



GROUP PROJECT  
MAHARISHI INTERNATIONAL UNIVERSITY

# BIG DATA TECHNOLOGY FINAL PROJECT

M I U   U N I V E R S I T Y



MAHARISHI INTERNATIONAL UNIVERSITY

# MEET OUR TEAM



**BINIAM GHEBREMEDHIN**  
SOFTWARE ENGINEER



**NATSSNET GHIRMAY**  
SOFTWARE ENGINEER



**HIDRI GHEBRENTSÄE**  
SOFTWARE ENGINEER



# INTRODUCTION

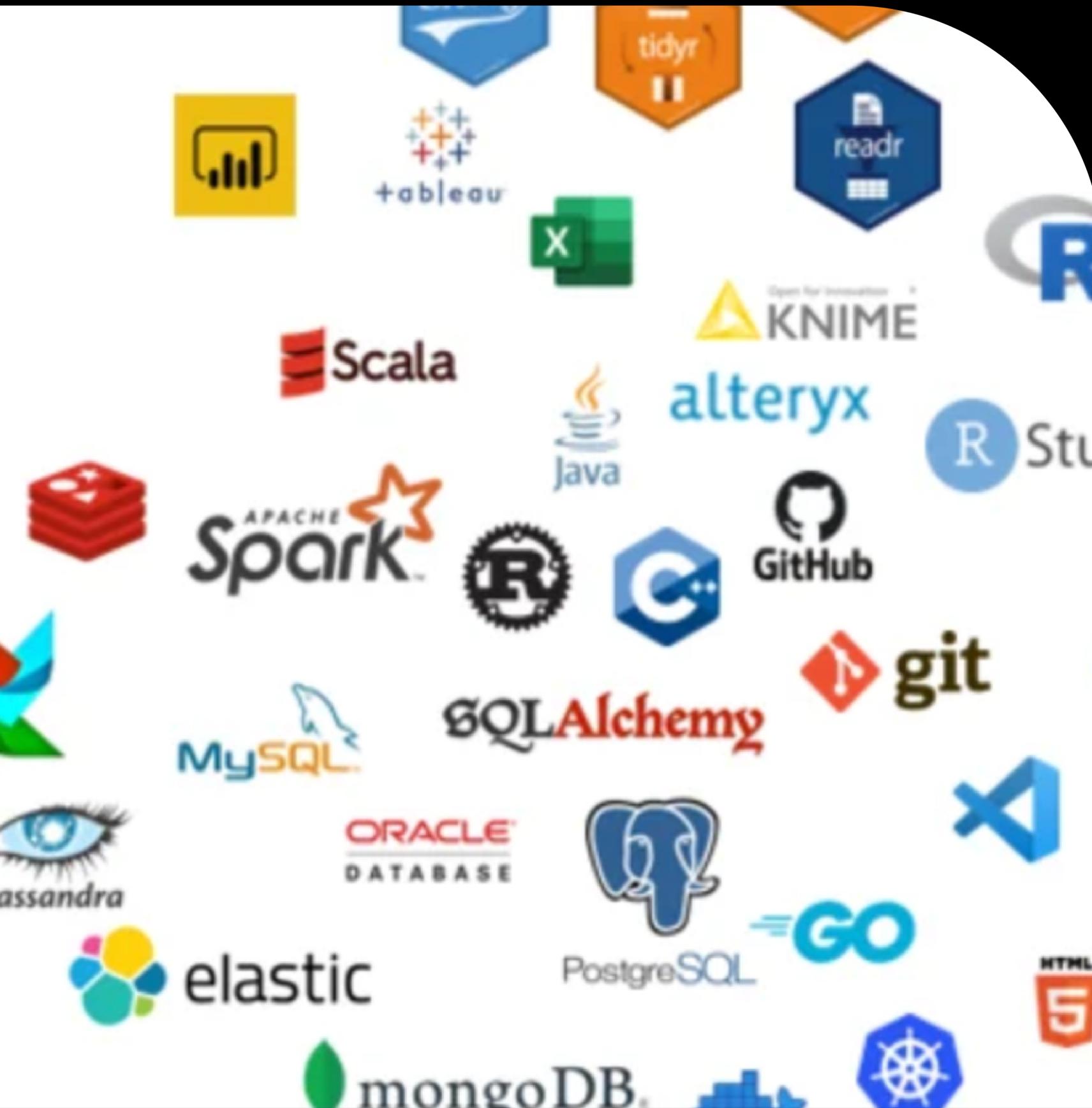
This project is all about the design and implementation of a process to process real-time data from an eCommerce store and visualization of the processed data.

