'negative'}  
10        return data  
11
```

● ● ●

● ● ●

Memuat Data ke PostgreSQL

```
1 def load(data, table_name):  
2     conn_string = 'postgresql://airflow:airflow@postgres:5432/data_warehouse'  
3     engine = create_engine(conn_string)  
4     data.to_sql(table_name, con=engine, if_exists='append', index=False)  
5     print("Successfully loaded to Postgres")
```

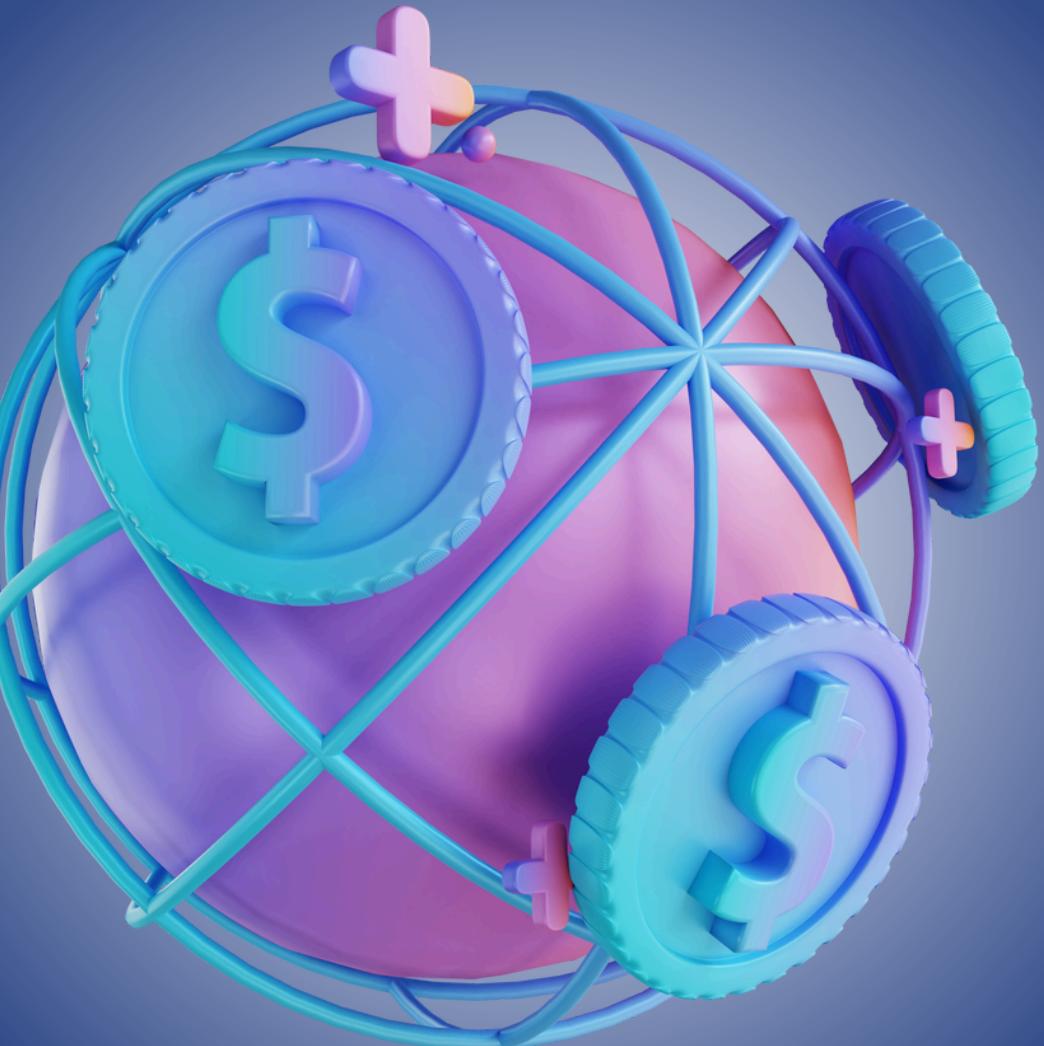
● ● ●

```
1 def run_analysis():  
2     mongo_uri = os.getenv("MONGO_URI")  
