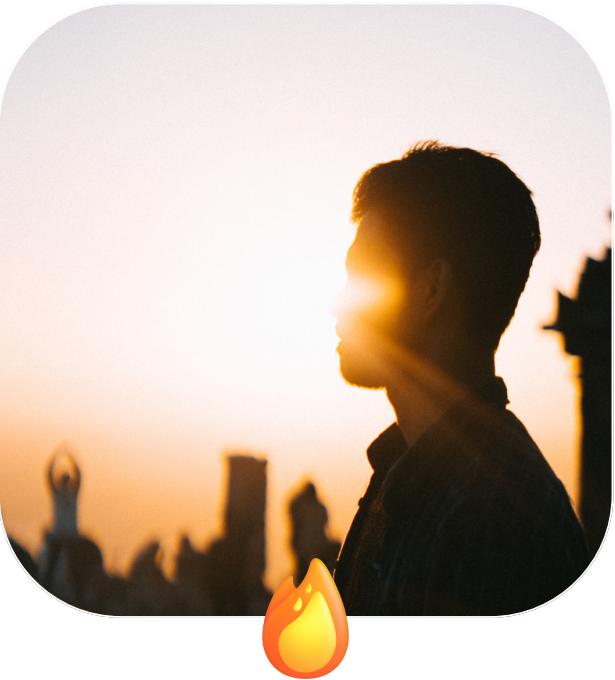


Kelompok 13

Final Project Presentation

Modeling And Optimization Techniques In Data Warehousing Bootcamp

Profile Team



Adnanfadhil Yaser
Universitas Pembangunan
Nasional Veteran Yogyakarta



Fajar Mulia Ananda
Universitas Gunadarma



Sela Aziza
Universitas Muhammadiyah
Kotabumi



Putri Sitti Naima
Institut Teknologi Sepuluh
Nopember

Background

Data Engineer!

US online retail company menjual produk pelanggan umum langsung ke pelanggan dari berbagai pemasok di seluruh dunia. Tantangannya adalah membangun infrastruktur data menggunakan data yang dihasilkan yang dibuat untuk mencerminkan data dunia nyata dari perusahaan. Inilah tugasnya :

- **ETL/ELT Job Creation menggunakan Airflow**
- **Data Modeling di Postgres**
- **Dashboard Creation dengan Data Visualisasi**