It involves the processing and analyzing of the data and storing it in a scalable distributed database model like HBase for further usage of the data.

The data analysis is done by sparkSQL by executing queries and then the processed data is stored in SQL and NoSQL DBs.

# SUPERCHARGE YOUR E-COMMERCE ANALYTICS WITH PYSPARK





## TECHNOLOGIES USED

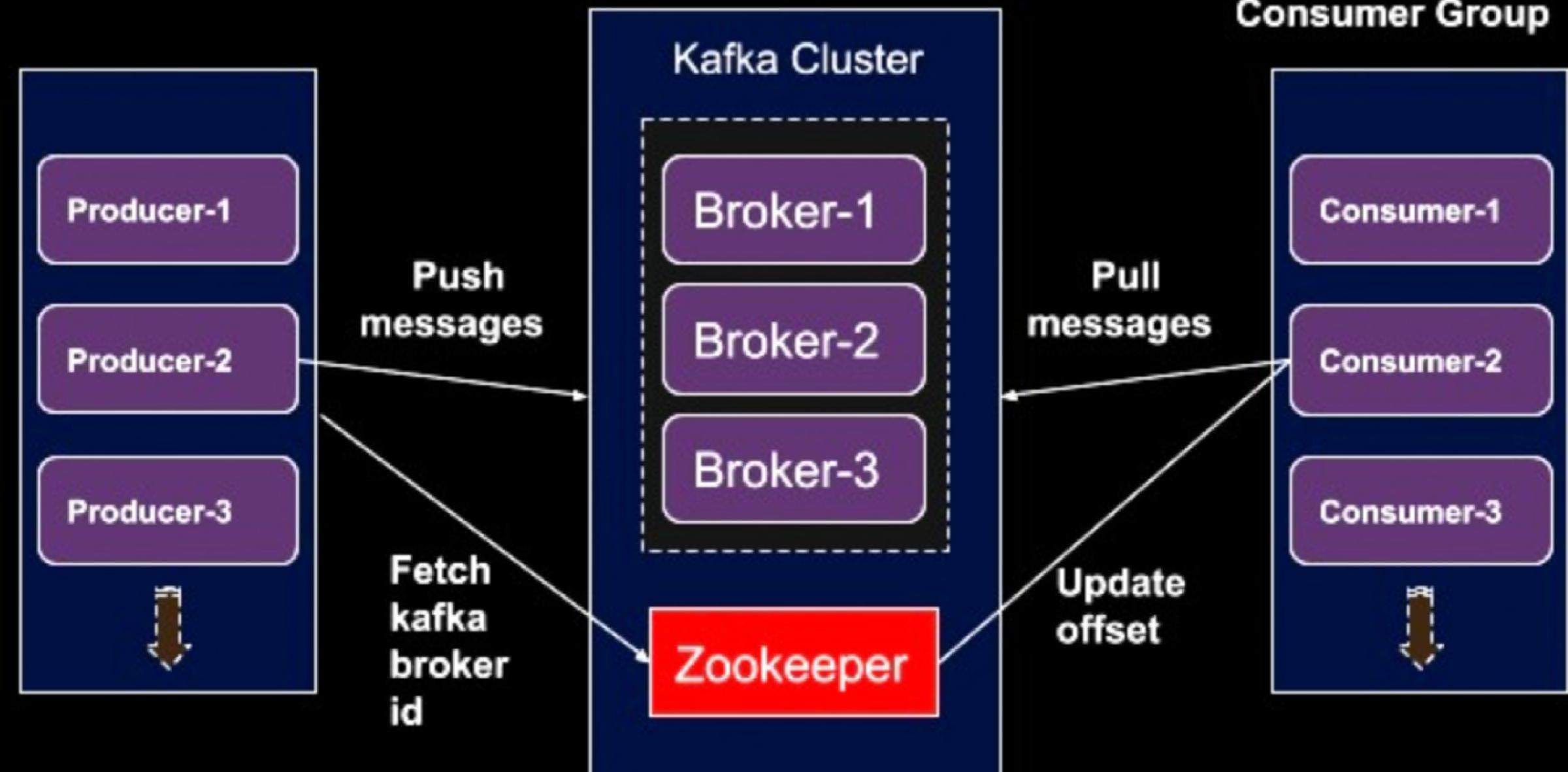
KAFKA  
HBASE  
SPARKSQL  
SPARAK STREAMING  
HIVE  
HADOOP  
DOCKER  
JAVA  
PYTHON  
INTELLIJ  
VSCODE  
ECLIPSE