3     db = get_mongo_client(mongo_uri)  
4     news_collection = db['news']['finnhub_news']  
5  
6     news = [x for x in news_collection.find()]
```

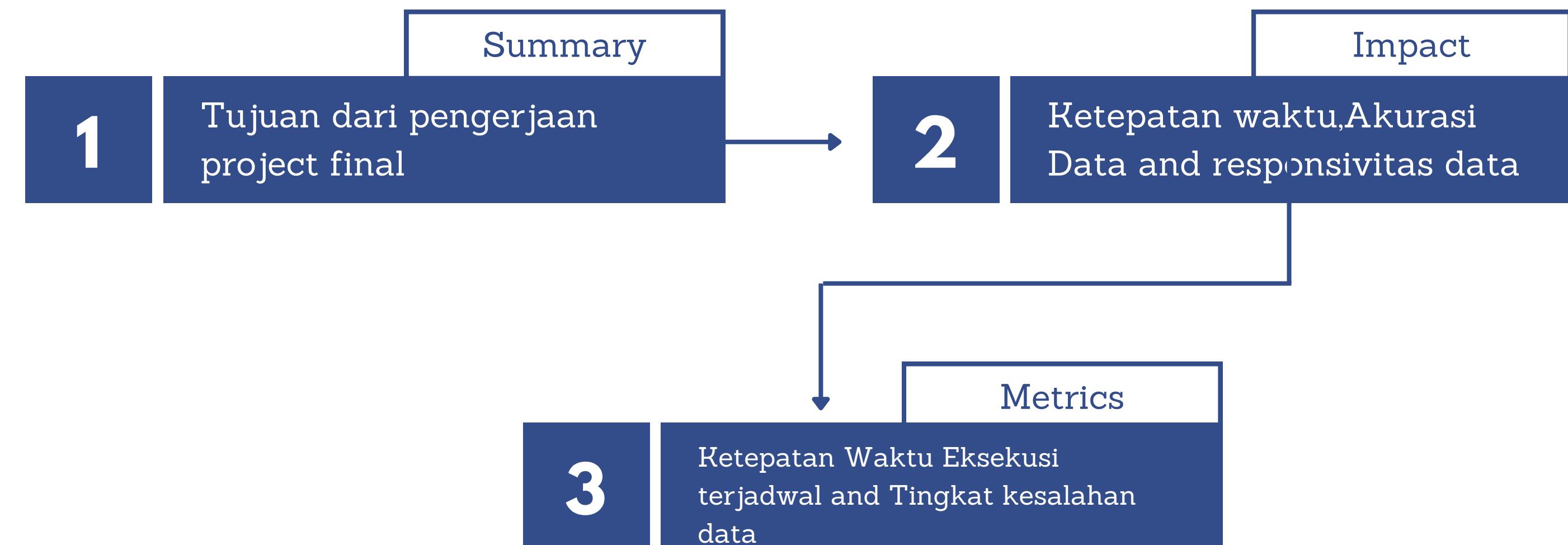
Inisialisasi Koneksi MongoDB dan Pengambilan Data

CHAPTER 3

Results



Summary, Impact and Metrics



Summary

bertujuan untuk mengotomatisasi proses pengumpulan dan analisis data sentiment dari Finnhub dengan menggunakan pipeline Airflow, Pipeline tersebut mengintegrasikan beberapa komponen penting.



