# E-Commerce Data Processing

**Deep Patel**

**October 2024**

# Big Data Problem

- **Handling large Volumes of data in Various form at high Velocity**

- **User Interaction log.**
  - **Approx.. Daily Log Size:**
  - **1,000,000 users × 100 interactions/user = 100,000,000 entries/day**
  - **Total daily log size: 100,000,000 entries × 80 bytes/entry = 8,000,000,000 bytes**
  - **8,000,000,000 bytes ÷ 1,024² = 7,629 MB/day (approximately 7.63 GB/day)**

1

# Business Use Case

- **Customer Activity by Region and Loyalty Tier**
- **Top Products by Category and Region**
- **Interaction Trend Over Time**
- **Historical Purchase Trends**
- **Customer Age and Interaction Analysis**

# Tech Stack

Apache Spark for batch processing.

Hive for metadata management.

HDFS for storage.

PostgreSQL for source data.

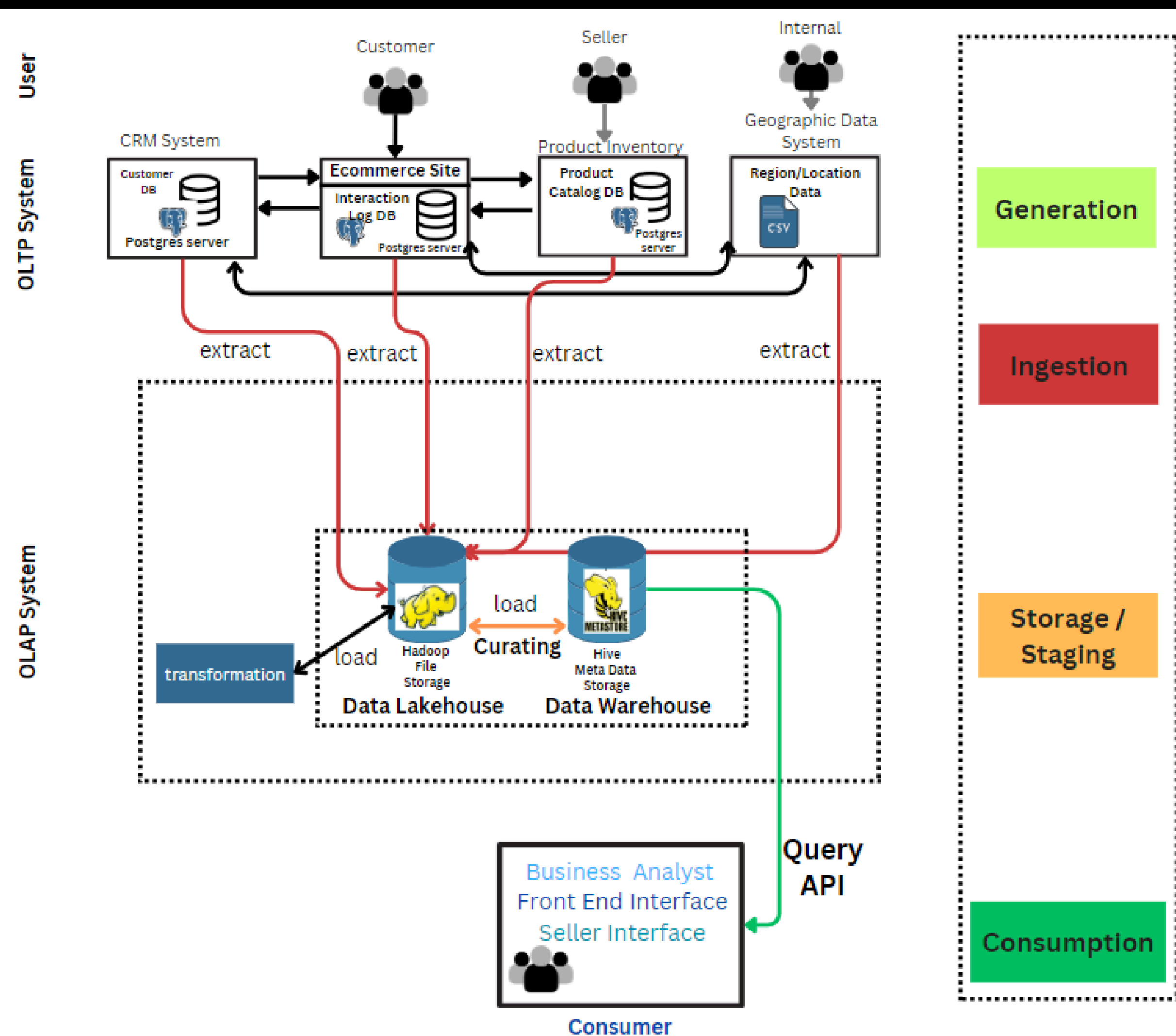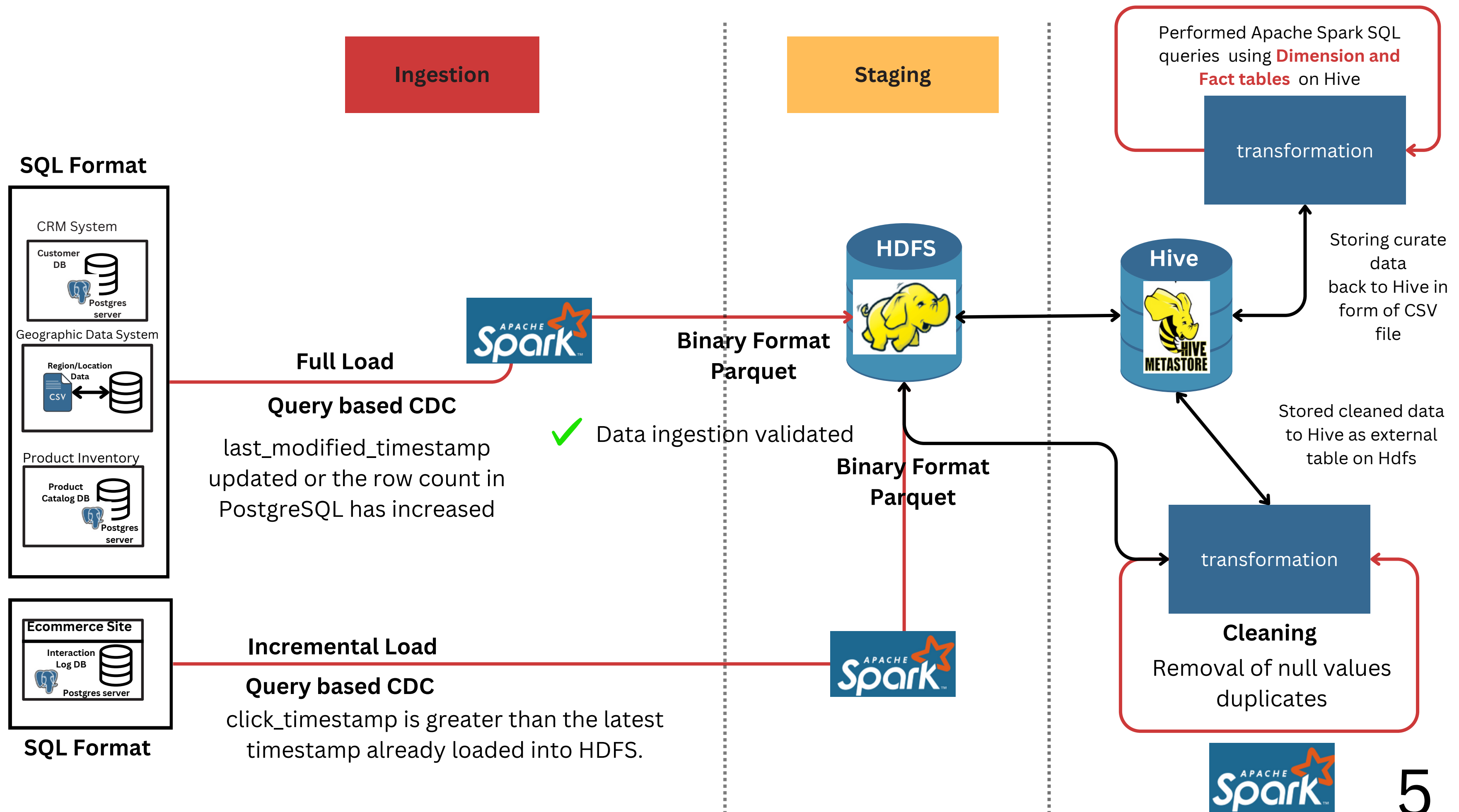matplotlib, seaborn (visualization)
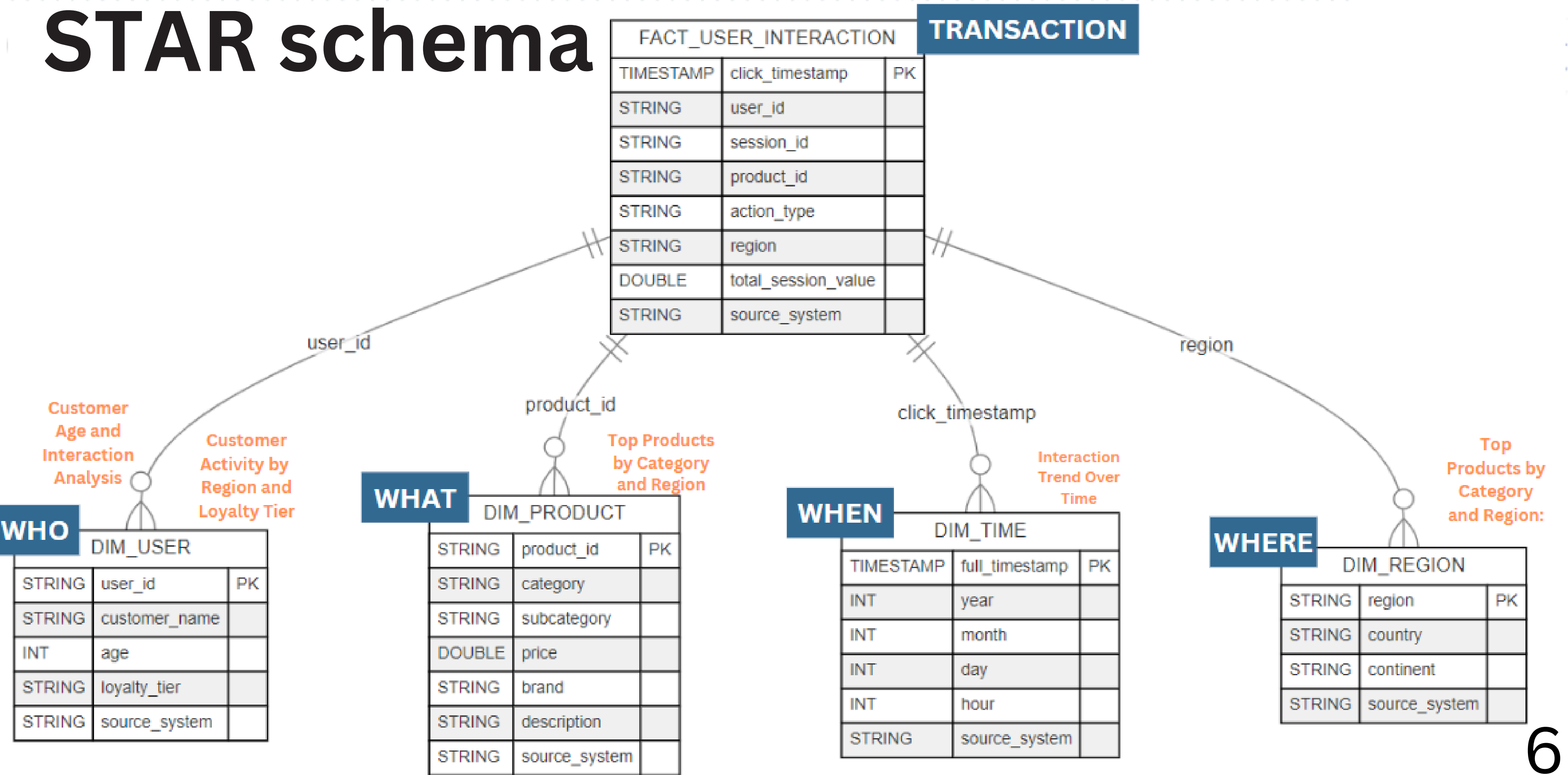
**Key Integrations:**