LEARN MORE



# KAFKA

## Kafka Architecture



Apache Kafka is an open-source stream processing platform developed by LinkedIn and donated to the Apache Software Foundation. It is designed for high-throughput, fault-tolerant, and scalable messaging systems. Kafka is commonly used for building real-time data pipelines and streaming applications.

Kafka's architecture consists of producers, topics, brokers, partitions, and consumers. Producers publish messages to topics, which are partitions of data. Brokers manage the storage and communication of these topics, and consumers subscribe to topics to receive the messages.



DOCKER IS A PLATFORM FOR DEVELOPING, SHIPPING, AND RUNNING APPLICATIONS IN CONTAINERS. CONTAINERS ARE LIGHTWEIGHT, PORTABLE UNITS THAT PACKAGE APPLICATIONS AND THEIR DEPENDENCIES TOGETHER, ENSURING CONSISTENCY ACROSS DIFFERENT ENVIRONMENTS. HERE'S A BRIEF EXPLANATION OF DOCKER AND ITS KEY USES:

01

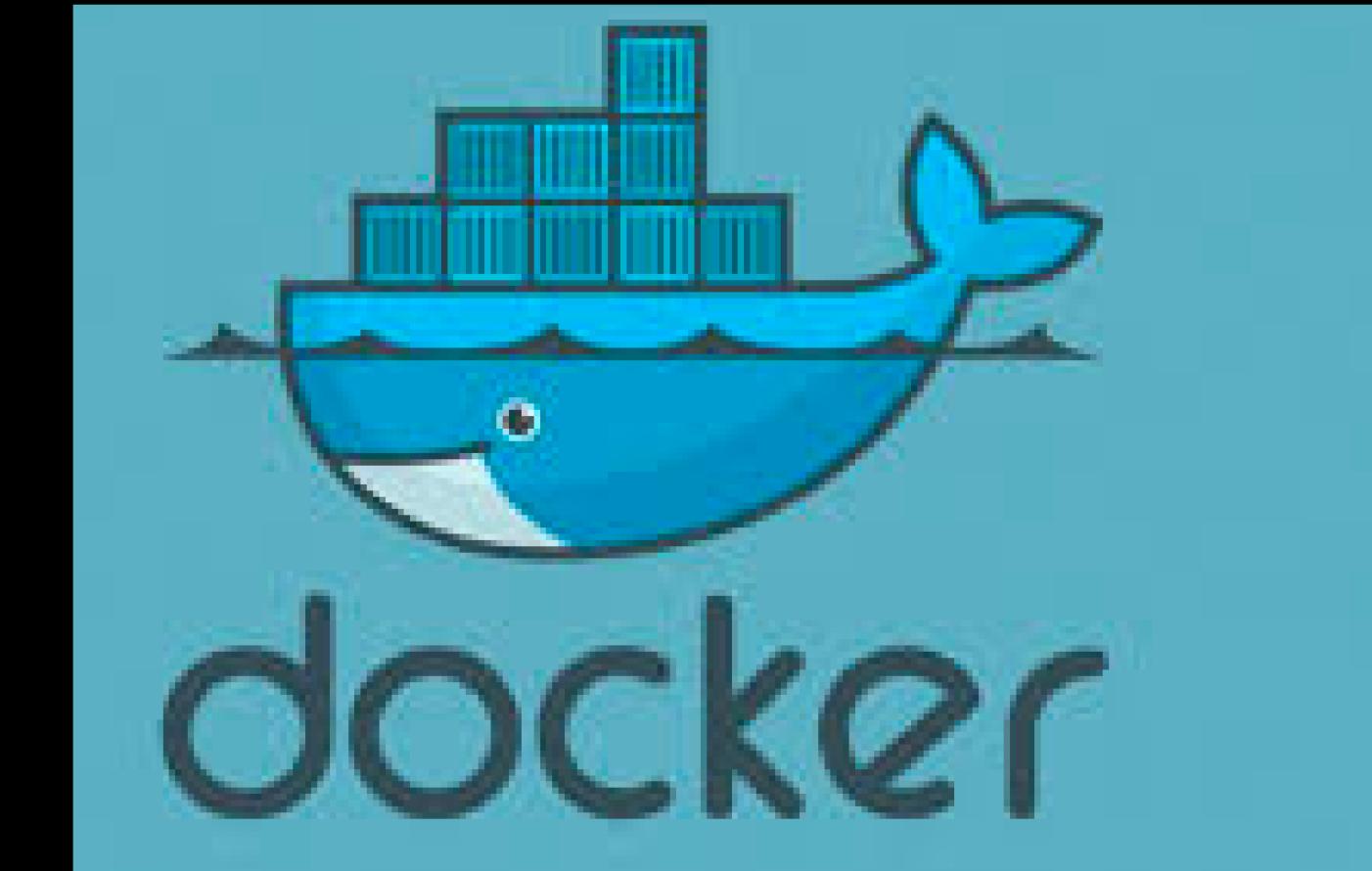
Isolation: Docker provides application isolation by encapsulating them in containers, preventing conflicts between applications and their dependencies.

02

Portability: Containers can run on any system that supports Docker, ensuring applications work consistently across various environments, from development to production.

03

Rapid Deployment: Docker allows quick deployment of applications, reducing the time from development to deployment significantly.