Finnhub
pengambilan
data



PostgreSQL
penyimpanan
data akhir

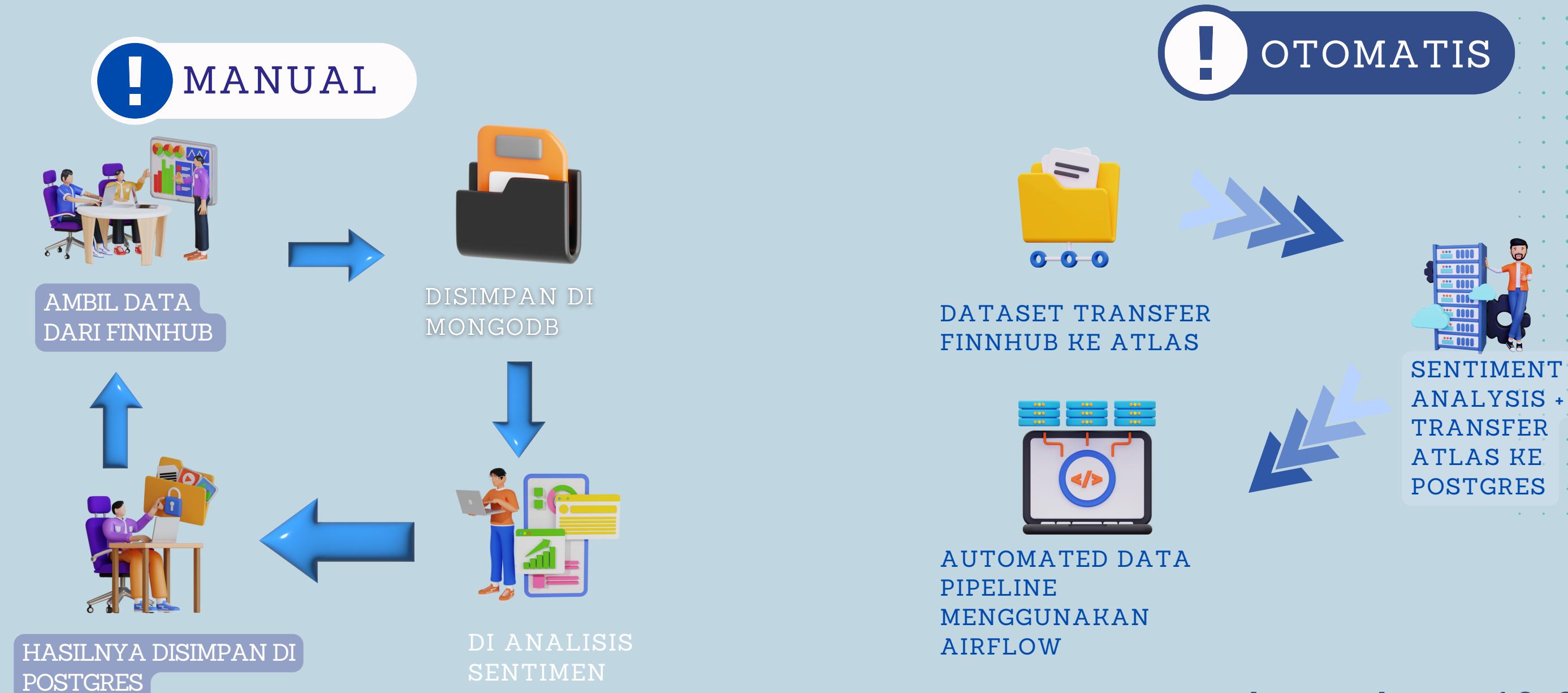


MongoDB
penyimpanan
data sementara



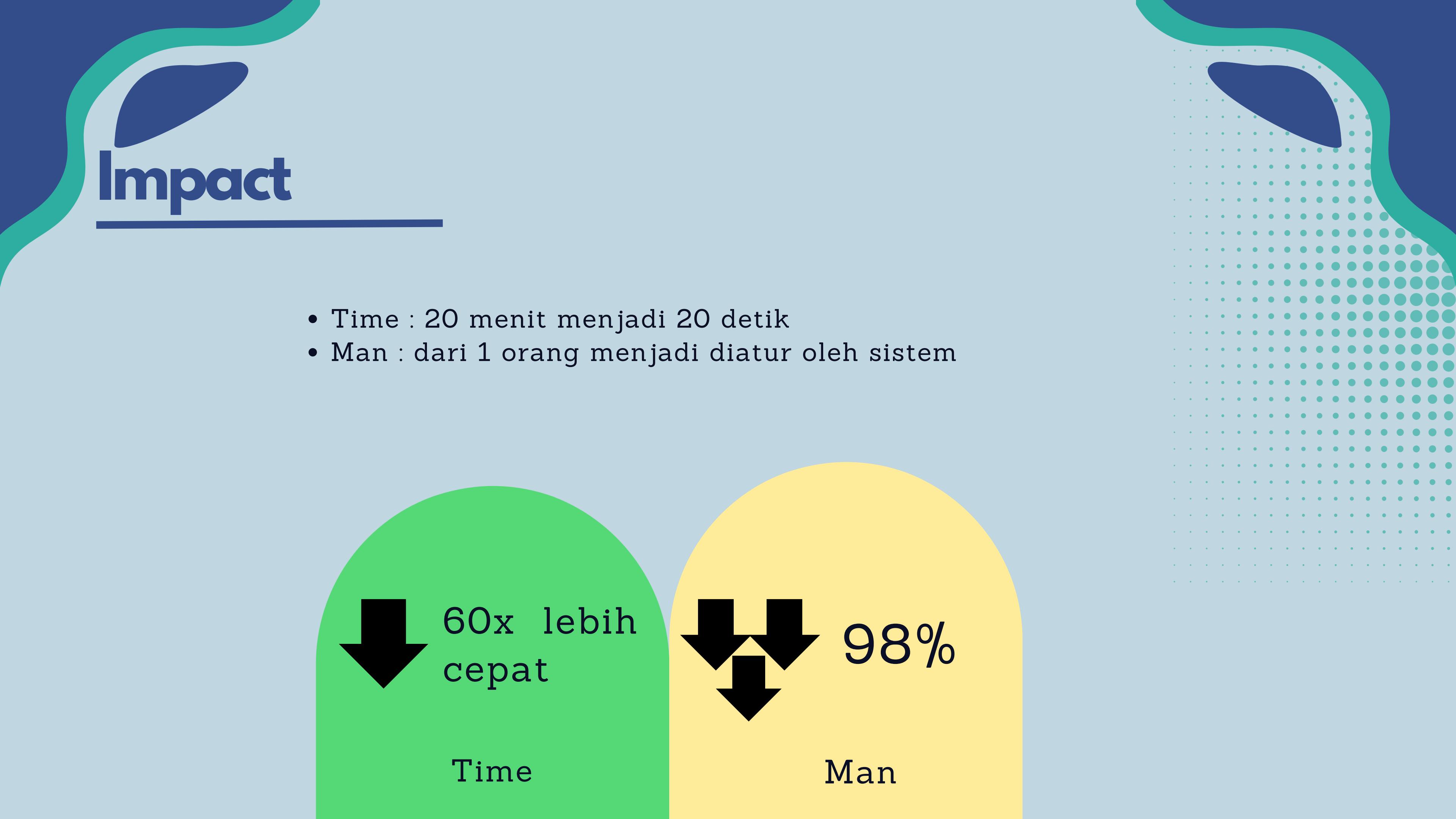
TextBlob
analisis
sentiment

Impact



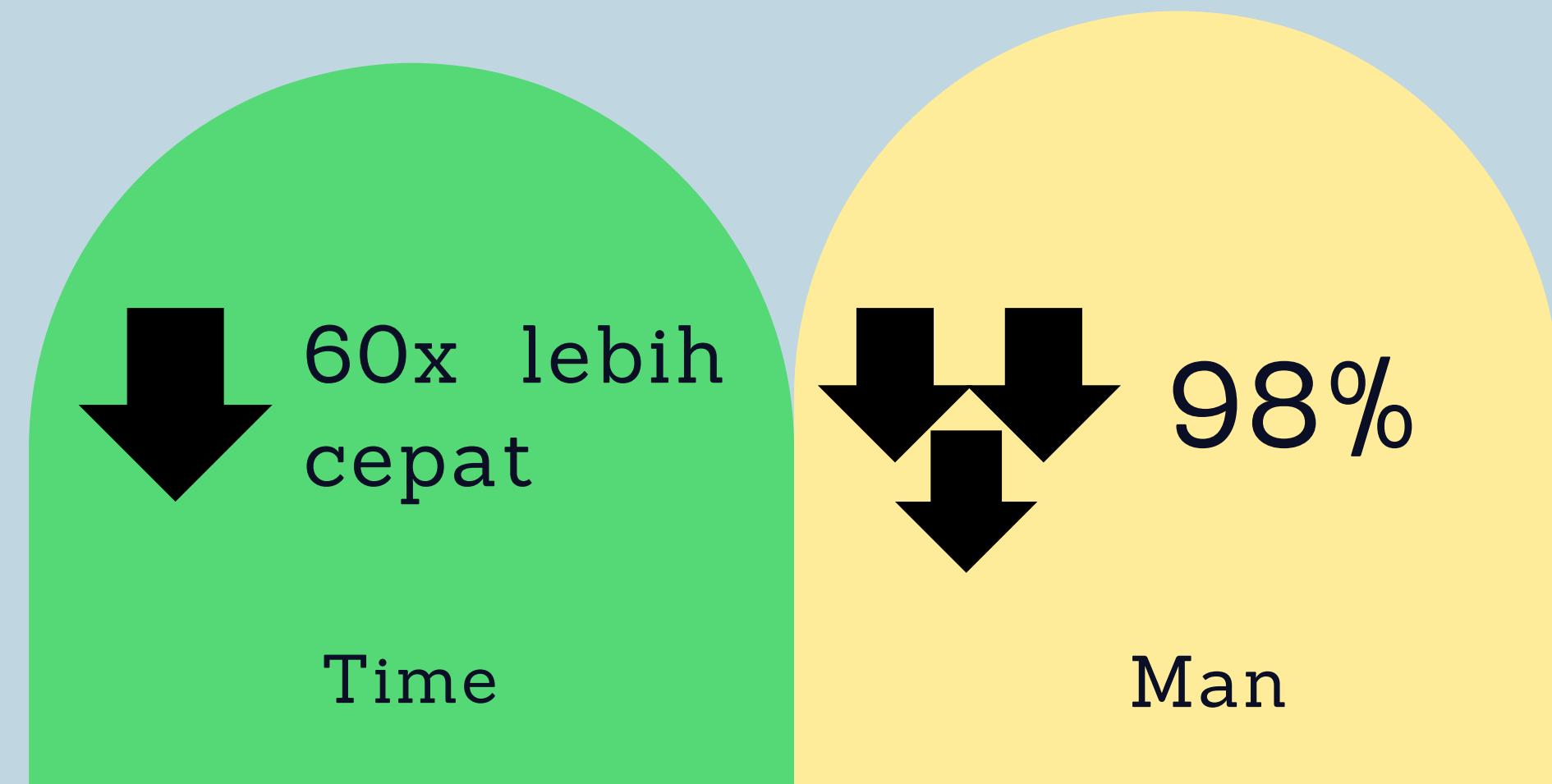
penggunaan waktu sekitar 20-40 menit untuk analysis secara manual oleh individu

penggunaan waktu sekitar 10-30 detik untuk menganalisis sentiment, dengan penjadwalan airflow



Impact

- Time : 20 menit menjadi 20 detik
- Man : dari 1 orang menjadi diatur oleh sistem



Impact



Efisiensi Waktu



Akurasi Data



Responsivitas



Metrics