Stack

1

Data

Mendapatkan 5 format data yang berbeda

2

Airflow

Airflow DAG and mempersiapkan ETL Data

3

Data Modelling

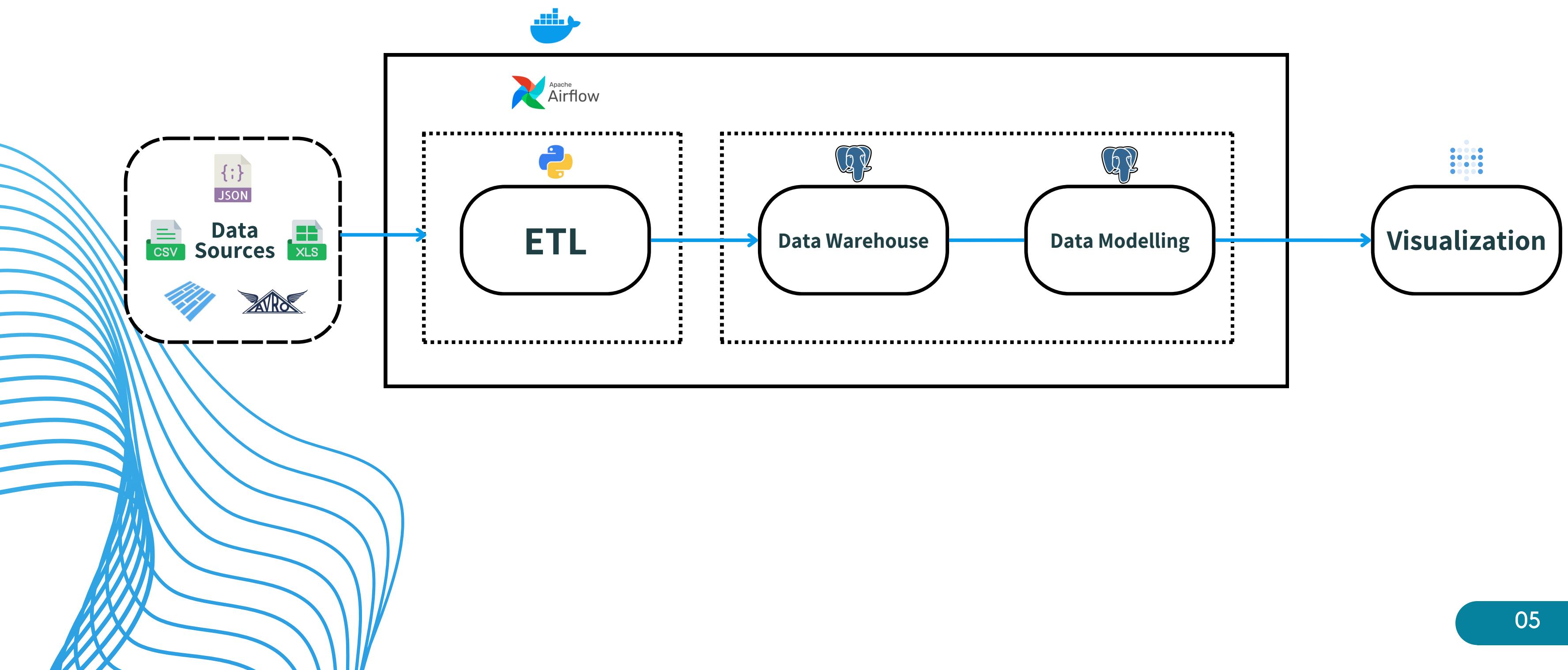
Struktur dari data modelling

4

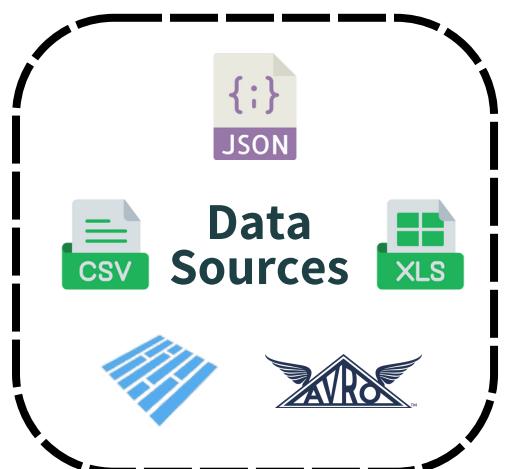
Visualisasi

Visualisasi data yang telah diproses

Data Pipeline



Data



```
requirements.txt X  
Dibimbing_FinalTest-main > requirements.txt  
1 pandas==2.1.0  
2 openpyxl  
3 fastparquet==0.8.2  
4 fastavro==1.4.6  
5 avro-python3==1.10.2  
6 pyarrow==5.0.0  
7 psycopg2-binary==2.9.1  
8 xlrd==1.2.0  
9 python-snappy==0.6.0
```

```
from datetime import datetime, timedelta
from airflow import DAG
from airflow.providers.postgres.hooks.postgres import PostgresHook
from airflow.operators.python import PythonOperator
import pandas as pd
from sqlalchemy import create_engine
import pyarrow.parquet as pq
import fastavro
import os
import json
```