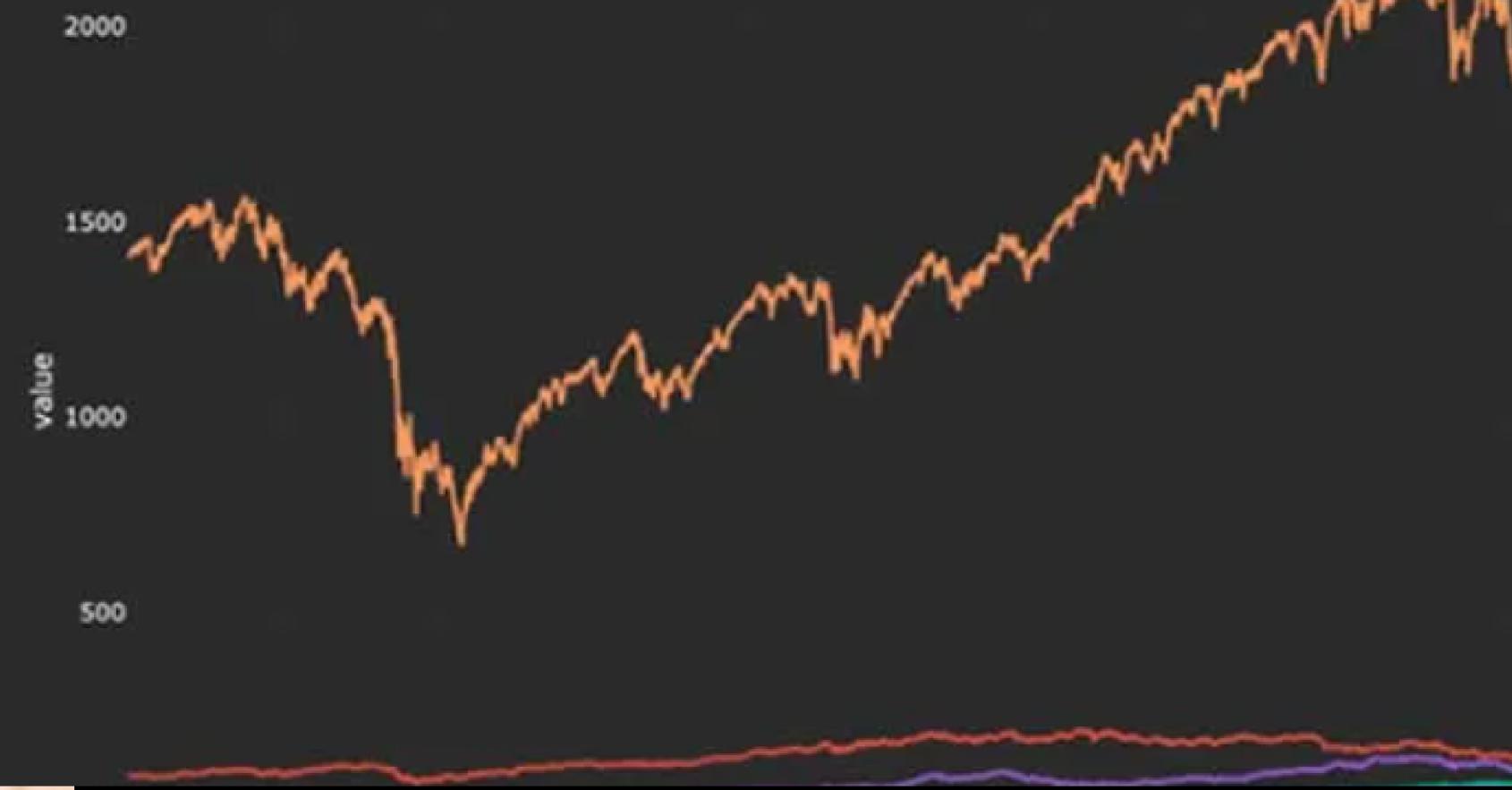
Collaboration: Docker simplifies collaboration between development and operations teams by providing a consistent environment for both. Developers can create containers locally, and operations teams can deploy them without worrying about compatibility issues.



