



Kampus
Merdeka
INDONESIA JAYA

MSIB



TEAM 1

FINAL PROJECT

Presentation

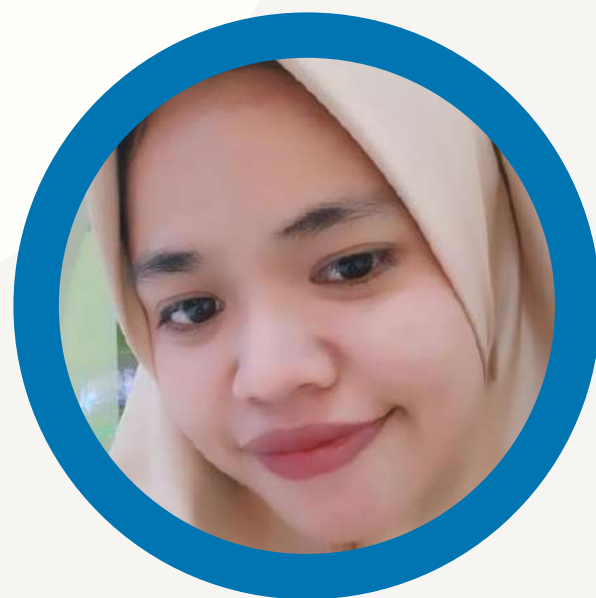
STUDI INDEPENDEN KAMPUS MERDEKA – Batch 6

Data Warehousing



HELLO!

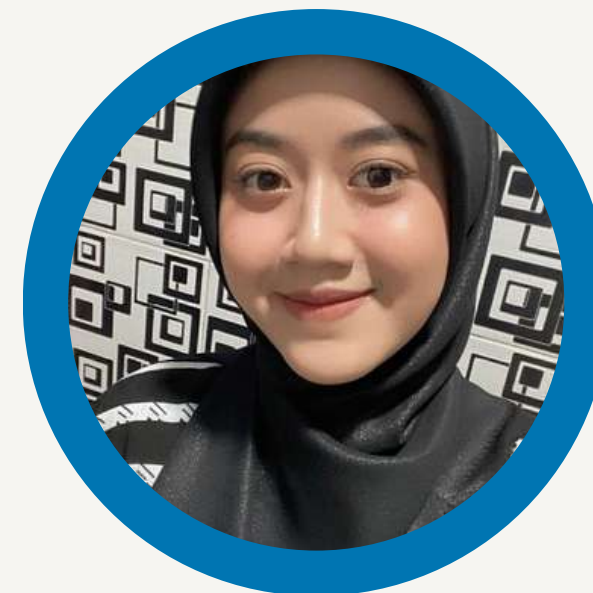
MEET OUR TEAM



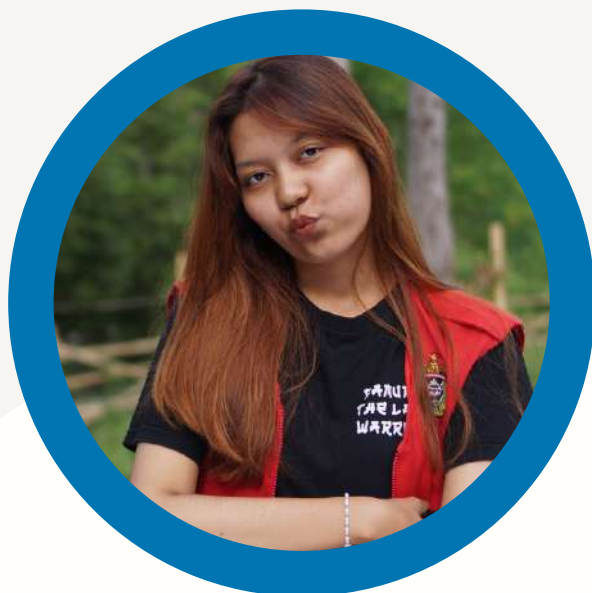
ERNAWATI



Seftia Nur'sukmawati



**Chika
Cardealfitri**



**Clarisa Eudia
Chesynanda**



Suci Rakhmayanti



Presentation



FINAL PROJECT



01.

INTRODUCTION

Penjelasan latar belakang masalah yang terjadi dalam membuat ETL, Data warehouse, Data Visualisasi.

02.

ETL

Extract, transform, and load (ETL) adalah proses menggabungkan data dari berbagai sumber ke dalam repositori pusat yang besar yang disebut gudang data.

03.

DATA MODELLING

Data modelling adalah proses mengumpulkan data dan mengubahnya menjadi diagram sederhana agar bisa menjadi informasi yang bisa digunakan sesuai kebutuhan bisnis.

04.

DATA VISUALISATION

Visualisasi data adalah proses menggunakan elemen visual seperti diagram, grafik, atau peta untuk merepresentasikan data. Visualisasi data menerjemahkan yang kompleks, bervolume tinggi, atau numerik menjadi representasi visual yang lebih mudah diproses.

1.

INTRODUCTION

Background

Congratulations on your first role as Data Engineer!. You are just hired at a US online retail company that sells general customer products directly to customers from multiple suppliers around the world. Your challenge is to **build-up a data infrastructure** using generated data crafted to mirror real-world data from leading tech companies. You will be provided with template projects to guide your work.



