



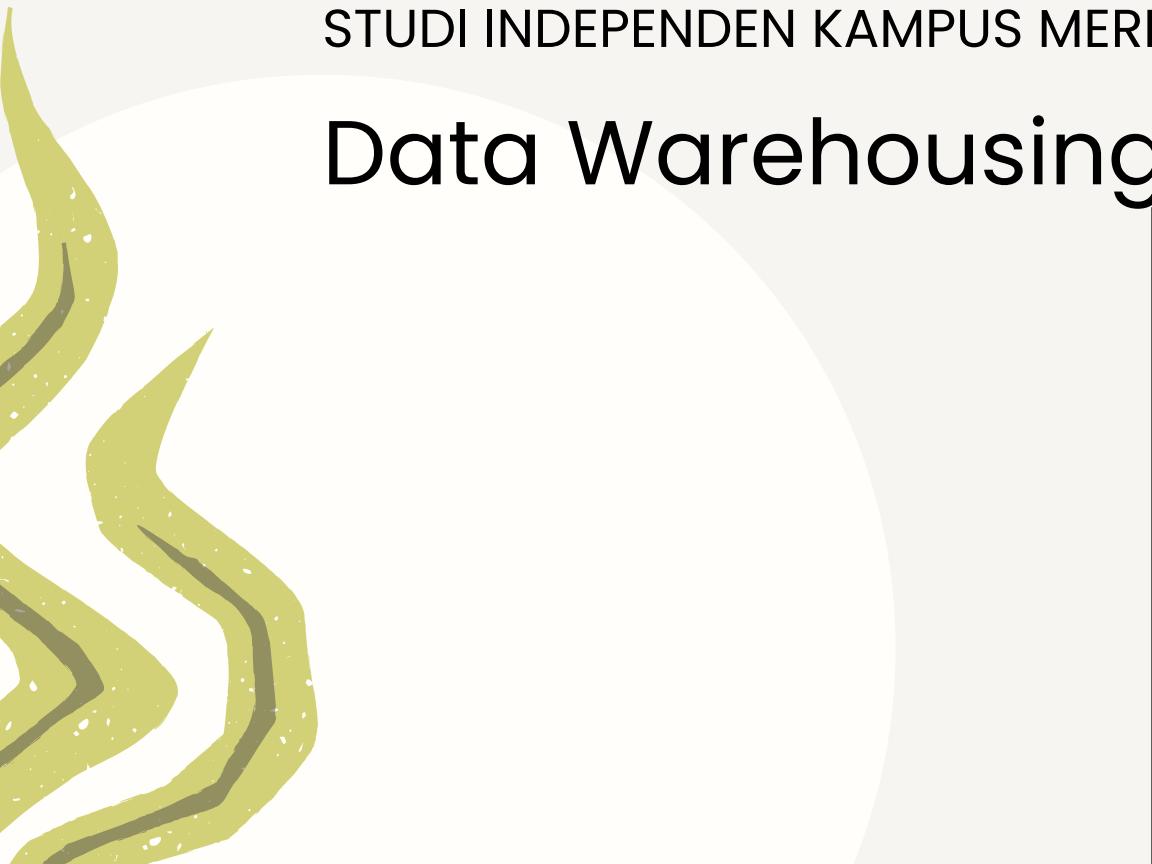
TEAM 1

FINAL PROJECT

Presentation

STUDI INDEPENDEN KAMPUS MERDEKA - Batch 6

Data Warehousing

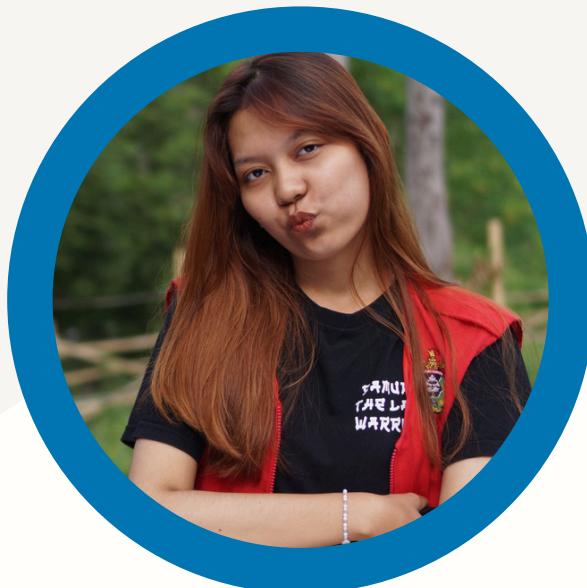


HELLO!

