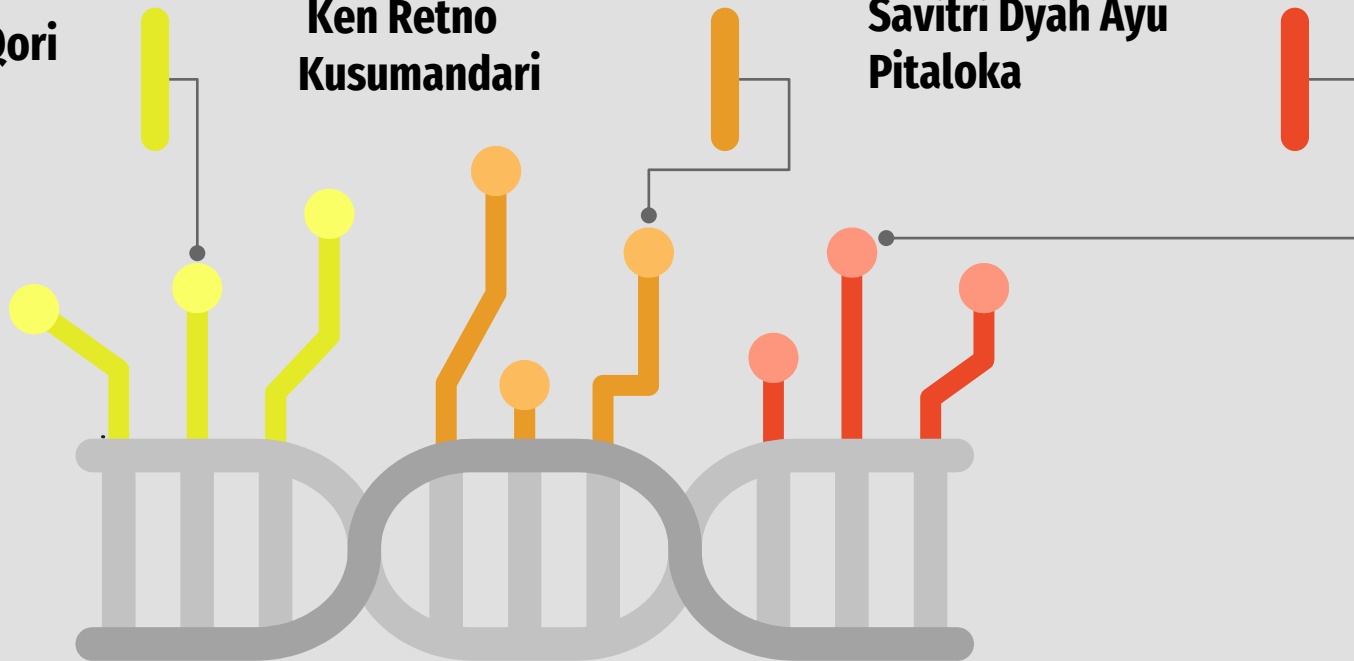# FINAL PROJECT

Kelompok 6
Fast Track Data Engineer
Digital Skola

# Nama Anggota
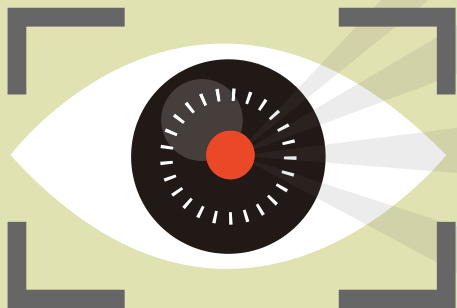
Ilham Abdi Al Qori

Ken Retno Kusumandari

Savitri Dyah Ayu Pitaloka

**Outline**

**Project 1** — Python and Business