14



Tasks

- ETL/ELT Job Creation using **Airflow**
- Data Modeling in **Postgres**
- Dashboard Creation with Data Visualization
- Craft a Presentation Based on Your Work

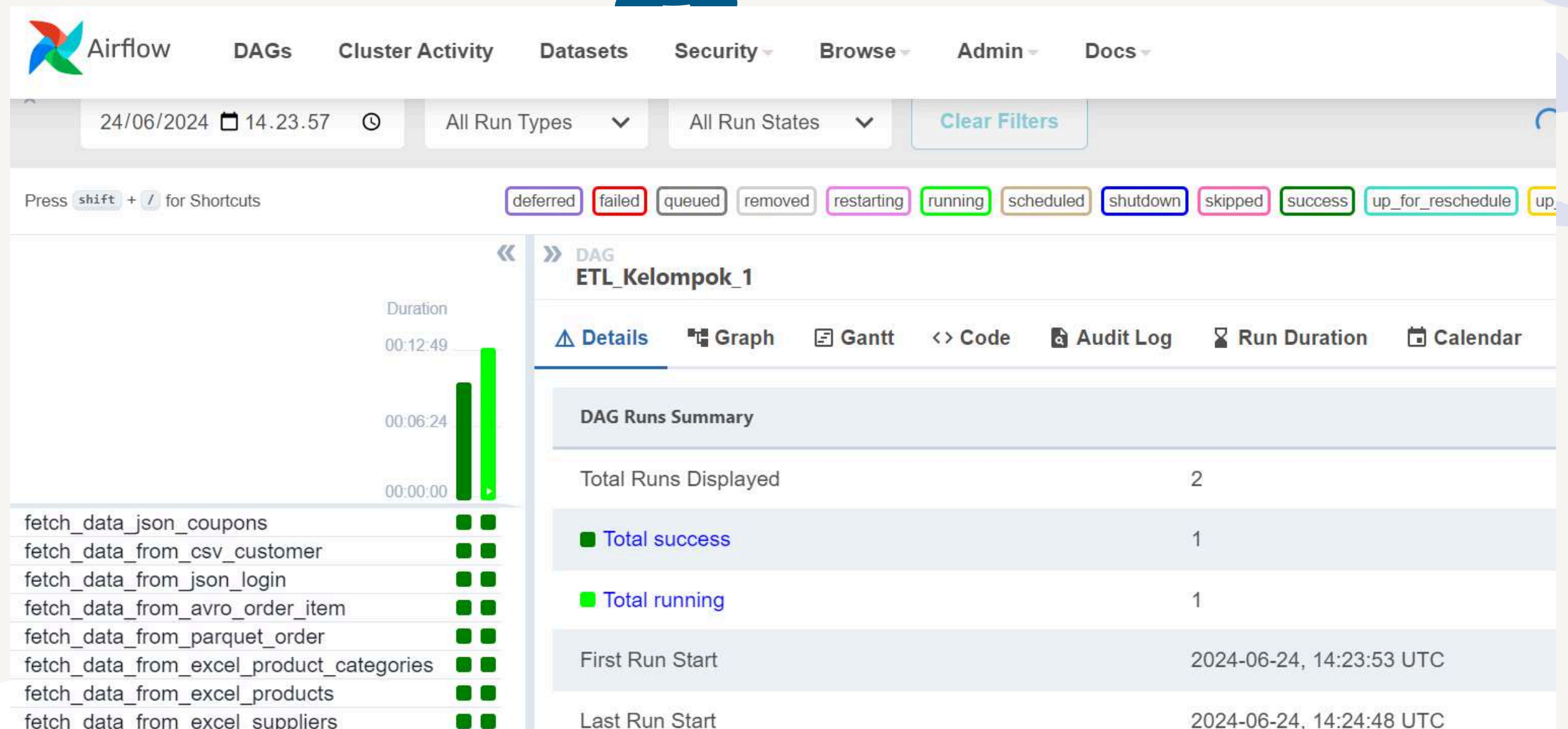
Notes: All data being used is generated data crafted to mirror real-world data from online retail companies.



2.