Implementasi project menghasilkan metrics bussines, dimana pada project ini mencapai tujuan yaitu responsivitas data dan juga ketepatan waktu

Metrics



Ketepatan Waktu Eksekusi Terjadwal (Scheduled Execution Time Accuracy)

- Secara manual : Dari 30 eksekusi yang dijadwalkan dalam satu bulan, hanya 20 yang dijalankan tepat waktu (65.5%).
- Secara otomatis : Dari 30 eksekusi yang dijadwalkan dalam setiap hari, dijalankan tepat waktu pada jam 00.00.

Keterlambatan dalam Pipeline Data (Latency in Data Pipeline)

- Secara manual : Rata-rata waktu pemrosesan adalah 20 - 30 menit.
- Secara otomatis : Rata-rata waktu pemrosesan berkurang menjadi 20 - 30 detik

Result for Project



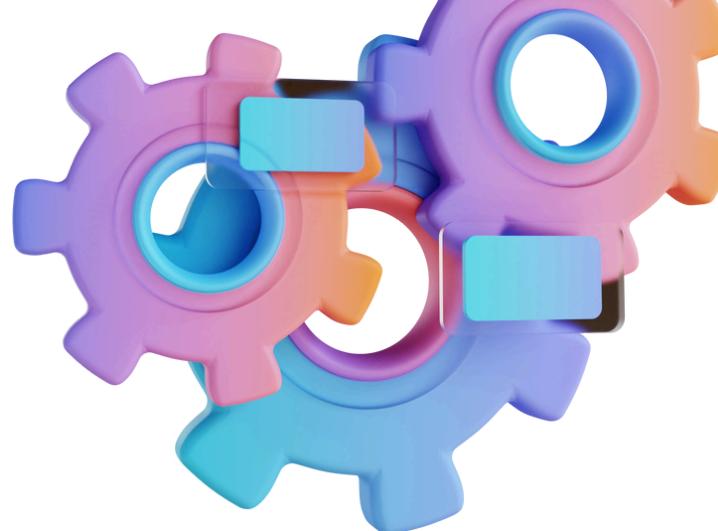
Connect in mongoDB

```
mahda@LAPTOP-6QU2AGRG:/mnt/c/Users/HP/STUPEN/Airflow_Data_Pipeline_for_Sentiment_Analysis/plugins$ python3 finnhub_mongodb_loader.py
/home/mahda/.local/lib/python3.8/site-packages/pymongo/network.py:36: UserWarning: Failed to use the installed version of PyOpenSSL. Falling back to stdlib ssl, disabling OCSP support. This is likely caused by incompatible versions of PyOpenSSL < 23.2.0 and cryptography >= 42.0.0. Try updating PyOpenSSL >= 23.2.0 to enable OCSP.