Intelligence

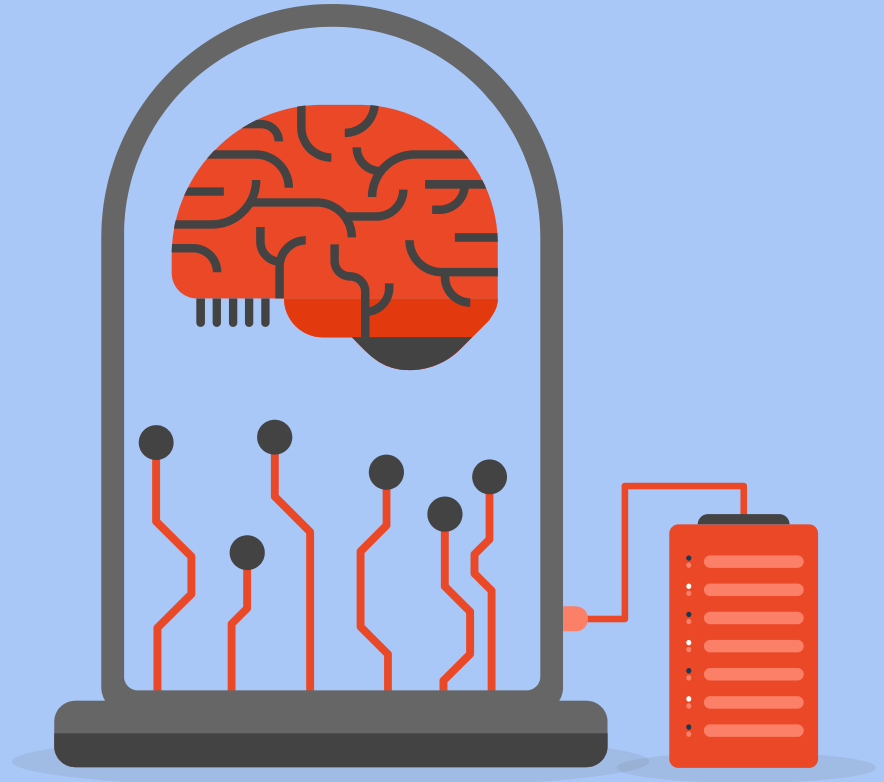**Project 2** — Integration Batch Projection with DBT

**Project 3** — Big Data Processing
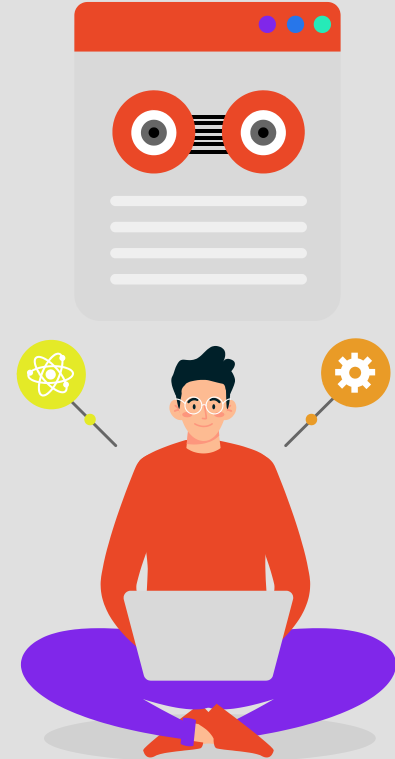
**Project 4** — Real-Time Processing
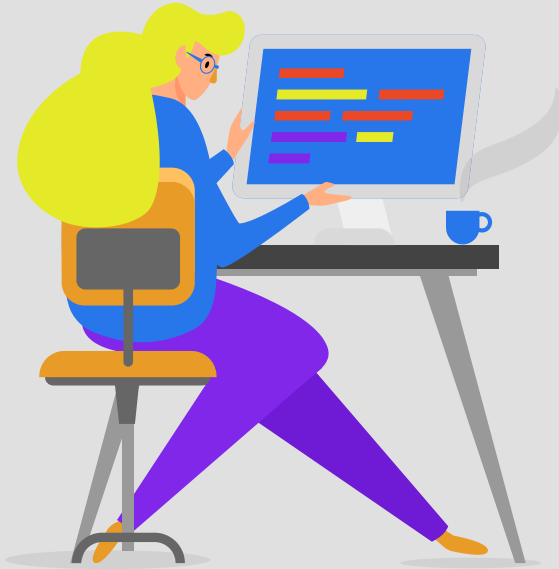
# PROJECT 1