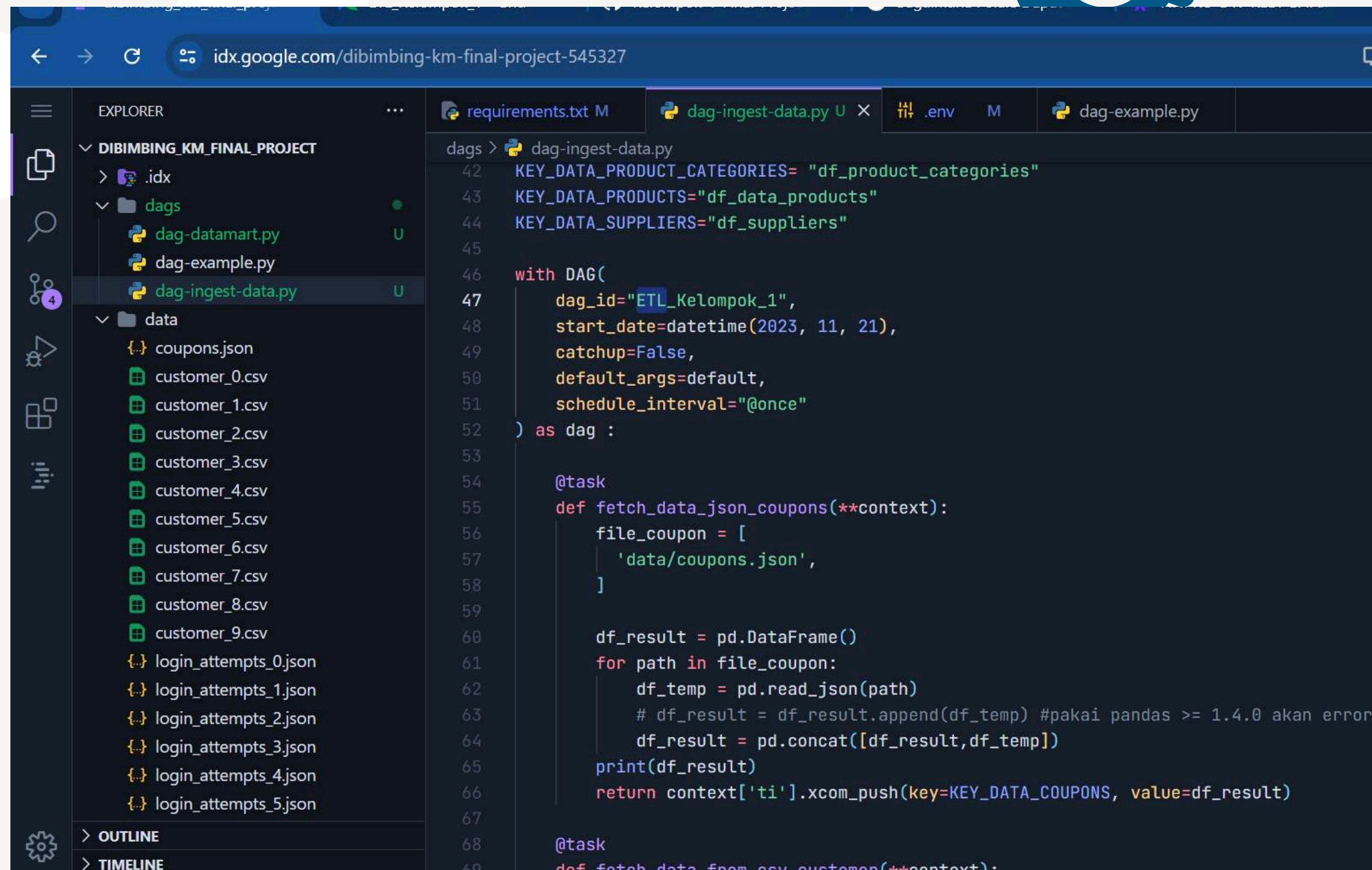
INGEST

proses pemindahan data dari satu atau beberapa sumber ke suatu penyimpanan. Data tersebut nantinya akan disimpan dan dianalisis lebih lanjut



3.

ETL



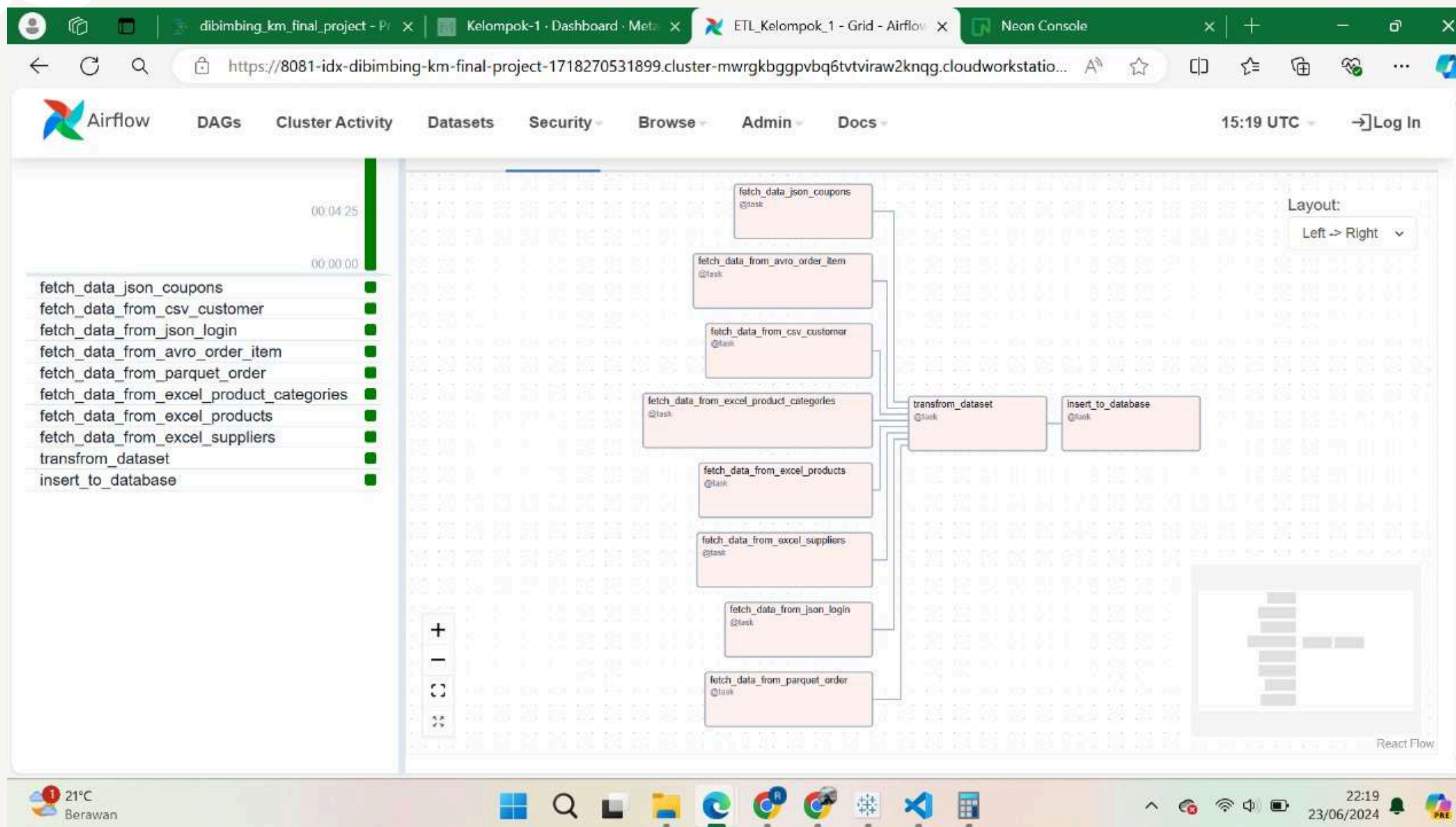
The screenshot shows a web-based IDE interface with a dark theme. The browser address bar displays `idx.google.com/dibimbing-km-final-project-545327`. The left sidebar contains an 'EXPLORER' panel showing a project structure for 'DIBIMBING_KM_FINAL_PROJECT'. It includes a '.idx' folder, a 'dags' folder with files 'dag-datamart.py', 'dag-example.py', and 'dag-ingest-data.py' (the latter is selected), and a 'data' folder containing various JSON and CSV files. The main editor area shows the code for 'dag-ingest-data.py'. The code defines a DAG with a specific ID and start date, and includes a task to fetch JSON coupon data using pandas.

```
requirements.txt M dag-ingest-data.py U .env M dag-example.py
dags > dag-ingest-data.py
42 KEY_DATA_PRODUCT_CATEGORIES= "df_product_categories"
43 KEY_DATA_PRODUCTS="df_data_products"
44 KEY_DATA_SUPPLIERS="df_suppliers"
45
46 with DAG(
47     dag_id="ETL_Kelompok_1",
48     start_date=datetime(2023, 11, 21),
49     catchup=False,
50     default_args=default,
51     schedule_interval="@once"
52 ) as dag :
53
54     @task
55     def fetch_data_json_coupons(**context):
56         file_coupon = [
57             'data/coupons.json',
58         ]
59
60         df_result = pd.DataFrame()
61         for path in file_coupon:
62             df_temp = pd.read_json(path)
63             # df_result = df_result.append(df_temp) #pakai pandas >= 1.4.0 akan error
64             df_result = pd.concat([df_result,df_temp])
65         print(df_result)
66         return context['ti'].xcom_push(key=KEY_DATA_COUPONS, value=df_result)
67
68     @task
69     def fetch_data_from_csv_customer(**context):
```