Batch processing

Query-based CDC

- Full and Incremental Loads.

Dimensional Modeling (Star Schema).

**Ingestion**

**Staging**
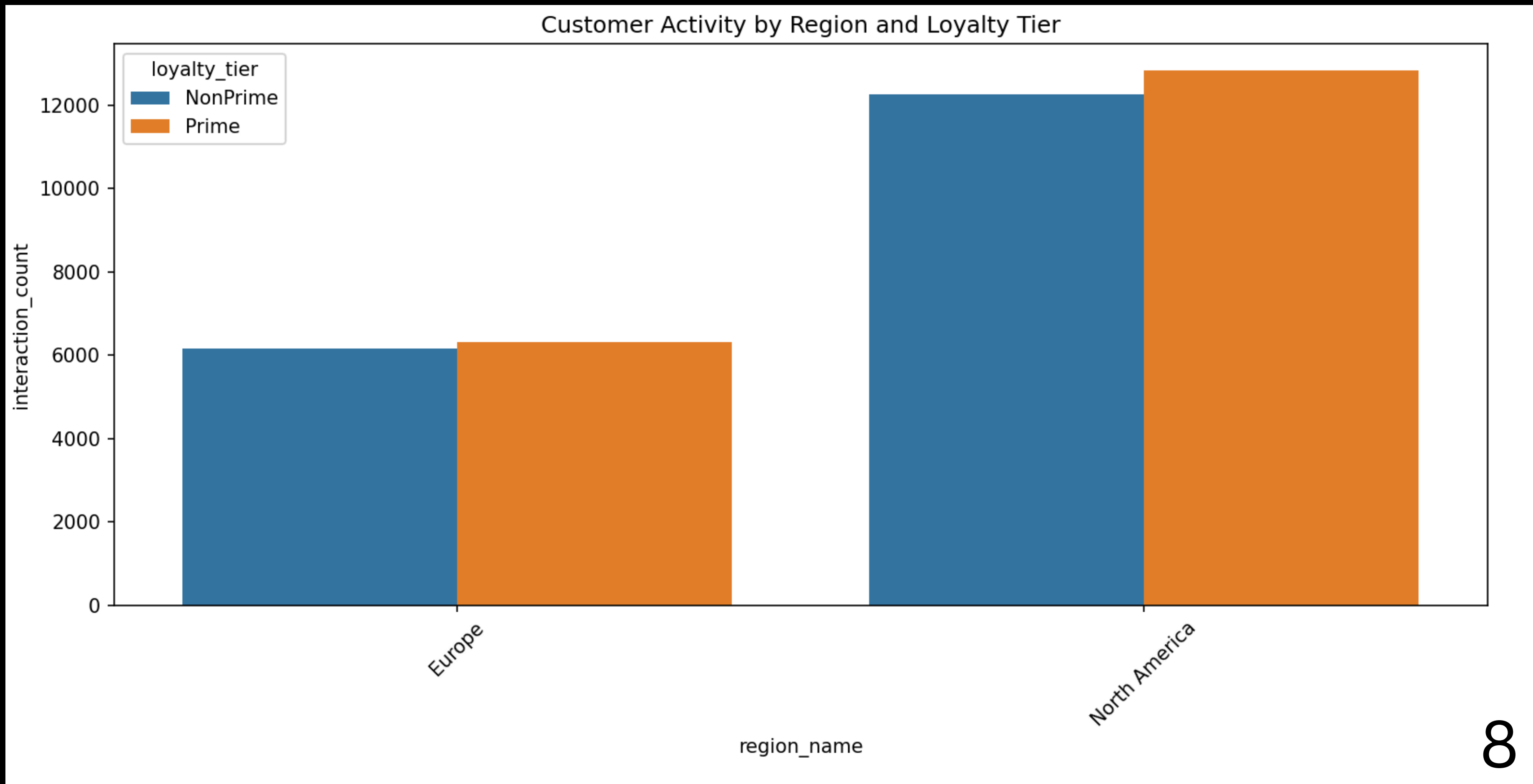
Performed Apache Spark SQL queries using **Dimension and Fact tables** on Hive

transformation

**SQL Format**

CRM System

Customer DB
Postgres server

Geographic Data System

Region/Location Data
CSV

Product Inventory

Product Catalog DB
Postgres server

**SQL Format**

**Full Load**

**Query based CDC**

last_modified_timestamp updated or the row count in PostgreSQL has increased

**HDFS**

**Binary Format Parquet**

✅ Data ingestion validated

**Binary Format Parquet**

**Hive**

HIVE METASTORE

Storing curate data back to Hive in form of CSV file

Stored cleaned data to Hive as external table on Hdfs

transformation

**Cleaning**
Removal of null values duplicates

Ecommerce Site

Interaction Log DB
Postgres server

**Incremental Load**

**Query based CDC**

click_timestamp is greater than the latest timestamp already loaded into HDFS.