  from pymongo import _csot, helpers, message, ssl_support
Pinged your deployment. You successfully connected to MongoDB!
Successfully load news to MongoDB
```

Connected in Postgres

```
mahda@LAPTOP-6QU2AGRG:/mnt/c/Users/HP/STUPEN/Airflow_Data_Pipeline_for_Sentiment_Analysis/plugins$ python3 sentiment_analysis_loader.py
/home/mahda/.local/lib/python3.8/site-packages/pymongo/network.py:36: UserWarning: Failed to use the installed version of PyOpenSSL. Falling back to stdlib ssl, disabling OCSP support. This is likely caused by incompatible versions of PyOpenSSL < 23.2.0 and cryptography >= 42.0.0. Try updating PyOpenSSL >= 23.2.0 to enable OCSP.
  from pymongo import _csot, helpers, message, ssl_support
Pinged your deployment. You successfully connected to MongoDB!
■ Summary 'Vanguard data casts a shadow on the economy' successfully analyzed
Summary 'The Federal Reserve should be able to cut interest rates with inflation coming down and the labor market cooling, but the timing of any move is uncertain, Federal Reserve governor Lisa Cook said Tuesday.' successfully analyzed
Berhasil dimuat ke PostgreSQL
Successfully loaded to Postgres
```

Result for Project



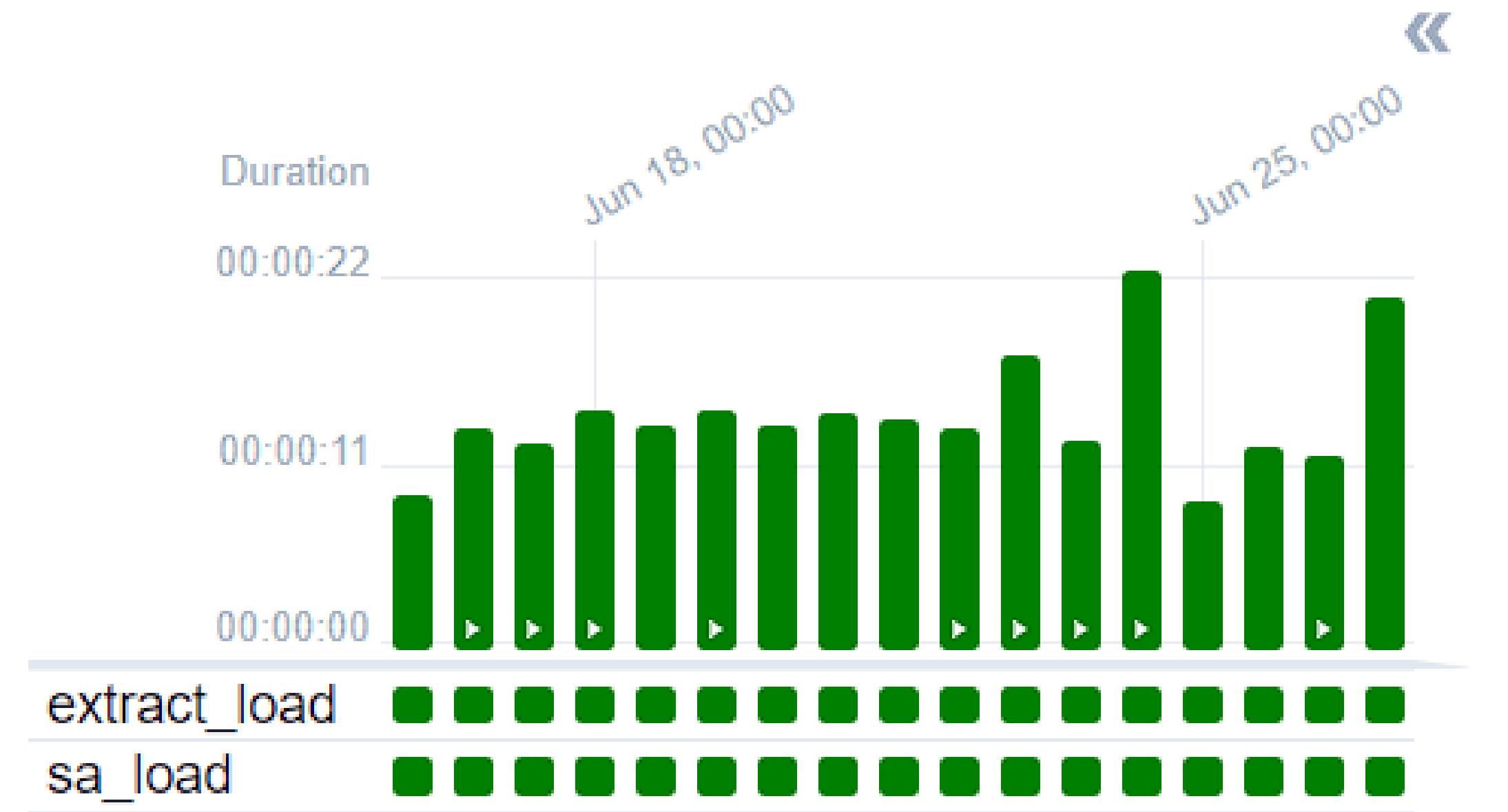
Connect in data_warehouse

```
data_warehouse=# \c data_warehouse
You are now connected to database "data_warehouse" as user "airflow".
data_warehouse=# \d
      List of relations
 Schema |           Name           | Type  | Owner
-----+-----+-----+-----+
 public | sentiment_news_analysis | table | airflow
(1 row)
```

displays data in the data_warehouse table

		text
		sentiment
	-----+-----+	
Each team gets about \$10 million for taking part in the European Championships – but none of that goes directly to the players.		positive
Have a clear road map for your estate so your heirs don't drown in paperwork and probate after your death		positive
Roaring Kitty weighed in on the event with fresh cryptic tweets.		positive
Oil futures edged higher Tuesday, on track to notch back-to-back gains.		positive

Result in otomatisasi