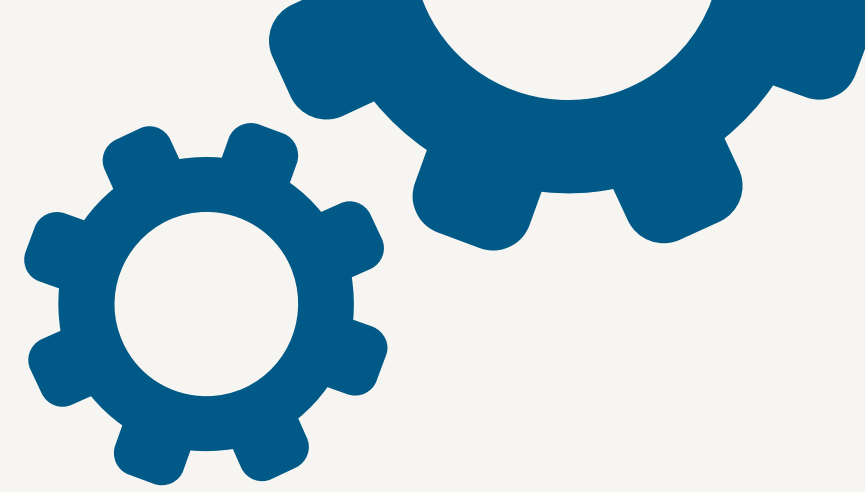
Menggunakan bahasa pemrograman python dan memakai library pandas dan polars untuk mengefektifkan ingestion data

4.

ETL



5. DATA MODELLING



A screenshot of the Apache Airflow web interface showing a DAG (Data As a Graph) named 'factdim_tabel_dag'. The interface includes a top navigation bar with tabs for DAGs, Cluster Activity, Datasets, Security, Browse, Admin, and Docs. The main content area displays the DAG's execution history on the left and a detailed graph view on the right. The graph shows a sequence of tasks, including 'create_dim_login_attempt_history', 'transfer_login_attempt_history', 'create_dim_suppliers', 'transfer_suppliers', 'create_dim_order_items', 'transfer_order_items', 'create_fact_order_details', 'transfer_order_details', 'create_dim_products', 'transfer_products', 'create_dim_orders', 'transfer_orders', 'create_dim_product_categories', 'transfer_product_categories', 'create_dim_coupons', 'transfer_coupons', and 'create_dim_customers', 'transfer_customers'. The tasks are arranged in a grid-like structure, indicating a complex dependency graph. The interface also shows a 'Layout' dropdown set to 'Left -> Right' and a 'Log In' button. The bottom status bar displays the system clock (22:20, 23/06/2024) and various system icons.

DATA MODELLING

tabel ETL yang telah berhasil di load, dilakukan transformasi dengan membentuk data modelling yang memiliki fact dan dimension table



```
dags > dag-factdim.py
16 def source_postgres_connection():
17     connection_string = URL.create(
18         'postgresql',
19         username='data_warehouse_owner',
20         password='SwEm1h0bTgWz',
21         host='ep-lively-tooth-a1rqoefa.ap-southeast-1.aws.neon.tech',
22         database='data_warehouse',
23         port=5432,
24         query={'sslmode': 'require'}
25     )
26     engine = create_engine(connection_string)
27     return engine
28
29 def target_postgres_connection():
30     connection_string = URL.create(
31         drivername='postgresql',
32         username='data_warehouse_owner',
33         password='PQJmnIdjYf02',
34         host='ep-noisy-river-a5fcgvg3.us-east-2.aws.neon.tech',
35         port=5432,
36         database='data_warehouse',
37         query={'sslmode': 'require'}
```

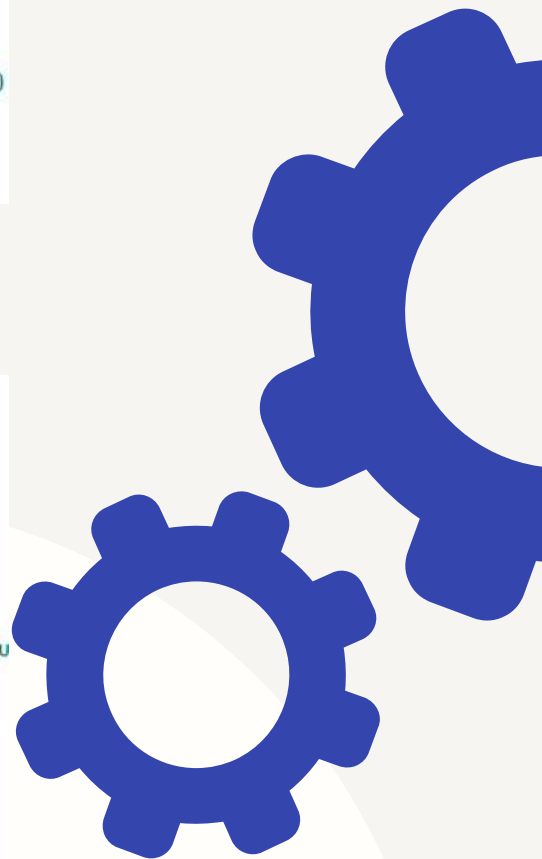
```
def create_dim_customers_table():
    logging.info("Creating dim_customers table in the target database.")
    try:
        engine = target_postgres_connection()
        with engine.connect() as connection:
            connection.execute("""
                CREATE TABLE IF NOT EXISTS dim_customers (
                    customer_id SERIAL PRIMARY KEY,
                    first_name VARCHAR,
                    last_name VARCHAR,
                    gender VARCHAR,
                    address VARCHAR,
                    zip_code VARCHAR
                );
            """)
        logging.info("Table 'dim_customers' has been successfully created in the new database.")
    except Exception as e:
        logging.error(f"Failed to create table 'dim_customers': {e}")
        raise
```