# Objective of The Project

**Creating a Business Intelligence Report with Python and PostgreSQL**

In this project, our goal is to gain practical insights into business intelligence (BI) by creating a comprehensive report. Specifically, we aim to achieve the following: Database Connection and Exploration Datamart Creation Visualization with Google Data Studio

# Pipeline

**01** Create Connection in Dbeaver
Postgre Connection for OLTP and DWH

**02** Understand Business Requirements
Carry out to determine the DWH design

**03** Create Relational Data Warehouse Design
In ERD form

**04** Create DDL Data Warehouse and Data Mart
Based on Business Requirements
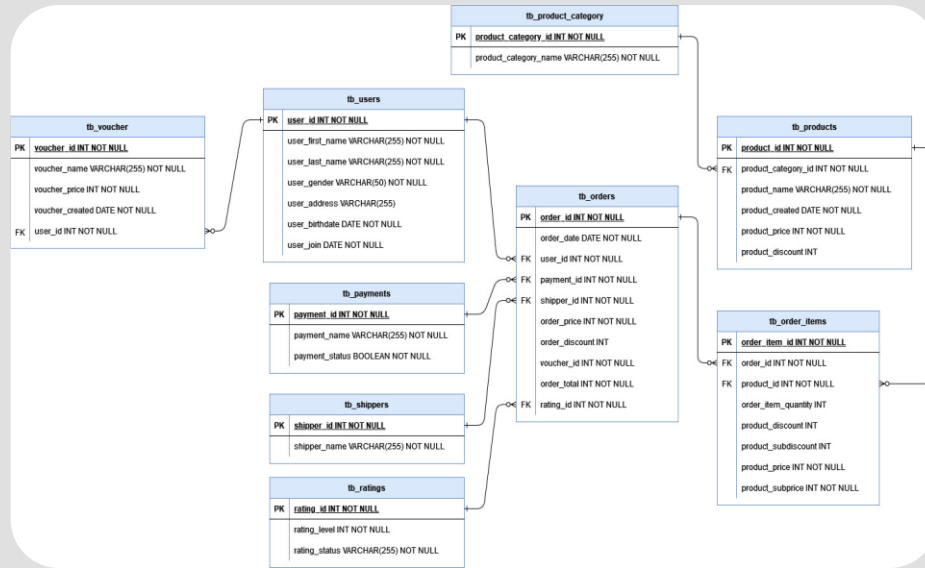
**05** Create ETL Script
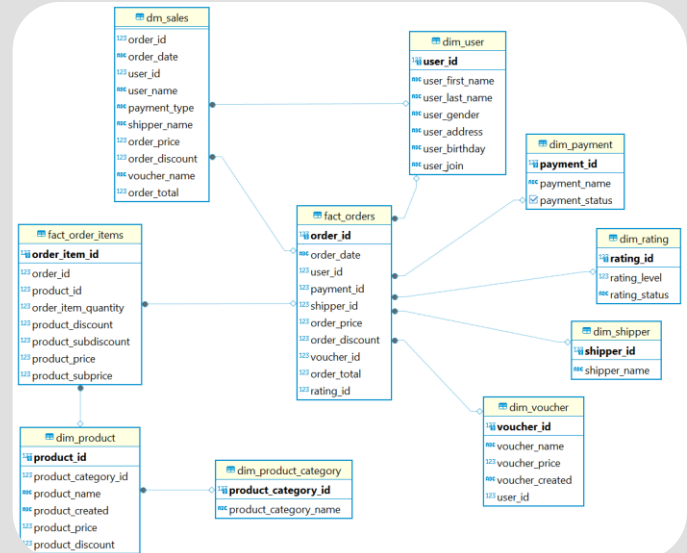Migrate data from OLTP to DWH