implementation
Airflow

Result for Project

displays in MongoDB atlas

The screenshot shows the MongoDB Atlas interface for a project named "Project 0". The left sidebar includes sections for Overview, Deployment (selected), Database (selected), Services, and Security. The main area displays the "news.finnhub_news" database with a storage size of 1.34MB and 4300 documents. A search bar and a "Find" button are available. Below the database name, it says "Generate queries from natural language in Compass". A query input field contains the placeholder "Type a query: { field: 'value' }". The results section shows two news documents:

```
_id: ObjectId('6671a9e014dbecc84bd71edd')
category : "top news"
datetime : 1718724300
headline : "Euro 2024 prize money: How much do the winners get?"
id : 7376435
image : "https://static2.finnhub.io/file/publicdatany/finnhubimage/market_watch_..."
related : ""
source : "MarketWatch"
summary : "Each team gets about $10 million for taking part in the European Champ..."
url : "https://www.marketwatch.com/story/euro-2024-prize-money-how-much-do-th..."
```



```
_id: ObjectId('6671a9e014dbecc84bd71ede')
category : "top news"
datetime : 1718723280
headline : "Your family will love you even more if you simplify your estate"
id : 7376432
```

Navigation buttons include "PREVIOUS", "NEXT", and "1-20 of many results". The bottom footer includes links for System Status, Status, Terms, Privacy, Atlas Blog, Contact Sales, and a copyright notice for 2024 MongoDB, Inc.

CHAPTER 4

Conclusion



Implementasi pipeline otomatis telah berhasil mencapai tujuan yang ditetapkan, yaitu menghemat waktu, meminimalisir kesalahan manusia, meningkatkan ketepatan waktu, dan meningkatkan responsivitas data. Dengan hasil ini, perusahaan dapat memperoleh insight yang lebih cepat dan akurat, memungkinkan pengambilan keputusan yang lebih baik dan lebih responsif terhadap perubahan pasar.



CHAPTER 5

Discussion



A 3D cartoon character with brown hair and a blue jacket, waving with one hand.

INFOTECH Thank You

LET'S DISCUSS!