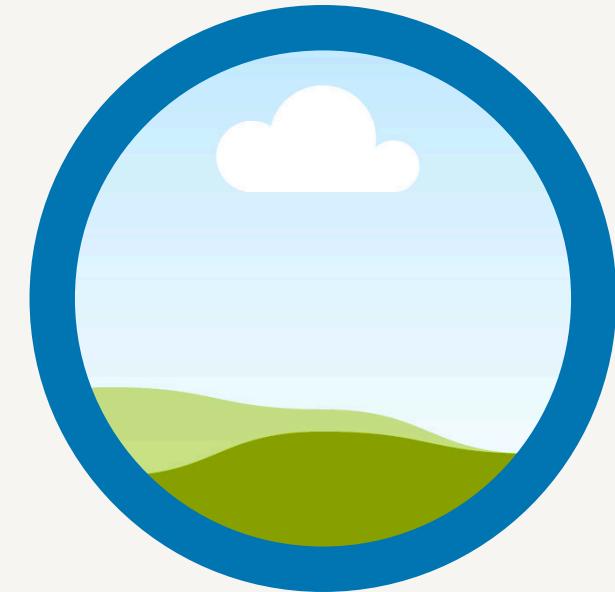
MEET OUR TEAM



INTRODUCTION



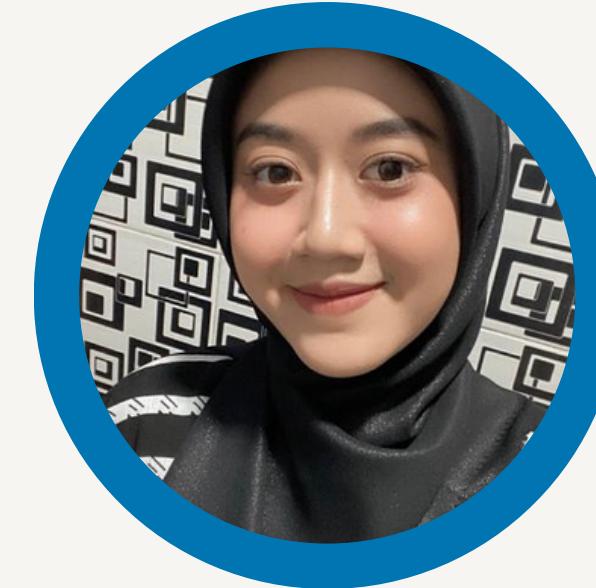
Clarisa Eudia
Chesynanda



INTRODUCTION



Suci Rakhmayanti



Chika
Cardealfitri



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

final project

Tugas akhir dari program msib bach 6
data warehousing.

Pada tugas ini menggunakan
idx.dev, neon.db, Github dan tableau



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

The screenshot shows the Airflow web interface for the DAG `ingest_data`. The main panel displays the `DAG Runs Summary` for two runs:

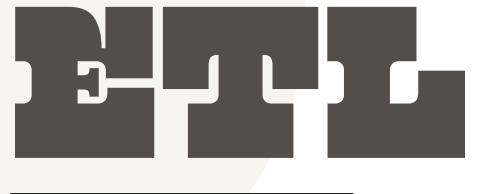
Run Type	Count
Total success	2

Key summary statistics for the runs:

- Total Runs Displayed: 2
- First Run Start: 2024-06-23, 06:44:31 UTC
- Last Run Start: 2024-06-23, 06:46:30 UTC
- Max Run Duration: 00:10:18
- Mean Run Duration: 00:09:31
- Min Run Duration: 00:08:44

The sidebar on the left lists the tasks in the DAG: `ingest_order_item`, `ingest_data_coupons`, `ingest_data_customer`, `ingest_login_attempts`, `ingest_product_category`, `ingest_products`, `ingest_supplier_data`, and `ingest_order_data`. Each task is represented by a small icon and a green vertical bar indicating its duration.