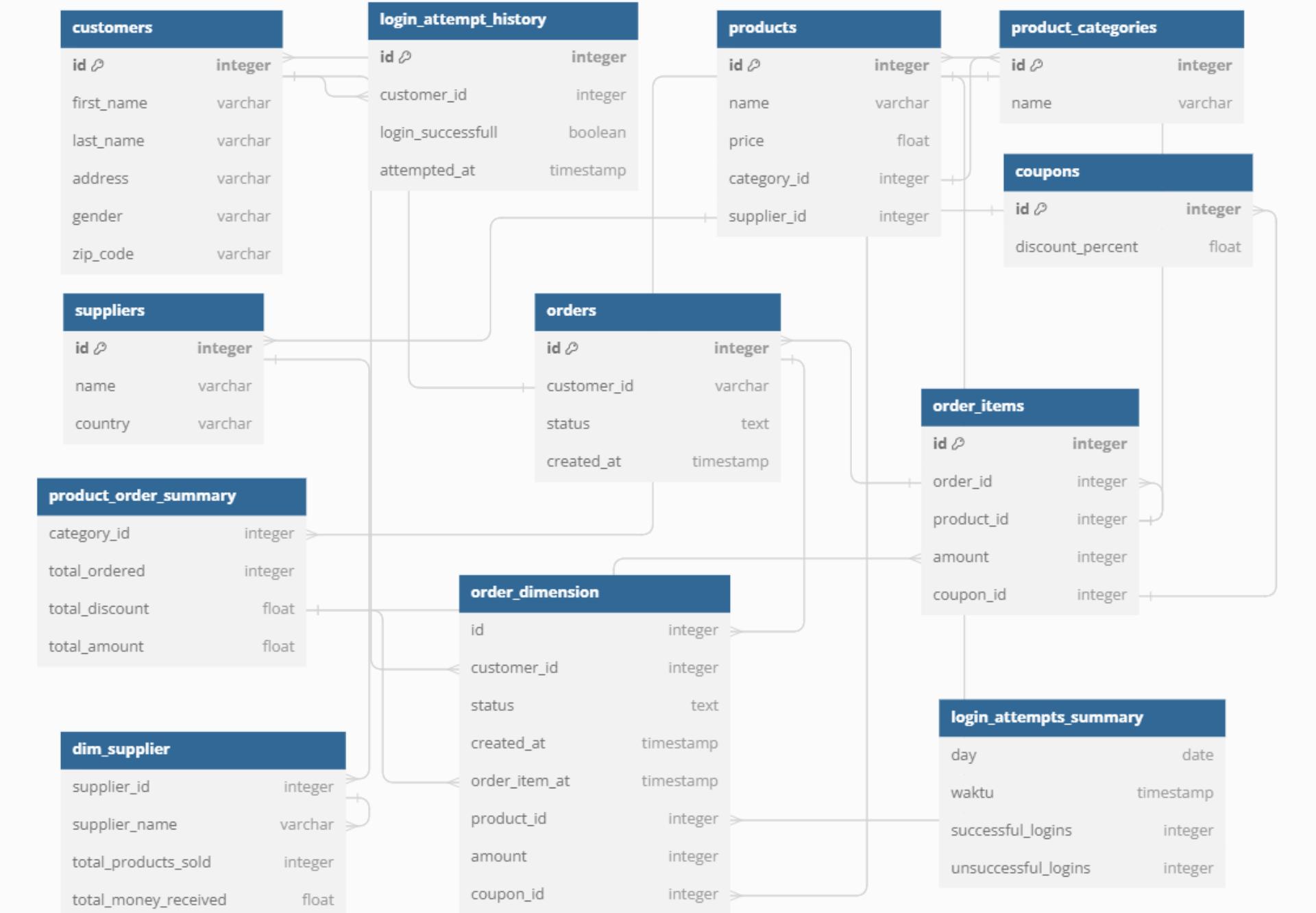
Yang dibutuhkan untuk mengolah lima data dengan format yang berbeda

Airflow

The screenshot shows the Airflow web interface with a list of Data Acquisition and Processing (DAG) tasks. Each task is represented by a row in the table, showing its name, owner, number of runs, last run's start time, and a series of green and red circles indicating run status across multiple execution slots.

DAG Name	Owner	Last Run (Start Time)	Status
create_dimension_tables_dag	Adrian	2023-01-01, 00:00:00	Green
extract_and_load_customer_data_to_postgres	airflow	2023-12-07, 07:24:10	Green (2 runs)
extract_and_load_xls_to_postgres	airflow	2023-12-07, 07:36:56	Red (1 run)
extract_load_login_attempts_to_postgres	airflow	2023-12-07, 06:57:08	Red (1 run)
Ingest-Category	Adrian	2023-01-01, 00:00:00	Green
Ingest-coupon	Adrian	2023-12-07, 07:24:20	Green (4 runs)
Ingest-Customers	Adrian	2023-01-01, 00:00:00	Green
Ingest-Geolocation	Adrian	2023-01-01, 00:00:00	Green
Ingest-Items	Adrian	2023-01-01, 00:00:00	Green
Ingest-login_attempts	Adrian	2023-01-01, 00:00:00	Green
Ingest-Order-Payment	Adrian	2023-01-01, 00:00:00	Green
Ingest-Orders	Adrian	2023-01-01, 00:00:00	Green
Ingest-Orders-Parquet	Adrian	2023-01-01, 00:00:00	Green (1 run)
Ingest-Product	Adrian	2023-01-01, 00:00:00	Green
Ingest-Review	Adrian	2023-01-01, 00:00:00	Green
Ingest-Sellers	Adrian	2023-01-01, 00:00:00	Green
Order_item	Adrian	2023-01-01, 00:00:00	Red (2 runs)
trigger_dags_in_order	Adrian	2023-01-01, 00:00:00	Green

Data Modelling

[Lihat disini](#)