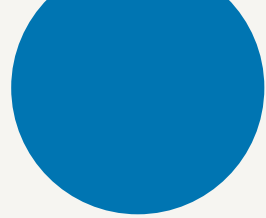
```
def transfer_customers_data():
    logging.info("Transferring data from customer to dim_customers.")
    try:
        source_engine = source_postgres_connection()
        target_engine = target_postgres_connection()

        # Read data from the source database
        with source_engine.connect() as source_conn:
            query = "SELECT id AS customer_id, first_name, last_name, gender, address, zip_code FROM cu"
            customers_df = pd.read_sql(query, source_conn)

        # Write data to the target database
        with target_engine.connect() as target_conn:
            customers_df.to_sql('dim_customers', target_conn, if_exists='replace', index=False)

        logging.info("Data has been successfully transferred from 'customer' to 'dim_customers'.")
    except Exception as e:
        logging.error(f"Failed to transfer data: {e}")
        raise
```





Menambahkan dua tabel baru, yaitu
dim_customer_detail dan dim_order_transaction,

```
# Define the DAG
with DAG(
    dag_id='modelling_dag',
    default_args=default_args,
    schedule_interval='@once',
) as dag:
    create_customer_detail_task = PostgresOperator(
        task_id='create_customer_detail_table',
        sql = """
            CREATE TABLE IF NOT EXISTS dim_customer_detail (
                customer_id INTEGER,
                full_name VARCHAR(255),
                gender VARCHAR(10)
            );
        """,
        postgres_conn_id='postgres_dw',
    )
```

```
load_customer_detail_task = PostgresOperator(
    task_id='load_customer_detail_table',
    sql="""
        INSERT INTO dim_customer_detail (customer_id, full_name, gender)
        SELECT customer_id, CONCAT(first_name, ' ', last_name) AS full_name,
            CASE
                WHEN gender = 'M' THEN 'Male'
                WHEN gender = 'F' THEN 'Female'
                ELSE NULL
            END AS gender
        FROM dim_customers;
    """,
    postgres_conn_id='postgres_dw',
)
```

```
create_order_transaction_task = PostgresOperator(
    task_id='create_order_transaction_table',
    sql = """
        CREATE TABLE IF NOT EXISTS dim_order_transaction (
            order_item_id INTEGER,
            transaction_date timestamp,
            product_name VARCHAR(100),
            price FLOAT,
            amount INTEGER
        );
    """,
    postgres_conn_id='postgres_dw',
)
```