3.



EXPLORER

- DIBIMBING_KM_FINAL_PROJECT
 - .idx
 - dags
 - dag-example.py
 - dag-ingest-data.py**
 - data
 - coupons.json
 - customer_0.csv
 - customer_1.csv
 - customer_2.csv
 - customer_3.csv
 - customer_4.csv
 - customer_5.csv
 - customer_6.csv
 - customer_7.csv
 - customer_8.csv
 - customer_9.csv
 - login_attempts_0.json
 - login_attempts_1.json
 - login_attempts_2.json
 - login_attempts_3.json
 - login_attempts_4.json
 - login_attempts_5.json
 - login_attempts_6.json

requirements.txt M makefile dag-ingest-data.py U .env M dag-e

```

dags > dag-ingest-data.py
    owner: 'airflow',
    'start_date': datetime(2023, 1, 12),
    'depends_on_past': False,
    'retries': 0,
    'retry_delay': timedelta(minutes=5),
}

def postgres_connection():
    connection_string = URL.create(
        'postgresql',
        username='data_warehouse_owner',
        password='SwEm1h0bTgWz',
        host='ep-lively-tooth-airqoefa.ap-southeast-1.aws.neon.tech',
        database='data_warehouse',
        port=5432,
        query={'sslmode': 'require'}
    )
    engine = create_engine(connection_string)
    return engine

def ingest_coupons():
    with open('data/coupons.json') as f:
        data = json.load(f)

    df = pd.DataFrame(data)
    df.to_sql("coupons", engine, if_exists="replace", index=False)

```

OUTLINE
TIMELINE

main* ⏪ 0 △ 0 Gemini Ln 44

EXPLORER

- DIBIMBING_KM_FINAL_PROJECT
 - .idx
 - dags
 - dag-example.py
 - dag-ingest-data.py**
 - data
 - coupons.json
 - customer_0.csv
 - customer_1.csv
 - customer_2.csv
 - customer_3.csv
 - customer_4.csv
 - customer_5.csv
 - customer_6.csv
 - customer_7.csv
 - customer_8.csv
 - customer_9.csv
 - login_attempts_0.json
 - login_attempts_1.json
 - login_attempts_2.json
 - login_attempts_3.json
 - login_attempts_4.json
 - login_attempts_5.json
 - login_attempts_6.json