Data Modelling

[Lihat disini](#)

product_order_summary	
category_id	integer
total_ordered	integer
total_discount	float
total_amount	float

Untuk melihat seberapa baik produk terjual, kategori apa yang paling diminati

order_dimension	
id	integer
customer_id	integer
status	text
created_at	timestamp
order_item_at	timestamp
product_id	integer
amount	integer
coupon_id	integer

Untuk melacak detail pesanan, pelanggan.

dim_supplier	
supplier_id	integer
supplier_name	varchar
total_products_sold	integer
total_money_received	float

Untuk menilai pemasok berkinerja dan berapa banyak uang yang diperoleh dari penjualan produk pemasok.

login_attempts_summary	
day	date
waktu	timestamp
successful_logins	integer
unsuccessful_logins	integer

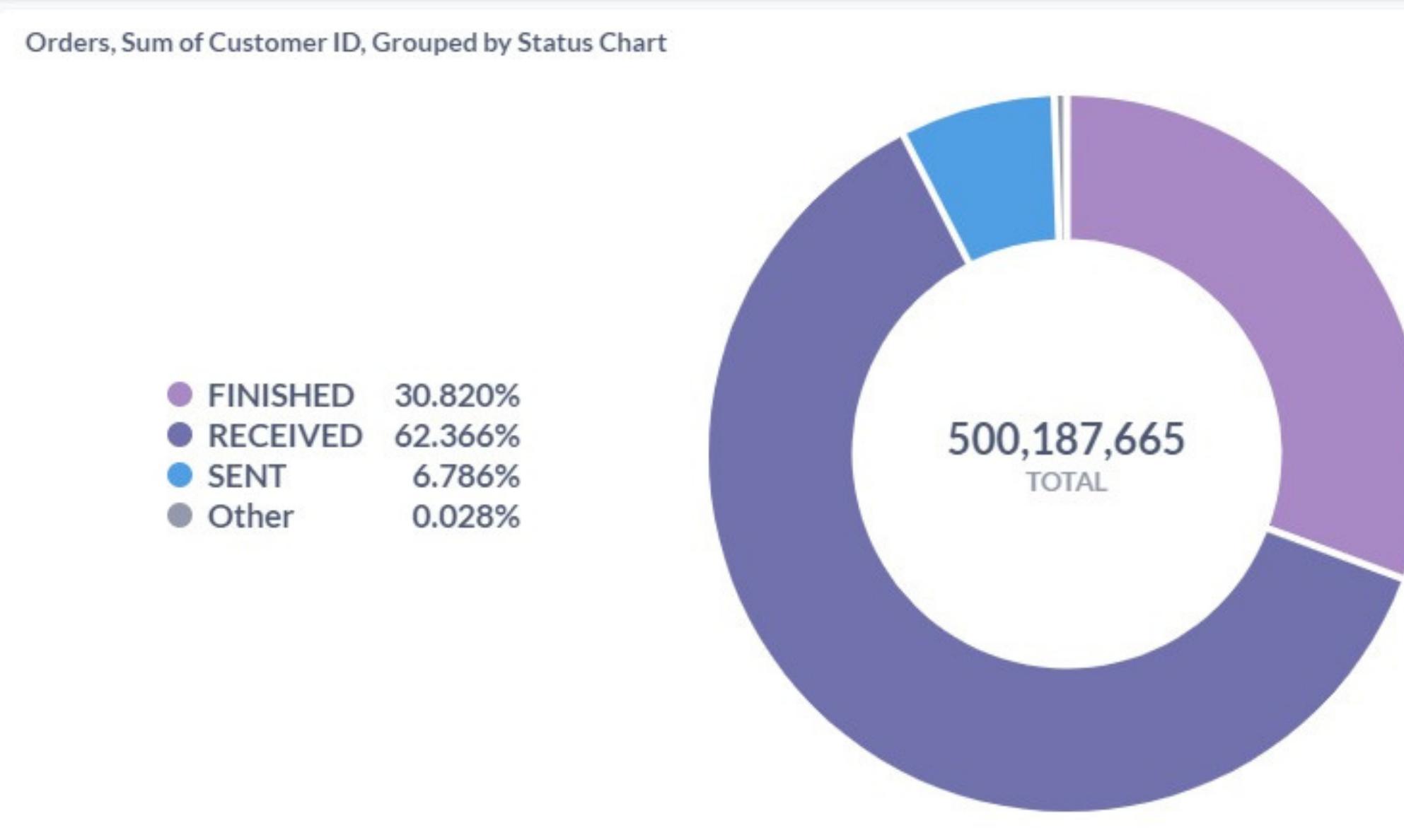
Untuk melihat pola aktivitas login

Visualisasi



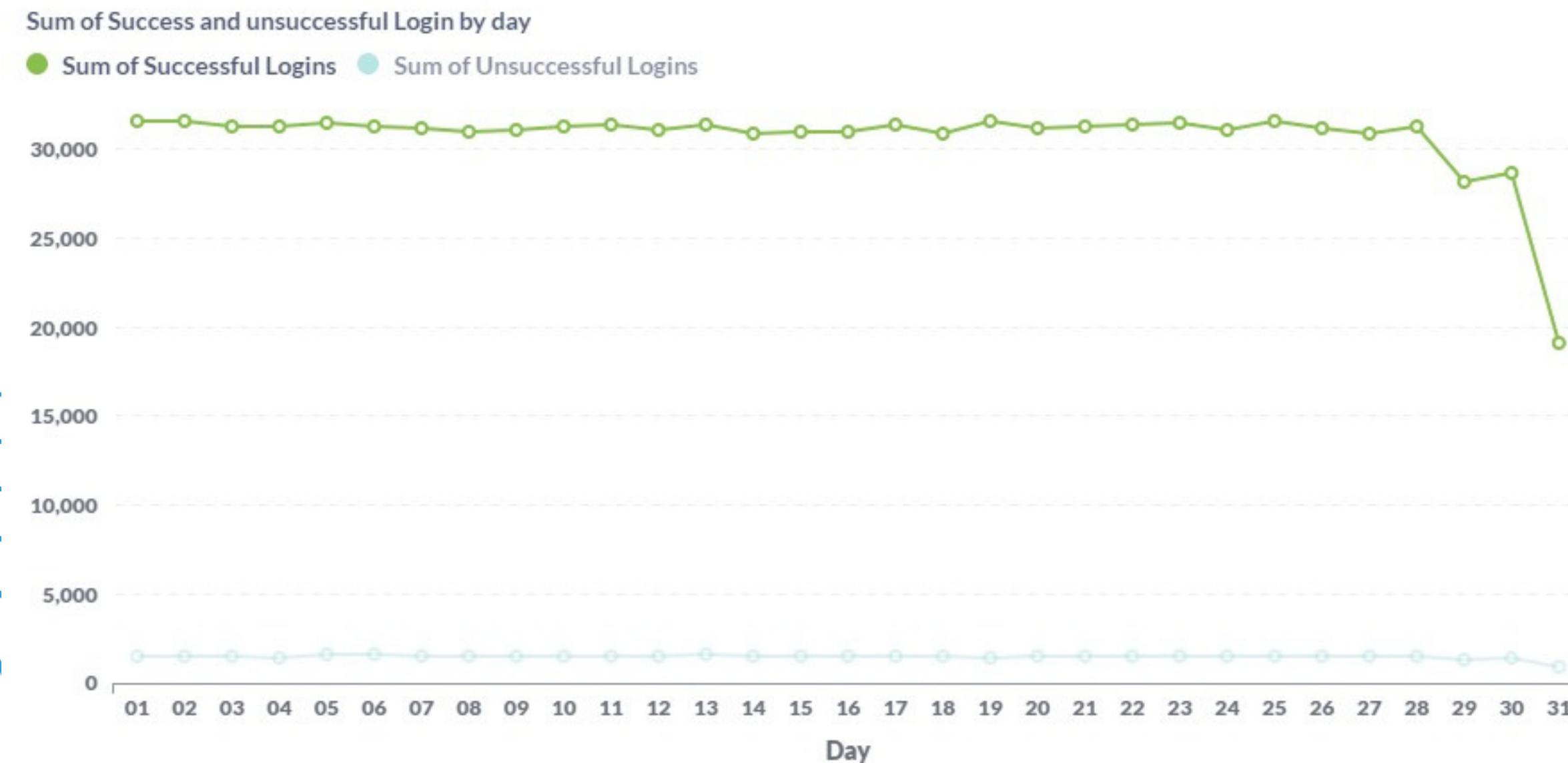
Hasil dari visualisasi terlihat bahwa kategori ‘food’ memiliki jumlah total order yang terbanyak mencapai \$18.8M, dan diikuti kategori ‘fashion’ mencapai \$ 12.5M. Kedua kategori ini yang memiliki pengaruh pendapatan pada perusahaan ini.

Visualisasi



Hasil dari visualisasi terlihat bahwa detail pesanan pada bulan tersebut mayoritas sudah diterima dan selesai oleh pelanggan yaitu sebesar 93%.