**SQL Format**

5

# STAR schema



**TRANSACTION**

**FACT_USER_INTERACTION**

| | | |
|---|---|---|
| TIMESTAMP | click_timestamp | PK |
| STRING | user_id | |
| STRING | session_id | |
| STRING | product_id | |
| STRING | action_type | |
| STRING | region | |
| DOUBLE | total_session_value | |
| STRING | source_system | |

user_id

product_id

click_timestamp

region

*Customer Age and Interaction Analysis*

*Customer Activity by Region and Loyalty Tier*

*Top Products by Category and Region*

*Interaction Trend Over Time*

*Top Products by Category and Region:*

**WHO**

**DIM_USER**

| | | |
|---|---|---|
| STRING | user_id | PK |
| STRING | customer_name | |
| INT | age | |
| STRING | loyalty_tier | |
| STRING | source_system | |

**WHAT**

**DIM_PRODUCT**

| | | |
|---|---|---|
| STRING | product_id | PK |
| STRING | category | |
| STRING | subcategory | |
| DOUBLE | price | |
| STRING | brand | |
| STRING | description | |
| STRING | source_system | |

**WHEN**

**DIM_TIME**

| | | |
|---|---|---|
| TIMESTAMP | full_timestamp | PK |
| INT | year | |
| INT | month | |
| INT | day | |
| INT | hour | |
| STRING | source_system | |

