# Real Time E-Commerce Data Processing

1. created a real-time data pipeline for processing e-commerce data using Apache Kafka and Apache Cassandra. I was ingesting data from a CSV file using a Kafka producer, transforming the data using a Kafka consumer, and finally storing the processed data in a Cassandra table.



2. So first loaded the 'olist_orders_dataset.csv' into a pandas dataframe and examined its structure and contents.

3. configured Confluent Kafka & created a Kafka topic, named 'ecommerce-orders', to hold the e-commerce data.

4. developed a kafka idempotence producer(idempotence to prevent duplicates) in python that reads the data from the pandas dataframe and publishes it to the 'ecommerce-orders' Kafka topic. The key for each message was a combination of the 'customer_id' and 'order_id' fields from the dataset.



5. Set up Datastax Cassandra & created a keyspace, named 'ecommerce', for storing the e-commerce data.

```python
from cassandra.cluster import Cluster
from cassandra.auth import PlainTextAuthProvider
import json

# This secure connect bundle is autogenerated when you download your SCB,
# if yours is different update the file name below
cloud_config= {
    'secure_connect_bundle': 'secure-connect-ecommerce-db.zip'
}

# This token JSON file is autogenerated when you download your token,
# if yours is different update the file name below
with open("ecommerce_db-token.json") as f:
    secrets = json.load(f)

CLIENT_ID = secrets["clientId"]
CLIENT_SECRET = secrets["secret"]

auth_provider = PlainTextAuthProvider(CLIENT_ID, CLIENT_SECRET)
cluster = Cluster(cloud=cloud_config, auth_provider=auth_provider)
session = cluster.connect()

row = session.execute("select release_version from system.local").one()
if row:
    print(row[0])
else:
    print("An error occurred.")
```

4.0.0.6816

```python
try:
    query = 'use ecommerce'
    session.execute(query)
    print('Inside the ecommerce')
except Exception as err:
    print('Exception Occured while using Keyspace : ',err)
```

6. Then designed a table, named 'orders', within the 'ecommerce' keyspace. This table reflected the schema of the incoming data and included additional columns for the derived features: 'OrderHour' and 'OrderDayOfWeek'. The data model had 'customer_id' as the partition key and 'order_id' and 'order_purchase_timestamp' as clustering keys.

```python
try:
    query = """CREATE TABLE orders (
                order_id uuid,
                customer_id uuid,
                order_status text,
                order_purchase_timestamp timestamp,
                order_approved_at timestamp,
                order_delivered_carrier_date timestamp,
                order_delivered_customer_date timestamp,
                order_estimated_delivery_date timestamp,
                OrderHour int,
                OrderDayOfWeek text,
                PRIMARY KEY ((customer_id), order_id,
                order_purchase_timestamp))
            """
    session.execute(query)
    print('Table created inside the keyspace')
except Exception as err:
    print('Exception Occured while creating the table : ',err)
```

```
Table created inside the keyspace
```

7. Developed a Kafka consumer (with having a consumer group) in Python that subscribes to the 'ecommerce-orders' topic. The consumer derived two new columns 'PurchaseHour'/'OrderHour' and 'PurchaseDayOfWeek'/'OrderDayOfWeek', then ingesting transformed data into the 'orders' table in Cassandra.

8. While inserting data into the Cassandra 'orders' table, ensured that the write operations maintain quorum consistency.
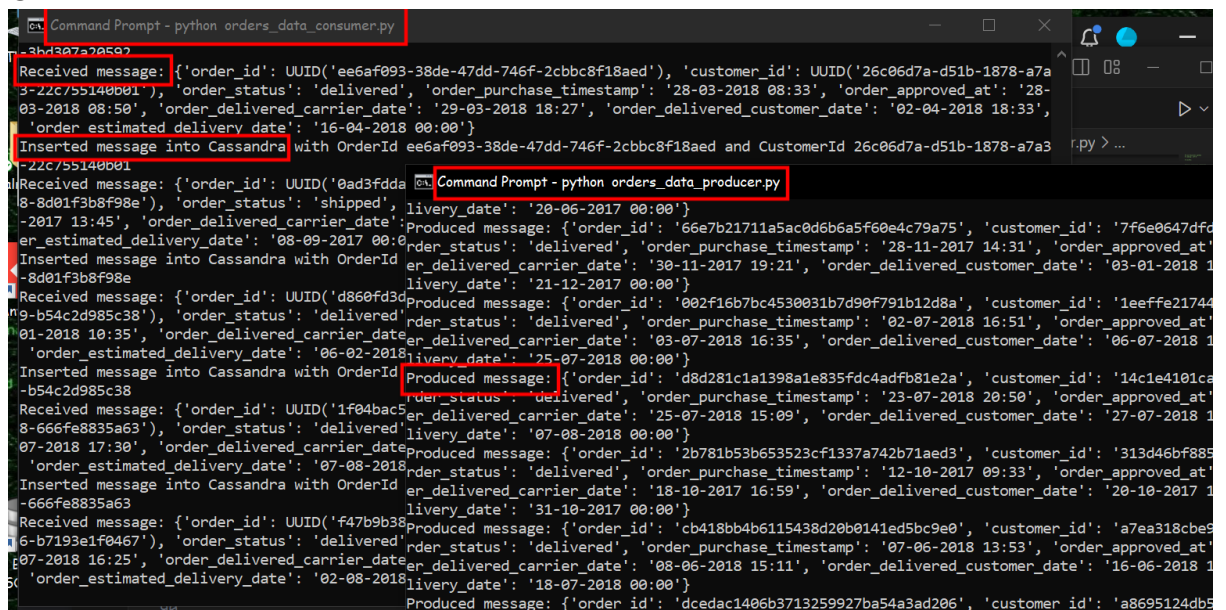


```
# Insert transformed data into Cassandra db
# Prepare the [Loading...] atement with dynamic values
insert_query [Loading...] prepare("INSERT INTO orders (order_id, customer_id, order_status, order_purchase_timestamp,
order_approved_at, order_delivered_carrier_date, order_delivered_customer_date, order_estimated_delivery_date, OrderHour,
OrderDayOfWeek) VALUES (?, ?, ?, ?, ?, ?, ?, ?, ?, ?)")

# Bind the dynamic values to the prepared statement
bound_statement = insert_query.bind((value['order_id'], value['customer_id'], value['order_status'], opt_to_datetime,
oaa_to_datetime, odcd_to_datetime, odcud_to_datetime, oedd_to_datetime, orderHour, orderDayOfWeek))

# Execute the insert statement with the specified consistency level
bound_statement.consistency_level = ConsistencyLevel.QUORUM
session.execute(bound_statement)

print(f"Inserted message into Cassandra with OrderId {value['order_id']} and CustomerId {value['customer_id']}")
```

9. Tested the data pipeline end-to-end. Ran Kafka producer to ingest the data, then executed the Kafka consumer to process the data and inserted it into the Cassandra table.



10. Verified the data in the Cassandra table matches the processed data and that all transformations have been executed correctly.