requirements.txt M makefile dag-ingest-data.py

```

t5 = PythonOperator(
    task_id='ingest_product_category',
    python_callable=ingest_product_category,
    dag=dag
)

t6 = PythonOperator(
    task_id='ingest_products',
    python_callable=ingest_products,
    dag=dag
)

t7 = PythonOperator(
    task_id='ingest_supplier_data',
    python_callable=ingest_supplier_data,
    dag=dag
)

t8 = PythonOperator(
    task_id='ingest_order_data',
    python_callable=ingest_order_data,
    dag=dag
)

t1 >> t2 >> t3 >> t4 >> t5 >> t6 >> t7 >> t8

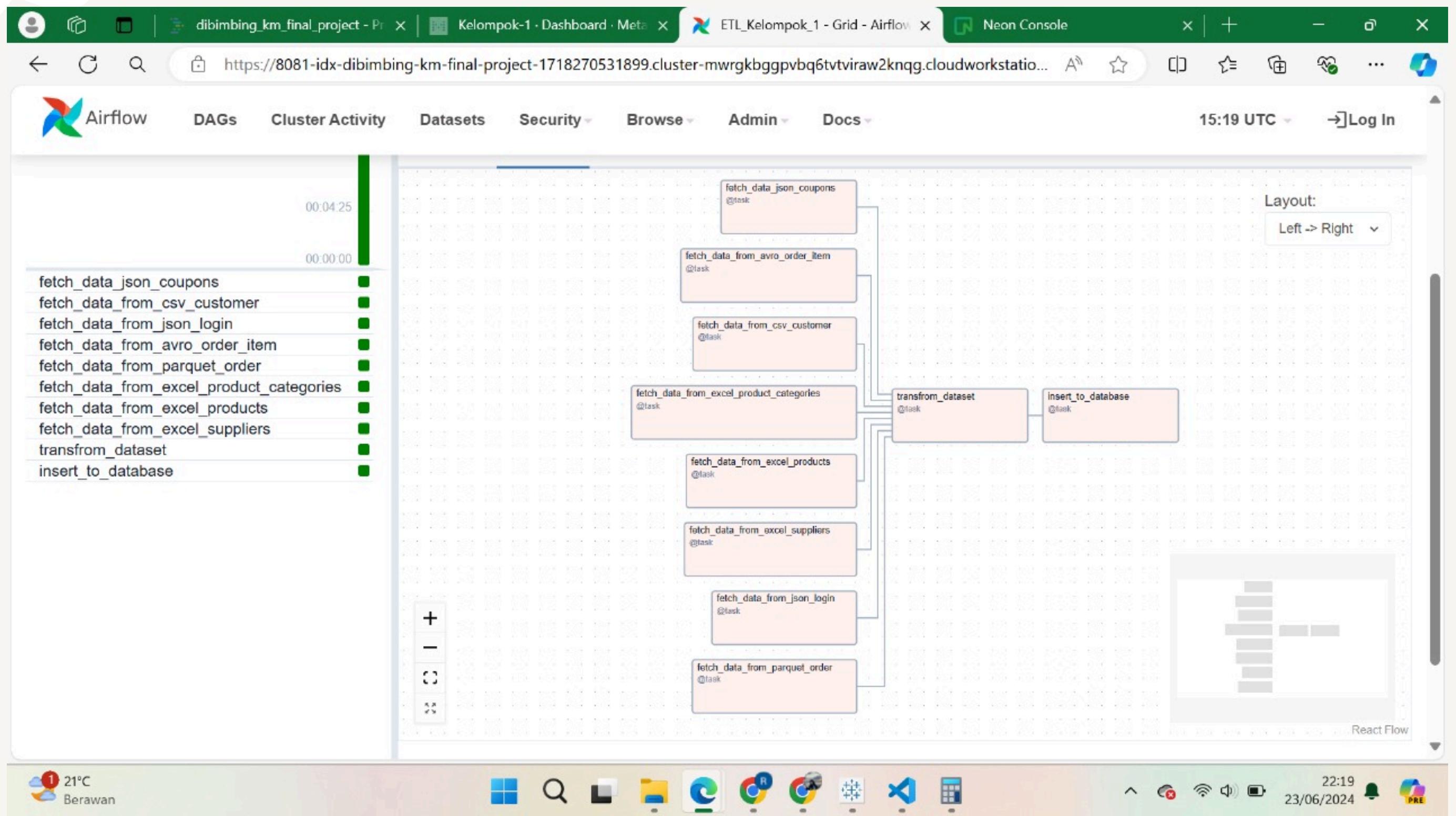
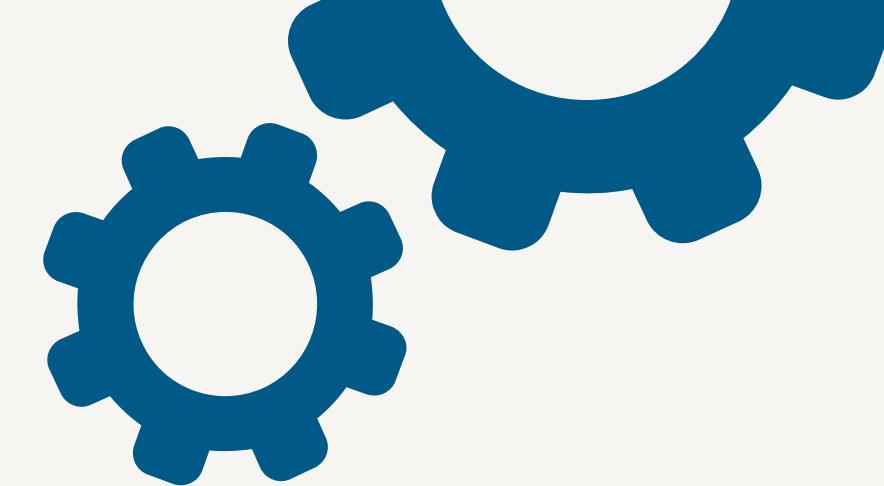
```