Visualisasi

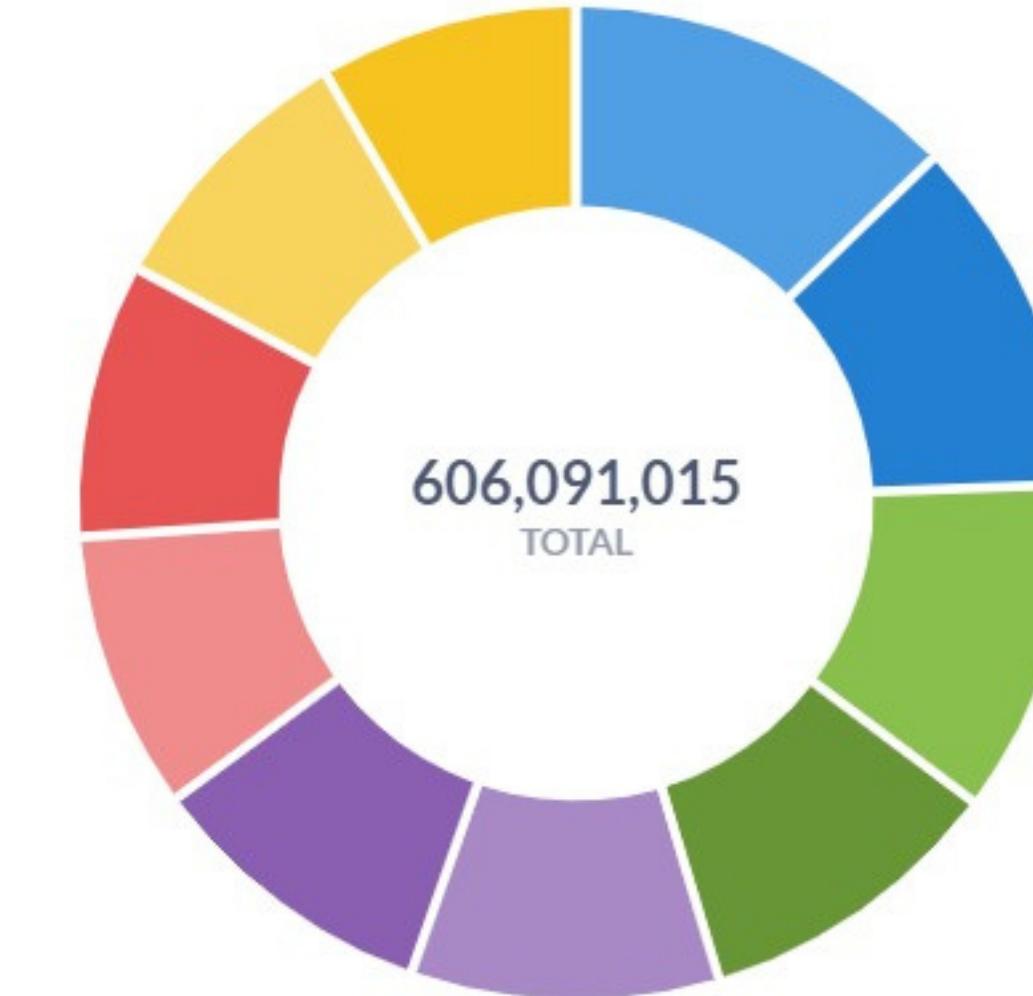


Hasil dari visualisasi terlihat bahwa siklus dari berhasil login dan tidak berhasil mengalami naik turun berdasarkan perilaku pengguna dalam menggunakannya.

Visualisasi

Amount Money Supplier

Supplier	Persentase
Cummings-Patterson	12.85%
Frazier, Peck and Saunders	11.66%
King-Anderson	10.89%
Cooley, Rivera and Greer	10.10%
Espinosa-Morrison	10.03%
Leblanc-Mann	9.49%
Pacheco-Garcia	9.01%
Collier-Sellers	8.86%
Ball Inc	8.75%
Smith, Shields and Parker	8.35%



Berdasarkan grafik 10 pemasok dana teratas, dapat disimpulkan bahwa Cummings-Patterson adalah pemasok dana terbesar, dengan pangsa aset sebesar 12,85%, diikuti oleh Frazier, Peck dan Saunders dengan pangsa aset sebesar 11,66%, dan King-Anderson dengan pangsa aset sebesar 10,89%.



Thank You

Kelompok 13