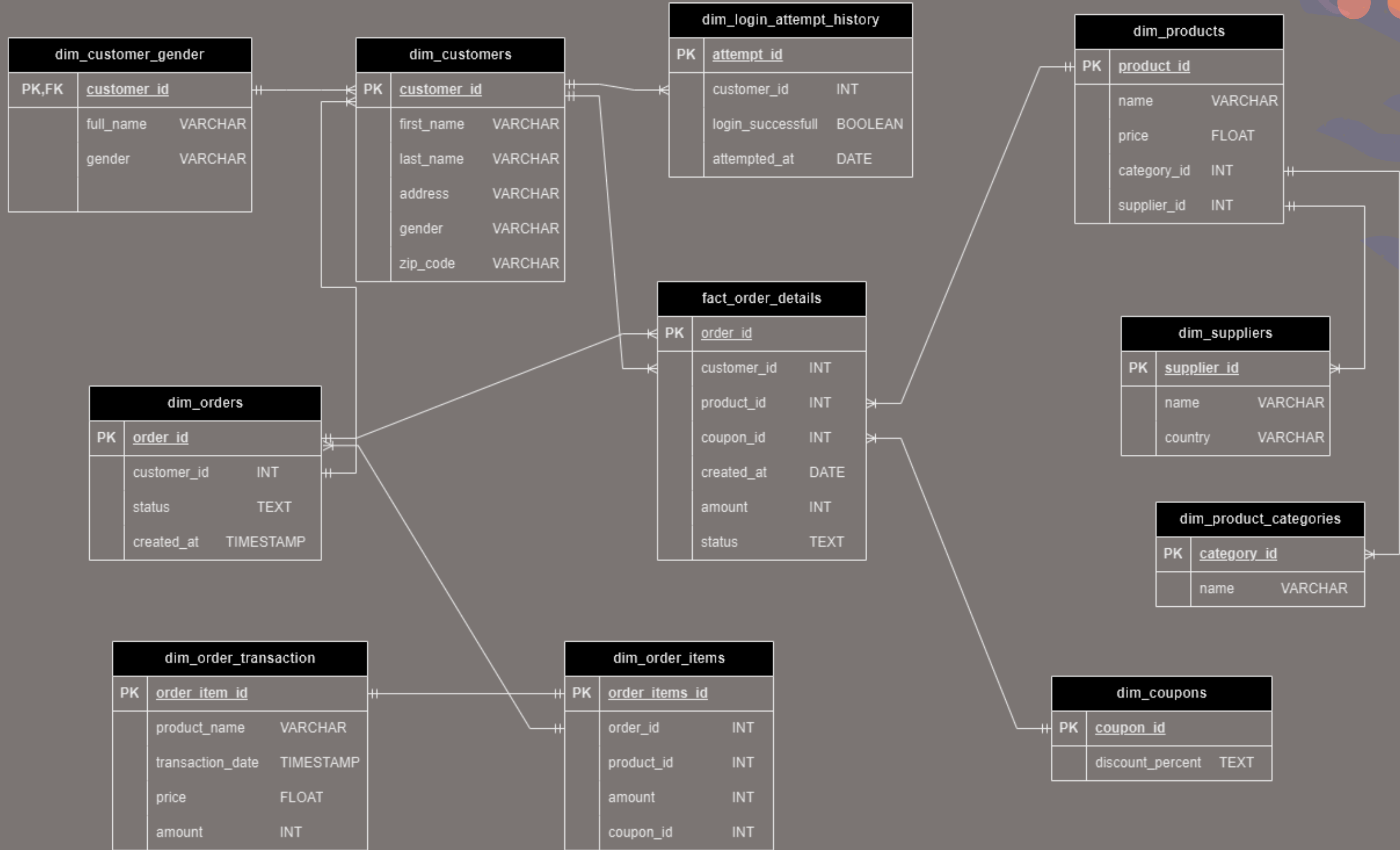
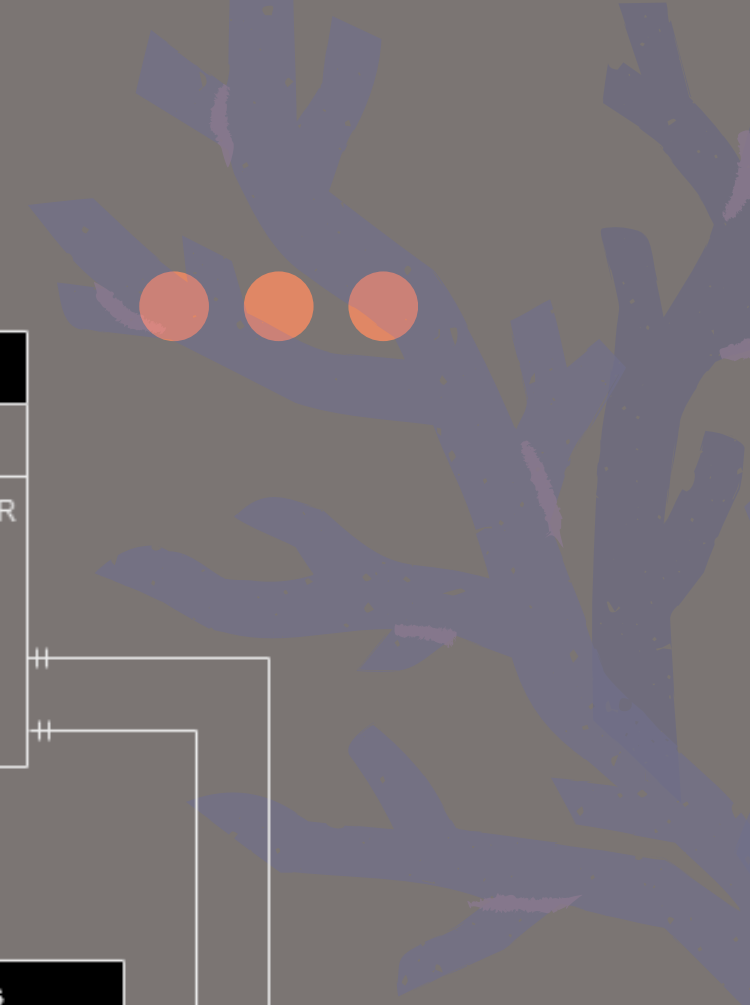
```
load_order_transaction_task = PostgresOperator(
    task_id='load_order_transaction_table',
    sql="""
        INSERT INTO dim_order_transaction (order_item_id, transaction_date, product_name, price, amount)
        SELECT oi.order_items_id, o.created_at, p.name, p.price, oi.amount
        FROM dim_order_items oi
        JOIN dim_products p ON oi.product_id = p.product_id
        JOIN dim_orders o ON oi.order_id = o.orders_id;
    """,
    postgres_conn_id='postgres_dw',
)
```

sehingga membentuk ERD seperti di selanjutnya...





ERD





DATA MART



Tables			
database: postgres_dw	⌵	<input type="checkbox"/>	customer
schema: public	⌵	<input type="checkbox"/>	0
data_mart	×	<input type="checkbox"/>	1
data_mart_coupon_effectiveness		<input type="checkbox"/>	2
<input checked="" type="checkbox"/> data_mart_customer_summary		<input type="checkbox"/>	3
data_mart_login_summary		<input type="checkbox"/>	4
data_mart_product_performance		<input type="checkbox"/>	5
data_mart_sales_summary		<input type="checkbox"/>	6
		<input type="checkbox"/>	7



DATA MART



```
requirements.txt M  finpro-dag.py U  data-mart-dag.py U x  .env M  dag-example.py
dags > data-mart-dag.py
45 df.to_sql('data_mart_sales_summary', engine, if_exists='replace', index=False)
46
47 # Fungsi untuk membuat data mart Customer Summary
48 def customer_summary():
49     engine = postgres_connection()
50     query = """
51     SELECT
52         c.id AS customer_id,
53         c.first_name || ' ' || c.last_name AS customer_name,
54         c.gender,
55         c.zip_code,
56         COUNT(DISTINCT o.id) AS total_orders,
57         SUM(p.price * oi.amount) AS total_spent,
58         AVG(p.price * oi.amount) AS average_order_value,
59         MAX(o.created_at) AS last_order_date,
60         COUNT(DISTINCT CASE WHEN lah.login_successful THEN lah.id END) AS successful_logins,
61         COUNT(DISTINCT CASE WHEN NOT lah.login_successful THEN lah.id END) AS failed_logins
62     FROM
63         customers c
64         LEFT JOIN orders o ON c.id = o.customer_id
65         LEFT JOIN order_items oi ON o.id = oi.order_id
66         LEFT JOIN products p ON oi.product_id = p.id
67         LEFT JOIN login_attempt_history lah ON c.id = lah.customer_id
68     GROUP BY
69         c.id, c.first_name, c.last_name, c.gender, c.zip_code
70     """
71     df = pd.read_sql(query, engine)
72     df.to_sql('data_mart_customer_summary', engine, if_exists='replace', index=False)
```

COUNT(DISTINCT o.id): Menghitung jumlah unik order untuk setiap pelanggan.