**06** Design the Dashboard
Customize Business Requirements and create in Looker Studio

# ERD OLTP and DWH

# Config and ETL Scripts

## Config

```python
#!python config       You, last month • config

oltp_conn_string = "postgresql://ftde01:ftde01-digitalskola@35.240.216.24:5432/ftde01_oltp"
warehouse_conn_string = "postgresql://ftde01:ftde01-digitalskola@35.240.216.24:5432/ftde01_dwh"

oltp_tables = {
    "users": "tb_users",
    "payments": "tb_payments",
    "shippers": "tb_shippers",
    "ratings": "tb_ratings",
    "vouchers": "tb_vouchers",
    "product_category": "tb_product_category",
    "products": "tb_products",
    "orders": "tb_orders",
    "order_items": "tb_order_items"
}

warehouse_tables = {
    "users": "dim_user",
    "payments": "dim_payment",
    "shippers": "dim_shipper",
    "ratings": "dim_rating",
    "vouchers": "dim_voucher",
    "product_category": "dim_product_category",
    "products": "dim_product",
    "orders": "fact_orders",
    "order_items": "fact_order_items"
}
```

## ETL

```python
#!python ETL       You, last month • etl

import pandas as pd
import sqlalchemy as sa

from config import oltp_conn_string, warehouse_conn_string, oltp_tables, warehouse_tables, dimension_colum

def create_tables():
    """Create tables in the data warehouse if they do not exist."""
    engine = sa.create_engine(warehouse_conn_string)
    with engine.connect() as conn:
        for ddl in ddl_statements.values():
            conn.execute(ddl)

def extract_data(table_name):
    """Extract data from a table in the OLTP database."""
    engine = sa.create_engine(oltp_conn_string)
    query = f"SELECT * FROM {oltp_tables[table_name]}"
    df = pd.read_sql(query, engine)
    print(f'Extract Data {oltp_tables[table_name]} Success')
    return df

def transform_data(df, target_table):
    """Transform the extracted data to match the schema of the target dimension table."""
    columns = dimension_columns.get(target_table)
    if columns:
        df = df[columns]
    print(f'Transform Data {target_table} Success')
```