OUTLINE
TIMELINE

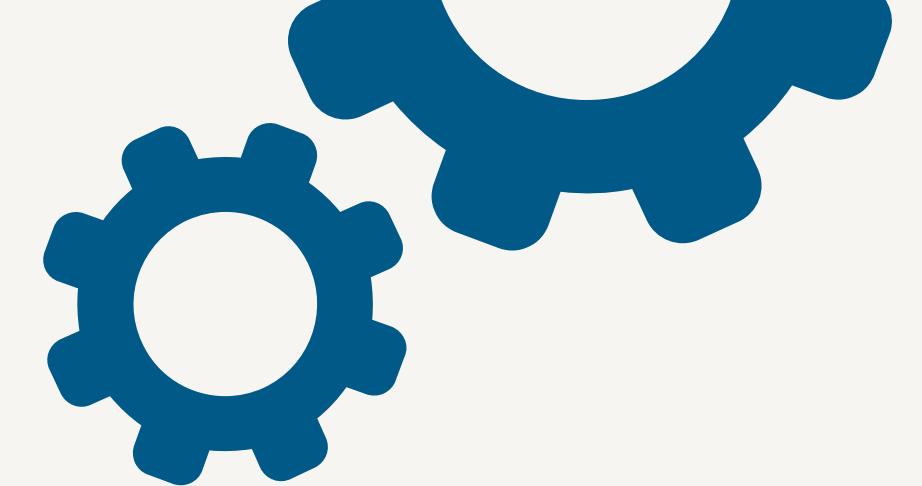
main* ⏪ 0 △ 0 Gemini Ln 44

4.

ETL



5. DATA MODELLING



The screenshot shows the Apache Airflow web interface. The top navigation bar includes links for 'dibimbing_km_final_project - Pr...', 'Kelompok-1 · Dashboard · Meta...', 'factdim_tabel_dag - Grid - Airflow...', 'Neon Console', and a search bar with the URL 'https://8081-idx-dibimbing-km-final-project-1718270531899.cluster-mwrgkbggpvbq6tvviraw2knqg.cloudworkstation...'. The main header features the Airflow logo, 'DAGs', 'Cluster Activity', 'Datasets', 'Security', 'Browse', 'Admin', 'Docs', the time '15:20 UTC', and a 'Log In' button.

The left sidebar displays a list of DAG tasks with their execution times:

- create_dim_coupons
- transfer_coupons
- create_dim_customers
- transfer_customers
- create_dim_login_attempt_history
- transfer_login_attempt_history
- create_dim_product_categories
- transfer_product_categories
- create_dim_products
- transfer_products
- create_dim_orders
- transfer_orders
- create_dim_suppliers
- transfer_suppliers
- create_dim_order_items
- transfer_order_items
- create_fact_order_details
- transfer_order_details

The main area shows a DAG graph with nodes arranged in a grid. Each node is a box labeled with a task name and 'PythonOperator'. The nodes are connected by arrows indicating dependencies. A legend on the right indicates that green boxes represent successful runs and orange boxes represent failed runs. A 'Layout' dropdown menu shows 'Left -> Right'.

The bottom of the screen shows a taskbar with icons for Windows, search, file, Microsoft Edge, Google Chrome, File Explorer, and other applications, along with system status indicators like battery level and signal strength.