SUM(p.price * oi.amount): Menghitung total yang dibelanjakan oleh pelanggan.

AVG(p.price * oi.amount): Menghitung rata-rata nilai order.

MAX(o.created_at): Mengambil tanggal order terbaru.

COUNT(DISTINCT CASE WHEN...): Menghitung jumlah login sukses dan gagal.

LEFT JOIN: Digunakan untuk memastikan semua pelanggan termasuk dalam hasil, bahkan jika mereka belum pernah melakukan order.

GROUP BY: Mengelompokkan hasil berdasarkan informasi pelanggan.

DATA MART

```
requirements.txt M  finpro-dag.py U  data-mart-dag.py U x  .env M  dag-example.py
dags > data-mart-dag.py
19
20 # Fungsi untuk membuat data mart Sales Summary
21 def sales_summary():
22     engine = postgres_connection()
23     query = """
24     SELECT
25         o.id AS order_id,
26         o.created_at AS order_date,
27         c.id AS customer_id,
28         c.first_name || ' ' || c.last_name AS customer_name,
29         p.id AS product_id,
30         p.name AS product_name,
31         pc.name AS product_category,
32         oi.amount AS quantity,
33         p.price * oi.amount AS total_price,
34         COALESCE(cp.discount_percent, 0) AS discount_percent,
35         (p.price * oi.amount) * (1 - COALESCE(cp.discount_percent, 0)/100) AS discounted_price
36     FROM
37         orders o
38         JOIN customers c ON o.customer_id = c.id
39         JOIN order_items oi ON o.id = oi.order_id
40         JOIN products p ON oi.product_id = p.id
41         JOIN product_categories pc ON p.category_id = pc.id
42         LEFT JOIN coupons cp ON oi.coupon_id = cp.id
43     """
44     df = pd.read_sql(query, engine)
45     df.to_sql('data_mart_sales_summary', engine, if_exists='replace', index=False)
46
```

COALESCE(cp.discount_percent, 0): Menggunakan nilai discount_percent jika ada, atau 0 jika tidak ada (NULL).

$(p.price * oi.amount) * (1 - \text{COALESCE}(cp.discount_percent, 0)/100)$: Menghitung harga setelah diskon.

JOIN customers c ON o.customer_id = c.id: Menggabungkan tabel customers dengan orders berdasarkan customer_id.

LEFT JOIN coupons cp ON oi.coupon_id = cp.id: Menggunakan LEFT JOIN untuk mengikutsertakan semua order_items, termasuk yang tidak memiliki kupon.



DATA MART



```
requirements.txt M  finpro-dag.py U  data-mart-dag.py U x  .env M  dag-example.py
dags > data-mart-dag.py
73
74 # Fungsi untuk membuat data mart Product Performance
75 def product_performance():
76     engine = postgres_connection()
77     query = """
78     SELECT
79         p.id AS product_id,
80         p.name AS product_name,
81         pc.name AS category_name,
82         s.name AS supplier_name,
83         s.country AS supplier_country,
84         COUNT(DISTINCT oi.order_id) AS total_orders,
85         SUM(oi.amount) AS total_quantity_sold,
86         SUM(p.price * oi.amount) AS total_revenue,
87         AVG(p.price) AS average_price
88     FROM
89         products p
90         JOIN product_categories pc ON p.category_id = pc.id
91         JOIN suppliers s ON p.supplier_id = s.id
92         LEFT JOIN order_items oi ON p.id = oi.product_id
93     GROUP BY
94         p.id, p.name, pc.name, s.name, s.country
95     """
96     df = pd.read_sql(query, engine)
97     df.to_sql('data_mart_product_performance', engine, if_exists='replace', index=False)
98
```

COUNT(DISTINCT oi.order_id): Menghitung jumlah order unik untuk setiap produk.

SUM(oi.amount): Menghitung total kuantitas produk yang terjual.

SUM(p.price * oi.amount): Menghitung total pendapatan dari produk.

AVG(p.price): Menghitung rata-rata harga produk.

LEFT JOIN order_items: Memastikan semua produk termasuk, bahkan jika belum pernah dipesan.





DATA MART



```
98
99 # Fungsi untuk membuat data mart Login Summary
100 def login_summary():
101     engine = postgres_connection()
102     query = """
103     SELECT
104         c.id AS customer_id,
105         c.first_name || ' ' || c.last_name AS customer_name,
106         COUNT(lah.id) AS total_login_attempts,
107         SUM(CASE WHEN lah.login_successful THEN 1 ELSE 0 END) AS successful_logins,
108         SUM(CASE WHEN NOT lah.login_successful THEN 1 ELSE 0 END) AS failed_logins,
109         MAX(lah.attempted_at) AS last_login_attempt
110     FROM
111         customers c
112         LEFT JOIN login_attempt_history lah ON c.id = lah.customer_id
113     GROUP BY
114         c.id, c.first_name, c.last_name
115     """
116     df = pd.read_sql(query, engine)
117     df.to_sql('data_mart_login_summary', engine, if_exists='replace', index=False)
118
```

COUNT(lah.id): Menghitung total percobaan login.

SUM(CASE WHEN...): Menghitung jumlah login sukses dan gagal.

MAX(lah.attempted_at): Mengambil waktu percobaan login terakhir.

LEFT JOIN: memastikan bahwa semua pelanggan dari tabel 'customers' dimasukkan ke dalam hasil, bahkan jika mereka tidak memiliki catatan login di tabel 'login_attempt_history'.





DATA MART



```
119 # Fungsi untuk membuat data mart Coupon Effectiveness
120 def coupon_effectiveness():
121     engine = postgres_connection()
122     query = """
123     SELECT
124         cp.id AS coupon_id,
125         cp.discount_percent,
126         COUNT(DISTINCT oi.order_id) AS orders_used,
127         COUNT(DISTINCT oi.product_id) AS products_discounted,
128         SUM(p.price * oi.amount) AS total_pre_discount_value,
129         SUM(p.price * oi.amount * (1 - cp.discount_percent/100)) AS total_post_discount_value,
130         SUM(p.price * oi.amount * (cp.discount_percent/100)) AS total_discount_amount
131     FROM
132         coupons cp
133         JOIN order_items oi ON cp.id = oi.coupon_id
134         JOIN products p ON oi.product_id = p.id
135     GROUP BY
136         cp.id, cp.discount_percent
137     """
138     df = pd.read_sql(query, engine)
139     df.to_sql('data_mart_coupon_effectiveness', engine, if_exists='replace', index=False)
140
```

COUNT(DISTINCT oi.order_id): Menghitung jumlah order yang menggunakan kupon.