# ETL Success

# Dashboard
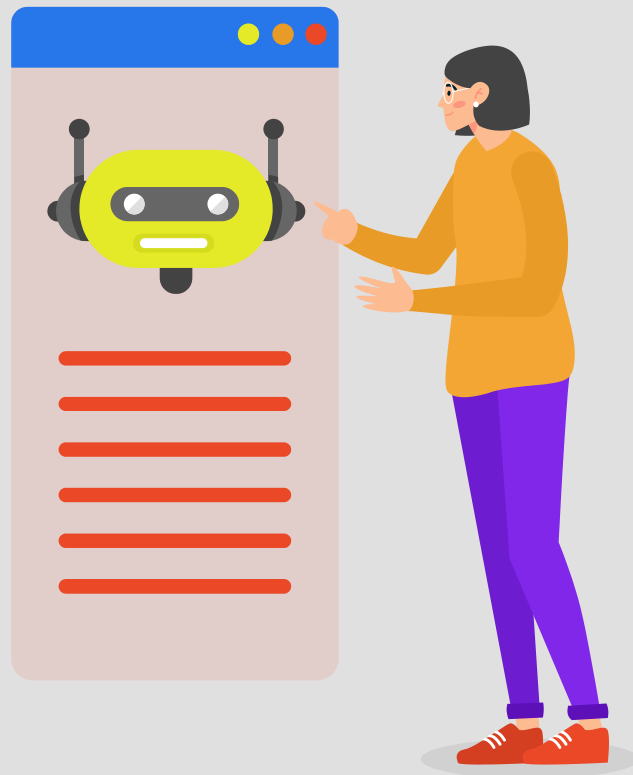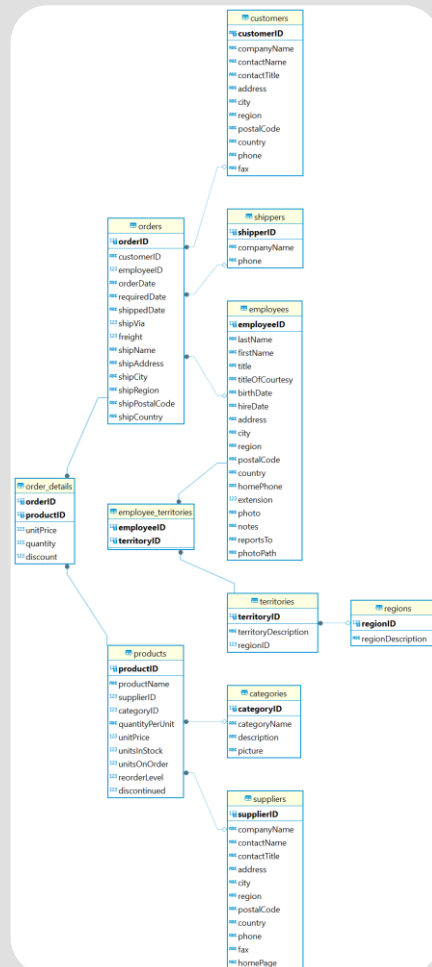
PROJECT 2

# Objective of The Project

**Leveraging dbt for Transforming Data in an ELT Workflow**

In this project, we aim to achieve a deep understanding of dbt (Data Build Tool) as a powerful transformation tool within the Extract, Load, Transform (ELT) methodology. Specifically, our objectives include:Data ExtractionTransformation with dbtIntegration with PostgreSQL
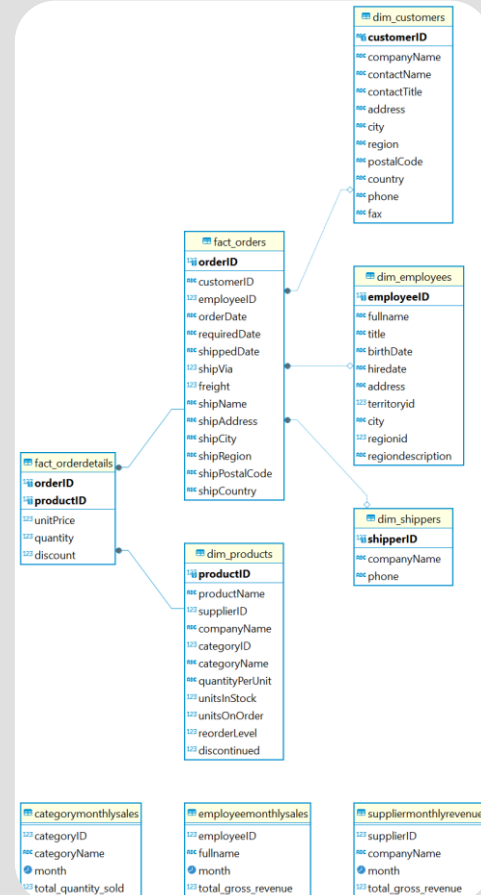
**OLTP**

**DWH**

# Data Mart

PROJECT 3

# Objective of The Project

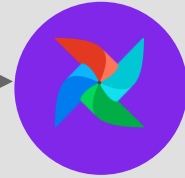Understanding and Implementing an End-to-End Data Pipeline with Apache Airflow, Docker, and PostgreSQL