**WHERE**

**DIM_REGION**

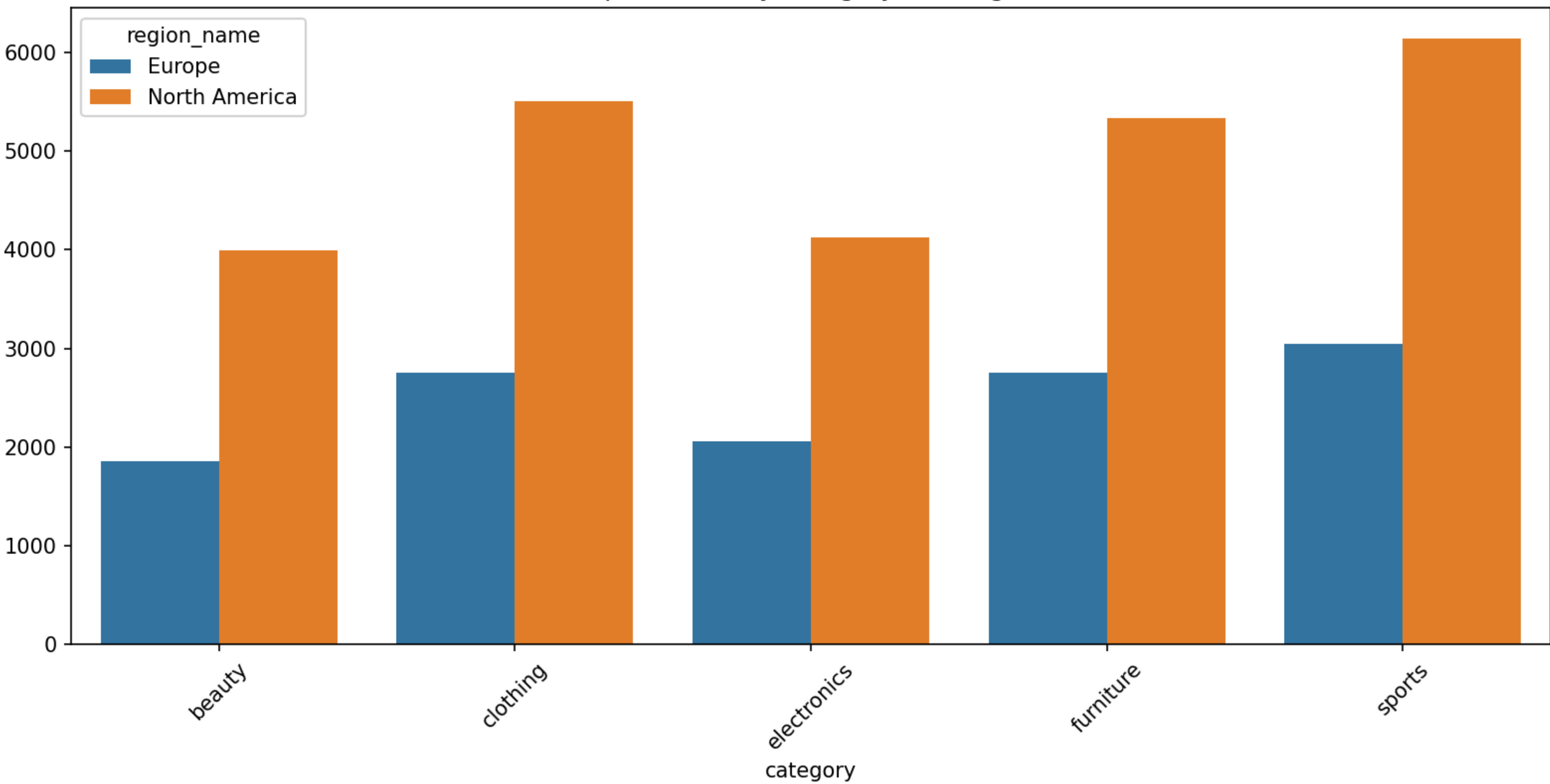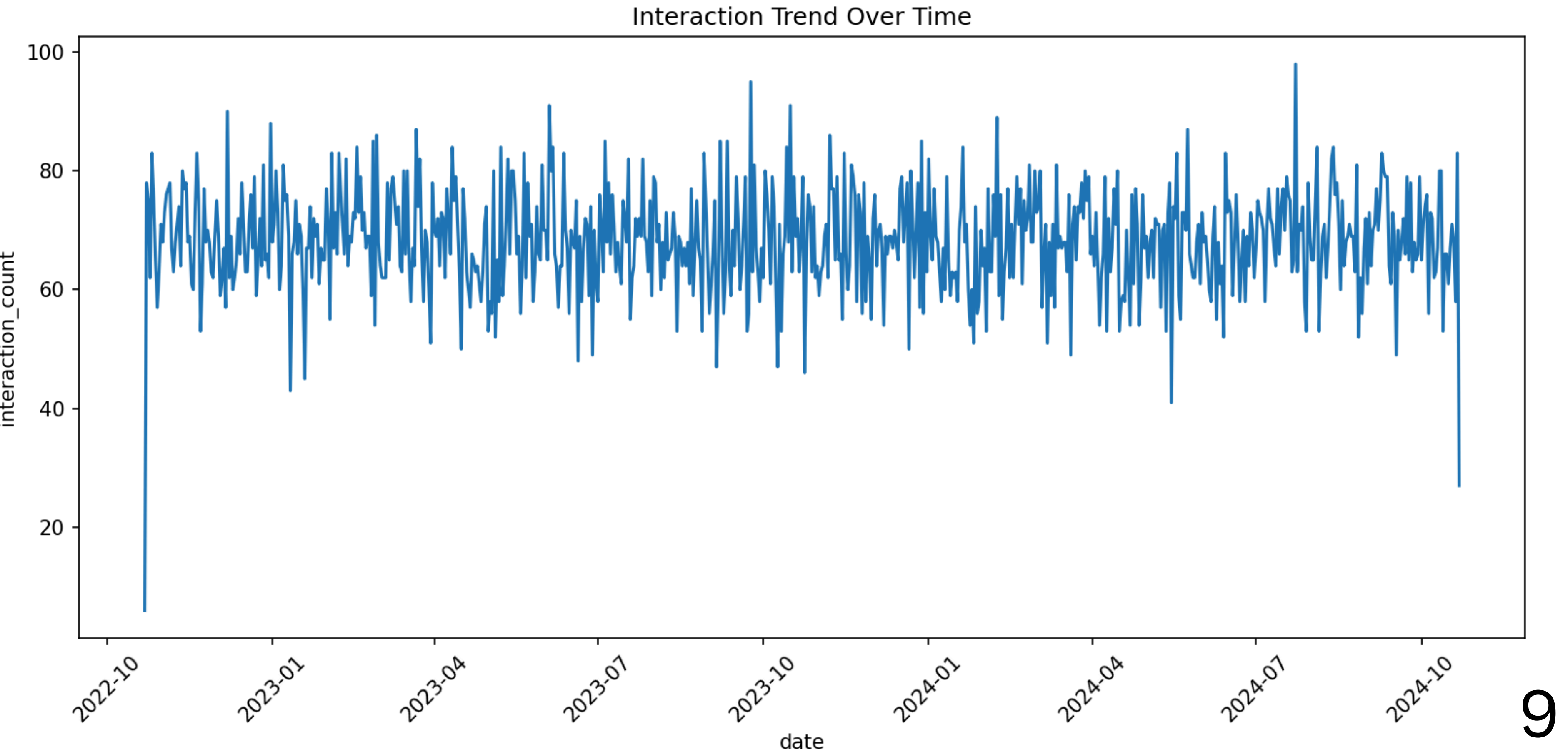| | | |
|---|---|---|
| STRING | region | PK |
| STRING | country | |
| STRING | continent | |
| STRING | source_system | |