DATA MODELING



```
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-alrqaefa.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-a5fcvgv3.us-east-2.aws.neon.tech',
35         port=5432,
36         database='data_warehouse',
37         query={'sslmode': 'require'}
```

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

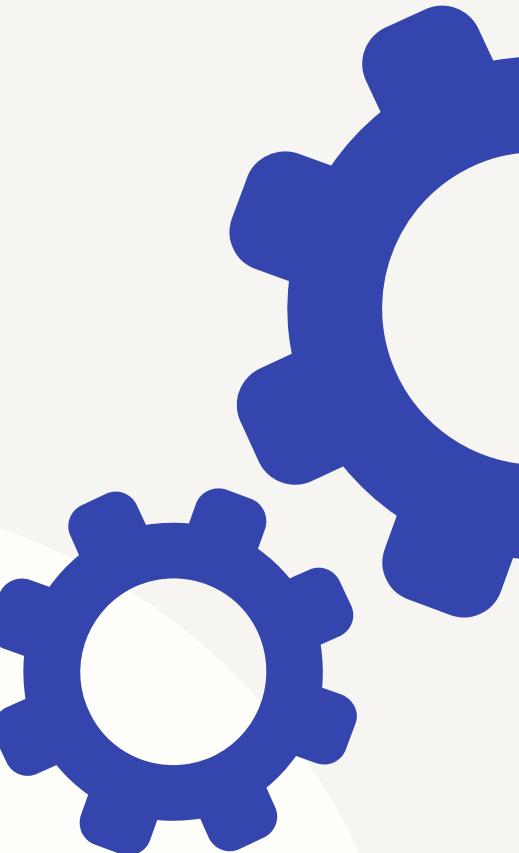
```
def create_dim_customers_table():
    logging.info("Creating dim_customers table in the target database.")
    try:
        engine = target_postgres_connection()
        with engine.connect() as connection:
            connection.execute("""
                CREATE TABLE IF NOT EXISTS dim_customers (
                    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
```



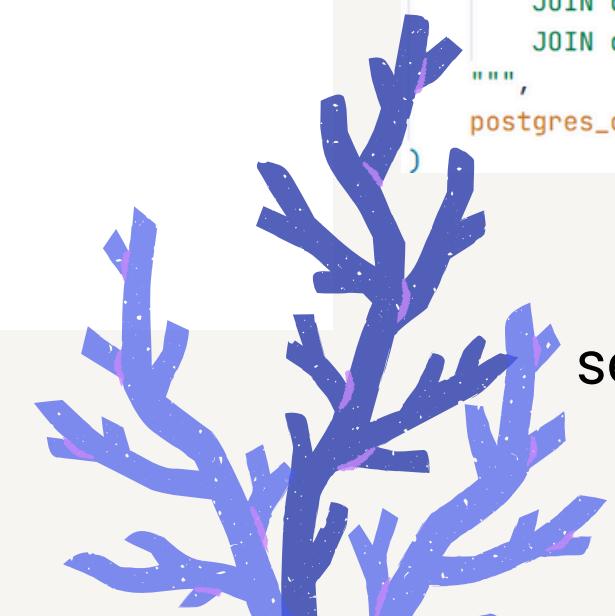
Menambahkan dua tabel baru, yaitu dim_customer_detail dan dim_order_transaction,

```
# 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',
)
```