COUNT(DISTINCT oi.product_id): Menghitung jumlah produk yang didiskon.

SUM(p.price * oi.amount): Menghitung total nilai sebelum diskon.

SUM(p.price * oi.amount * (1 - cp.discount_percent/100)): Menghitung total nilai setelah diskon.

SUM(p.price * oi.amount * (cp.discount_percent/100)): Menghitung total jumlah diskon yang diberikan.

JOIN: Hanya akan menampilkan kupon yang benar-benar digunakan dalam pesanan yang akan muncul



Data Visualization

DASHBOARD

Data dashboard adalah tampilan visual yang menyajikan berbagai jenis data dalam satu tempat secara terpusat.

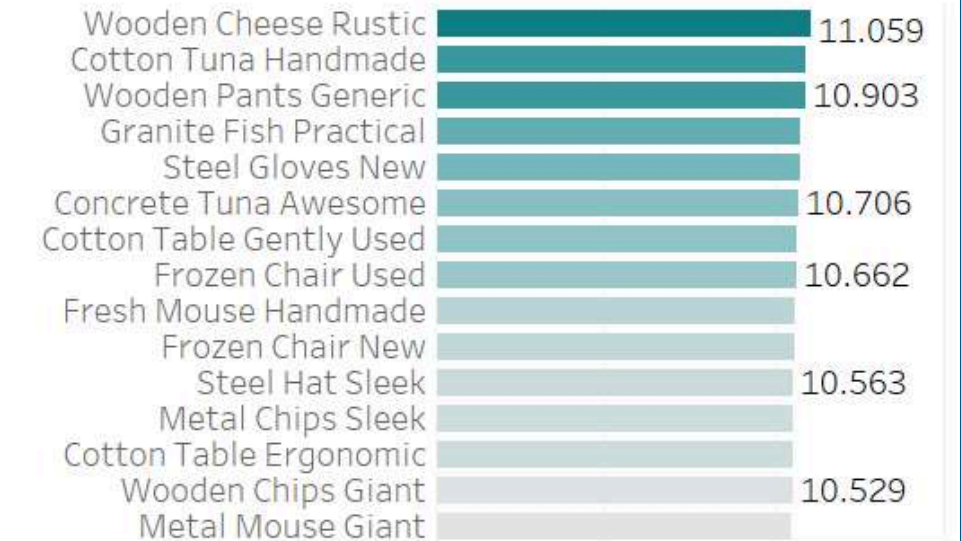
- **Supplier map** menunjukkan distribusi persebaran banyak order yang terjadi di setiap negara.
- **Top 15 Product** menunjukkan 15 produk yang memiliki penjualan tertinggi.
- **Average price** menunjukkan rata-rata harga pada setiap kategori produk.
- **Total revenue** menunjukkan besar keseluruhan pendapatan berdasarkan kategori produk.

Distribution Product Category

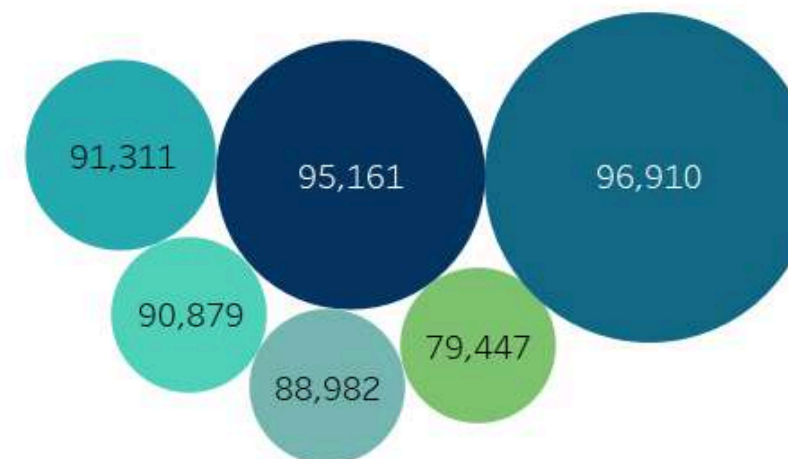
Supplier Map



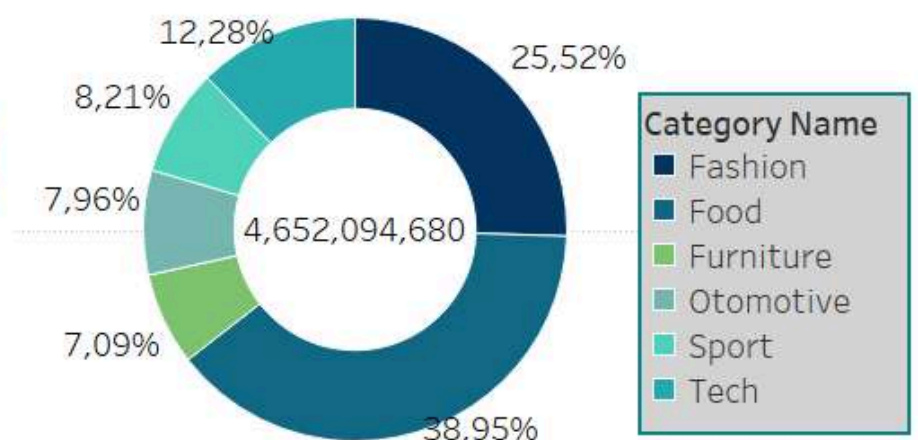
Top 15 Sold Products



Average Price



Total Revenue



THANK YOU!
FOR YOUR ATTENTION