# Pipeline



Source data from
kaggle

Insert raw data
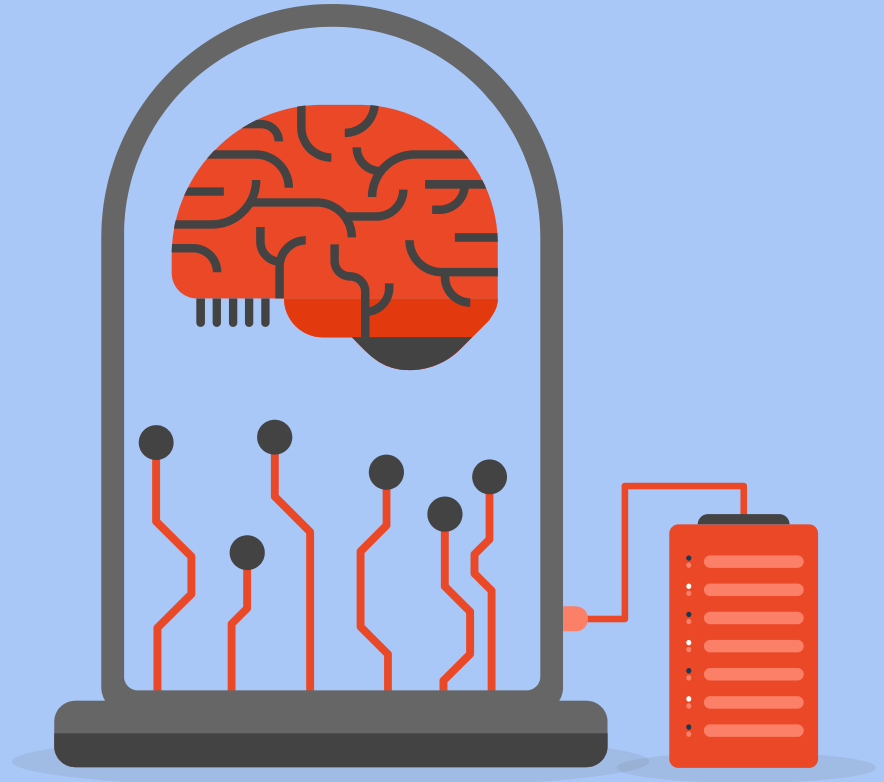from source into
postgre and create
dwh