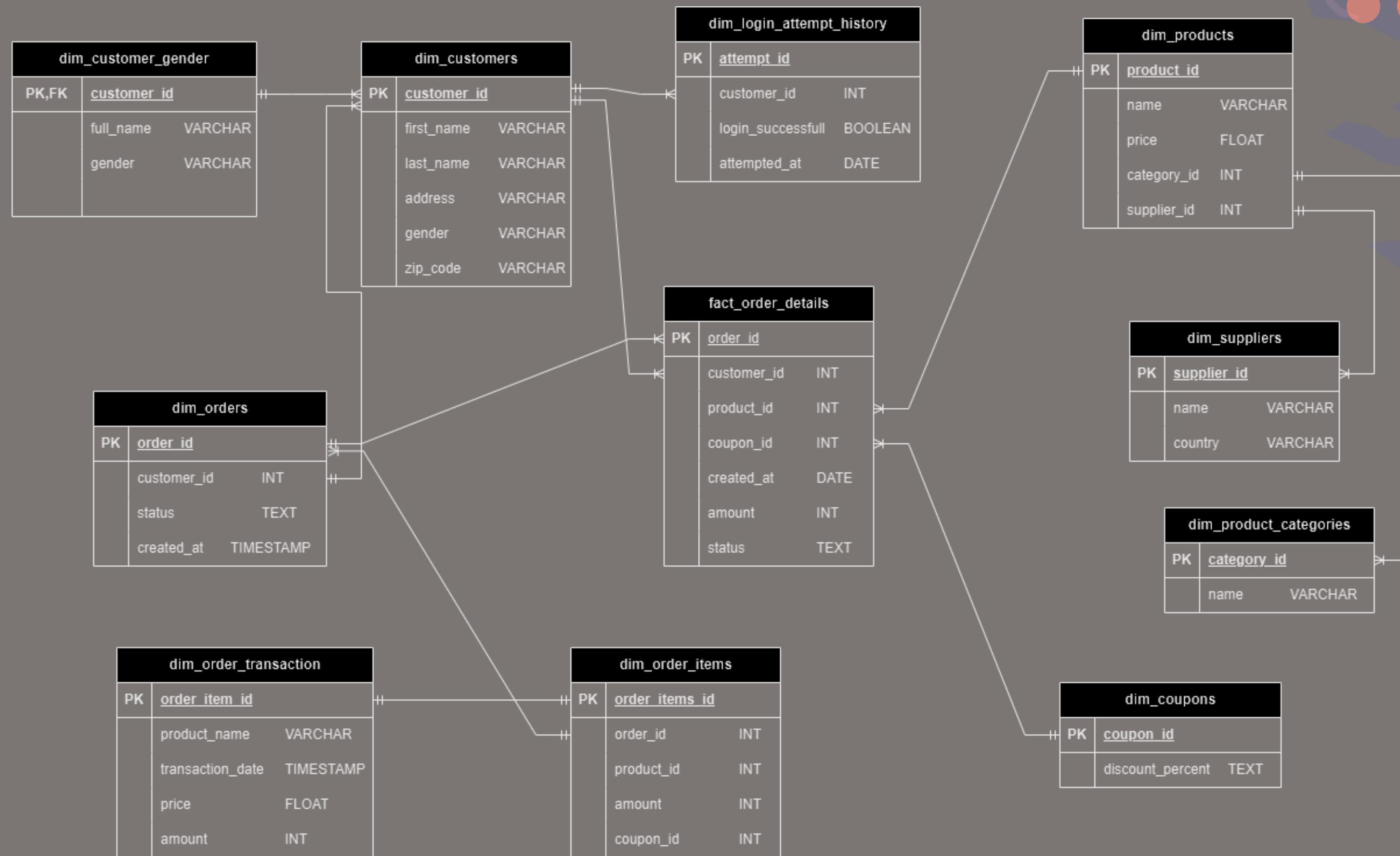
```
create_order_transaction_task = PostgresOperator(
    task_id='create_order_transaction_table',
    sql = """
        CREATE TABLE IF NOT EXISTS dim_order_transaction (
            order_item_id INTEGER,
            transaction_date timestamp,
            product_name VARCHAR(100),
            price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

schema: public

data_mart

- data_mart_coupon_effectiveness
- data_mart_customer_summary**
- data_mart_login_summary
- data_mart_product_performance
- data_mart_sales_summary

	customer_
0	
1	
2	
3	
4	
5	
6	
7	





DATA MART



```
... requirements.txt M finpro-dag.py U data-mart-dag.py U .env M dag-example.py
dags > data-mart-dag.py
45     dt.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
72     df = pd.read_sql(query, engine)
73     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 .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
45     df = pd.read_sql(query, engine)
46     df.to_sql('data_mart_sales_summary', engine, if_exists='replace', index=False)
```

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

$(p.price * oi.amount) * (1 -$

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)
```