6

# Analysis

- **Customer Activity by Region and Loyalty Tier**
- **Top Products by Category and Region**
- **Interaction Trend Over Time**
- **Historical Purchase Trends**
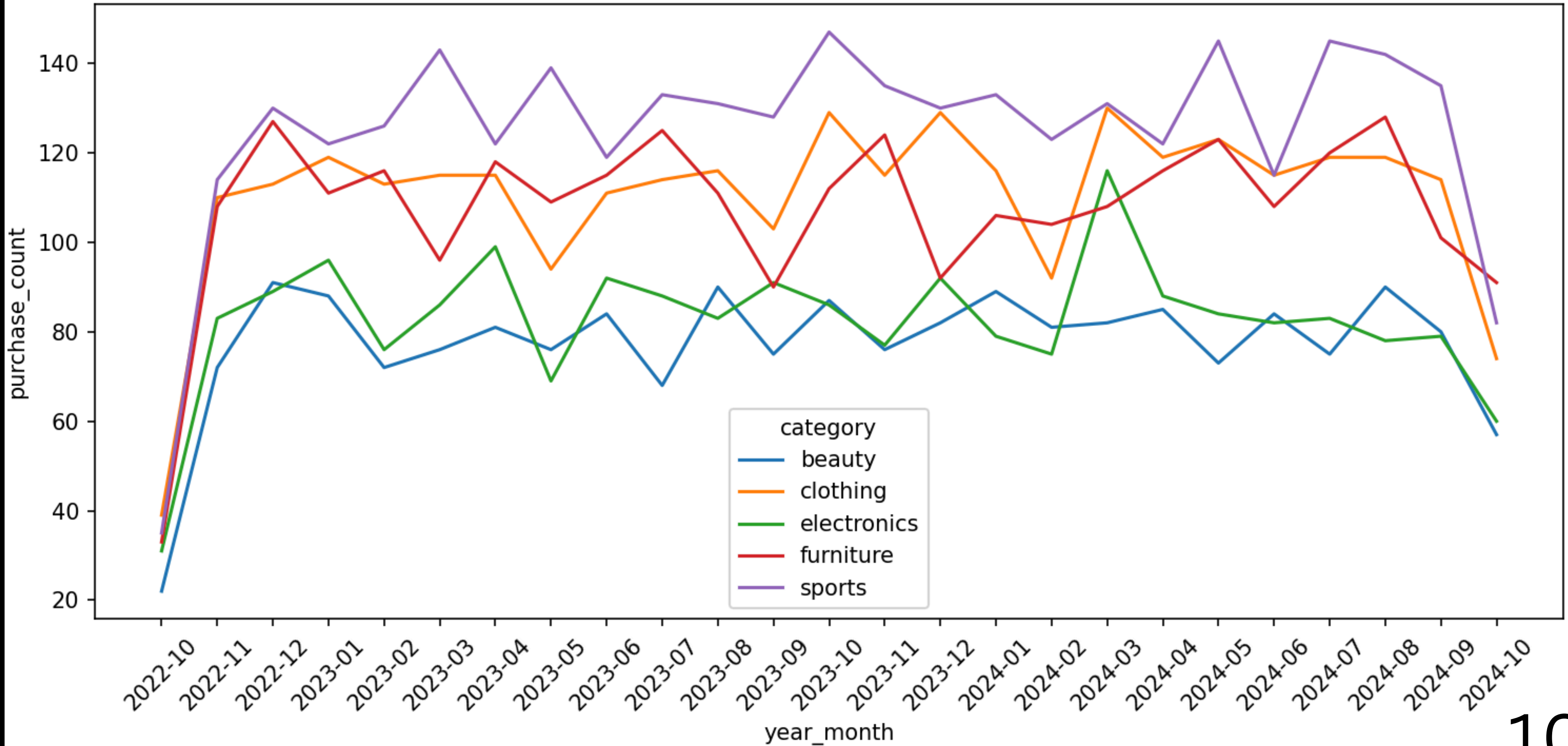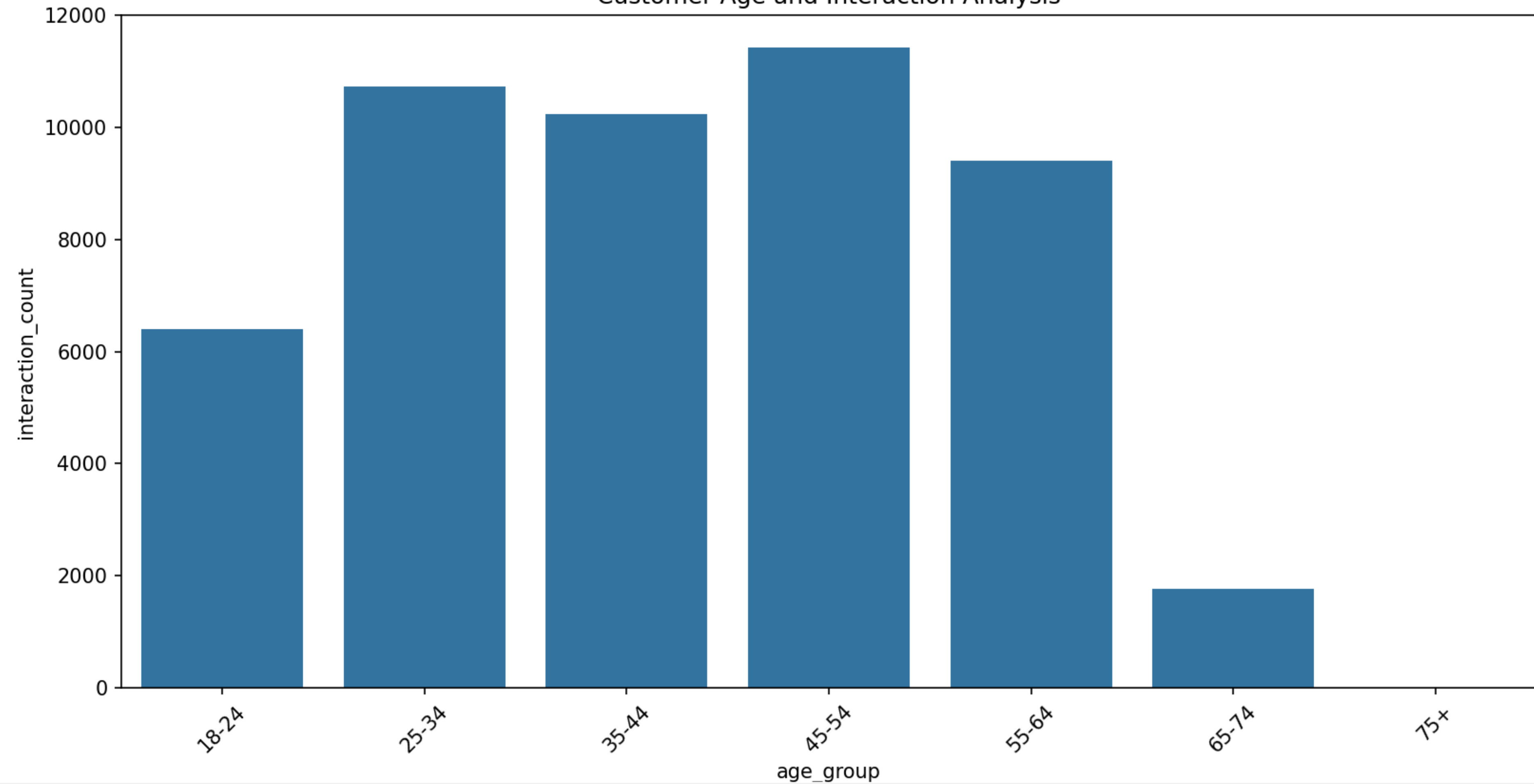- **Customer Age and Interaction Analysis**

Customer Activity by Region and Loyalty Tier

Top Products by Category and Region

Interaction Trend Over Time

9

Historical Purchase Trends by Category

10

Customer Age and Interaction Analysis

## Challenges:

- Performing Incremental load and full load using CSV file source - Solution: Use Postgres in Data Source Layer
- Encounter issues during the loading clean data to my dimension and fact table
- Permission relation issue in file path (Hdfs file path)
- Transformation ( Time stamp format mis-match can give invalid result in extraction )

12

# Future Extention

- Visualization - through power BI or tableau
- Kafka - (Real-time capturing processing user
  Interaction log)
- Automate the data pipeline  using Airflow
  transformation between layers
- Docker
- Airflow
- HBase (NoSQL) - Storing JSON based curated data

# QA

# Thank you