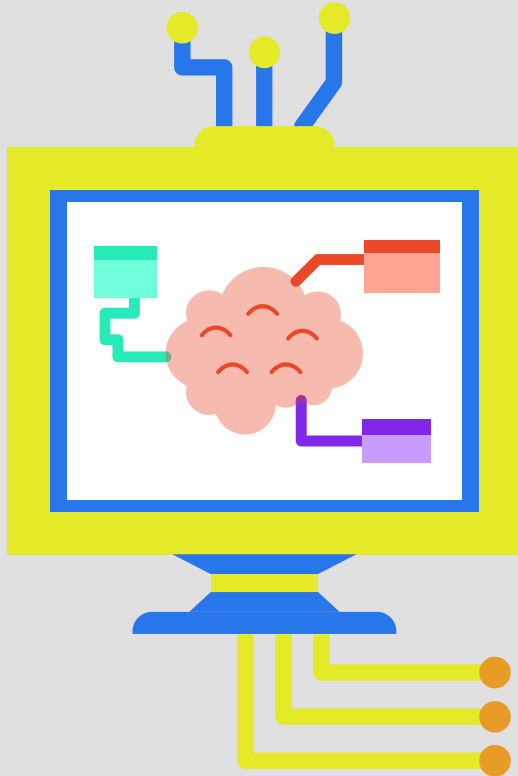Create dag airflow
with docker

# DAG AIRFLOW

# PROJECT 4

# Pipeline

**Producer**



**Consumer**

# Model Result

# Model Result in CSV

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | _id | step | type | amount | nameOrig | newbalanceOrig | nameDest | newbalanceDest | oldbalanceOrg | oldbalanceDest | predictionResult |
| 2 | 66910e0ae44630bf9f92c6bd | 1 | PAYMENT | 9839.64 | C1231006815 | 160296.36 | M1979787155 | 0 | 170136 | 0 | White List |
| 3 | 66910e0fe44630bf9f92c6be | 1 | PAYMENT | 1864.28 | C1666544295 | 19384.72 | M2044282225 | 0 | 21249 | 0 | White List |
| 4 | 66910e14e44630bf9f92c6bf | 1 | TRANSFER | 181 | C1305486145 | 0 | C553264065 | 0 | 181 | 0 | White List |
| 5 | 66910e19e44630bf9f92c6c0 | 1 | CASH_OUT | 181 | C840083671 | 0 | C38997010 | 0 | 181 | 21182 | White List |
| 6 | 66910e1ee44630bf9f92c6c1 | 1 | PAYMENT | 11668.14 | C2048537720 | 29885.86 | M1230701703 | 0 | 41554 | 0 | White List |
| 7 | 66910e23e44630bf9f92c6c2 | 1 | PAYMENT | 7817.71 | C90045638 | 46042.29 | M573487274 | 0 | 53860 | 0 | White List |
| 8 | 66910e28e44630bf9f92c6c3 | 1 | PAYMENT | 7107.77 | C154988899 | 176087.23 | M408069119 | 0 | 183195 | 0 | White List |
| 9 | 66910e2de44630bf9f92c6c4 | 1 | PAYMENT | 7861.64 | C1912850431 | 168225.59 | M633326333 | 0 | 176087.23 | 0 | White List |
| 10 | 66910e32e44630bf9f92c6c5 | 1 | PAYMENT | 4024.36 | C1265012928 | 0 | M1176932104 | 0 | 2671 | 0 | White List |
| 11 | 66910e37e44630bf9f92c6c6 | 1 | DEBIT | 5337.77 | C712410124 | 36382.23 | C195600860 | 40348.79 | 41720 | 41898 | White List |
| 12 | 66910e3ce44630bf9f92c6c7 | 1 | DEBIT | 9644.94 | C1900366749 | 0 | C997608398 | 157982.12 | 4465 | 10845 | White List |
| 13 | 66910e41e44630bf9f92c6c8 | 1 | PAYMENT | 3099.97 | C249177573 | 17671.03 | M2096539129 | 0 | 20771 | 0 | White List |
| 14 | 66910e46e44630bf9f92c6c9 | 1 | PAYMENT | 2560.74 | C1648232591 | 2509.26 | M972865270 | 0 | 5070 | 0 | White List |
| 15 | 66910e4be44630bf9f92c6ca | 1 | PAYMENT | 11633.76 | C1716932897 | 0 | M801569151 | 0 | 10127 | 0 | White List |
| 16 | 66910e50e44630bf9f92c6cb | 1 | PAYMENT | 4098.78 | C1026483832 | 499165.22 | M1635378213 | 0 | 503264 | 0 | White List |
| 17 | 66910e55e44630bf9f92c6cc | 1 | CASH_OUT | 229133.94 | C905080434 | 0 | C476402209 | 51513.44 | 15325 | 5083 | White List |
| 18 | 66910e5ae44630bf9f92c6cd | 1 | PAYMENT | 1563.82 | C761750706 | 0 | M1731217984 | 0 | 450 | 0 | White List |
| 19 | 66910e5fe44630bf9f92c6ce | 1 | PAYMENT | 1157.86 | C1237762639 | 19998.14 | M1877062907 | 0 | 21156 | 0 | White List |
| 20 | 66910e64e44630bf9f92c6cf | 1 | PAYMENT | 671.64 | C2033524545 | 14451.36 | M473053293 | 0 | 15123 | 0 | White List |
| 21 | 66910e69e44630bf9f92c6d0 | 1 | TRANSFER | 215310.3 | C1670993182 | 0 | C1100439041 | 0 | 705 | 22425 | White List |

# Thank You