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
00 def login_summary():
01     engine = postgres_connection()
02     query = """
03         SELECT
04             c.id AS customer_id,
05             c.first_name || ' ' || c.last_name AS customer_name,
06             COUNT(lah.id) AS total_login_attempts,
07             SUM(CASE WHEN lah.login_successful THEN 1 ELSE 0 END) AS successful_logins,
08             SUM(CASE WHEN NOT lah.login_successful THEN 1 ELSE 0 END) AS failed_logins,
09             MAX(lah.attempted_at) AS last_login_attempt
10         FROM
11             customers c
12             LEFT JOIN login_attempt_history lah ON c.id = lah.customer_id
13         GROUP BY
14             c.id, c.first_name, c.last_name
15             """
16     df = pd.read_sql(query, engine)
17     df.to_sql('data_mart_login_summary', engine, if_exists='replace', index=False)
18
```

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



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

```
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
139     df = pd.read_sql(query, engine)
140     df.to_sql('data_mart_coupon_effectiveness', engine, if_exists='replace', index=False)
```



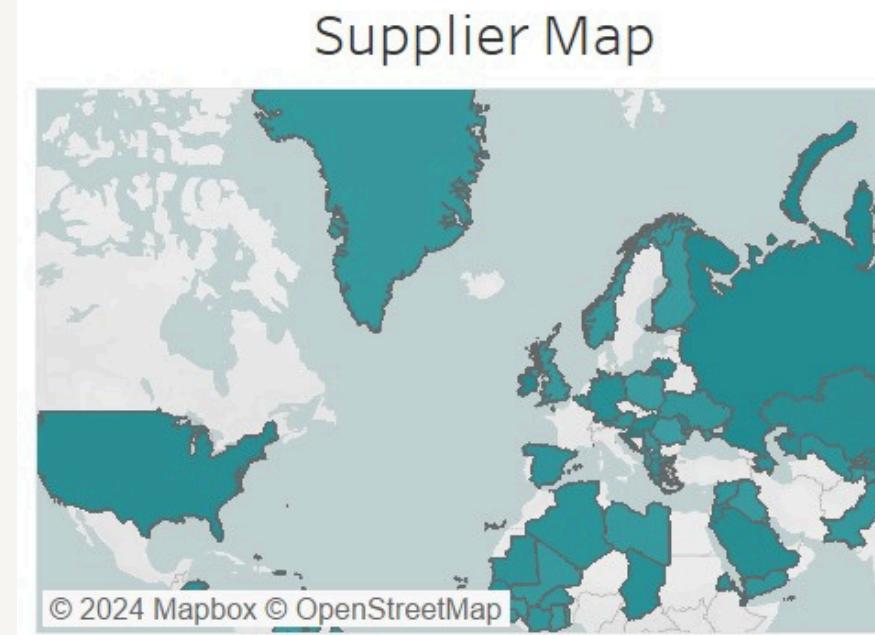
Distribution Product Category

DASHBOARD

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

- Supplier map menunjukkan distribusi 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



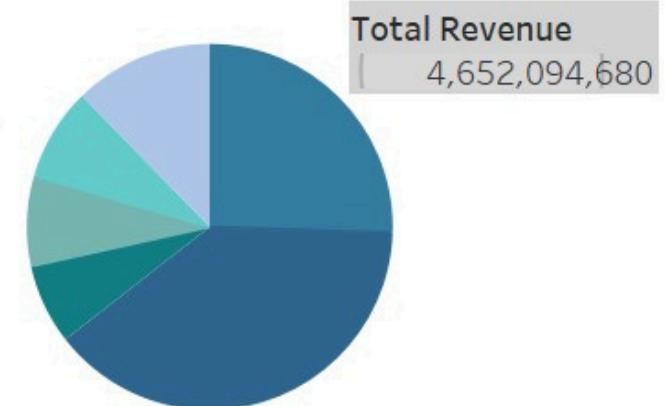
Top 15 Sold Products

Wooden Cheese Rustic	
Cotton Tuna Handmade	
Wooden Pants Generic	
Granite Fish Practical	
Steel Gloves New	
Concrete Tuna Awesome	
Cotton Table Gently Used	
Frozen Chair Used	
Fresh Mouse Handmade	
Frozen Chair New	
Steel Hat Sleek	
Metal Chips Sleek	
Cotton Table Ergonomic	
Wooden Chips Giant	
Metal Mouse Giant	

Average Price



Total Revenue



THANK YOU!
FOR YOUR ATTENTION