## DASH - STOCK PRICES

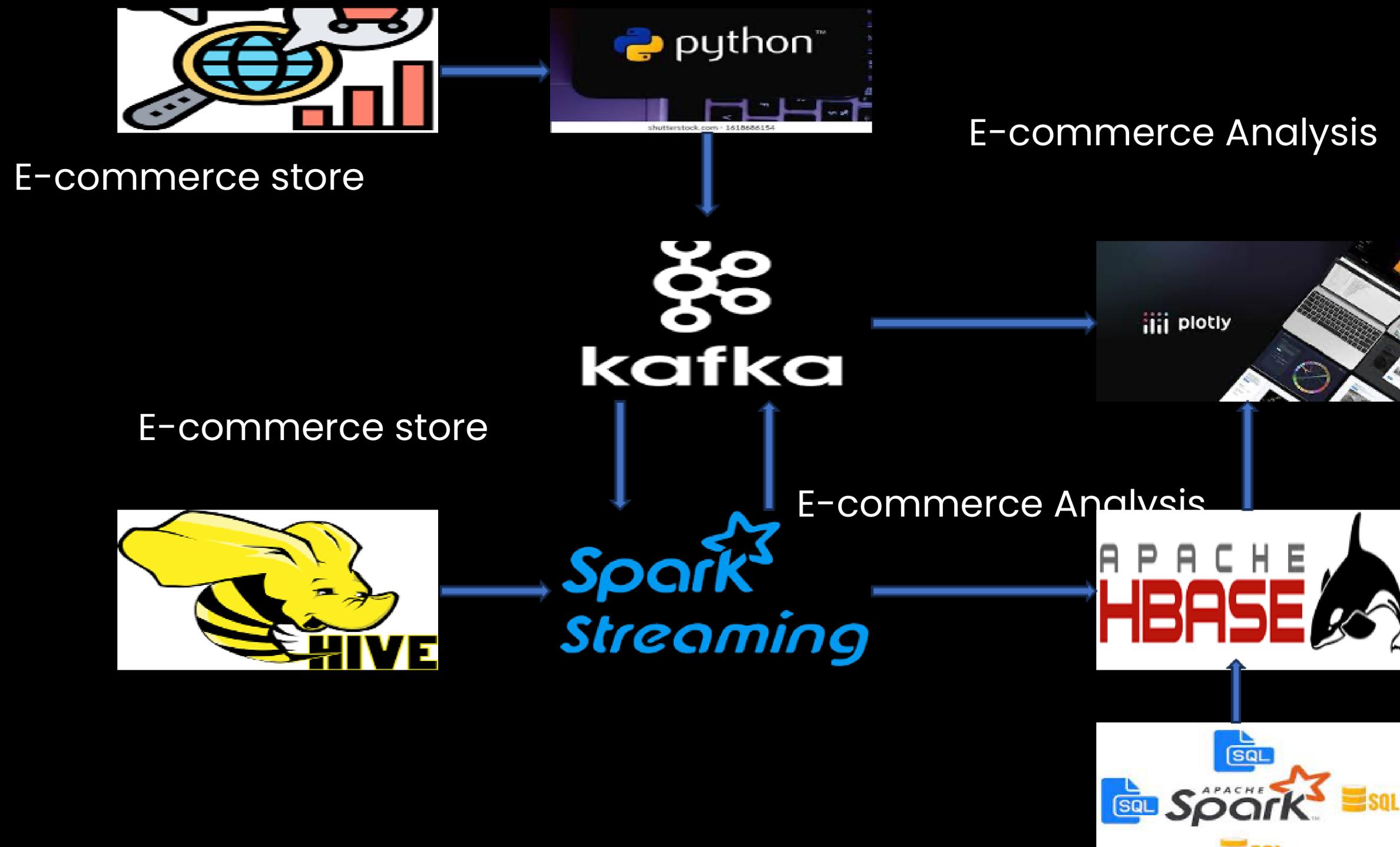
Visualising time series with Plotly - Dash.

Pick one or more stocks from the dropdown below.



- Plotly is a versatile Python graphing library.
- It creates interactive, high-quality graphs.
- Supports various chart types: lines, bars, 3D plots, and heat maps.
- Allows zooming, panning, and hovering over data points for insights.
- Operates seamlessly online and offline.
- In offline mode, visualize plots in local environments like Jupyter Notebooks.
- In online mode, it offers collaborative features for storing, sharing, and teamwork.

# PROJECT FLOW - STREAMING PIPELINE





Docker Desktop Update to latest

Search for local and remote images, containers, and more... **K**

Containers Images Volumes Dev Environments **BETA** Docker Scout **EARLY ACCESS** Learning Center Extensions **•** Telepresence Add Extensions

**stream** **jowilf/stream:0.0.1** 917fa6ef9fdf

STATUS Running (3 minutes ago)

Logs Inspect Terminal Files Stats

```
44415922427789416", "category_code": "computers.components.videocards", "brand": "msi", "price": "624.62", "user_id": "1515915625473223205", "user_session": "d07f8cf4-c0a6-4a0b-be8a-554844ff928d"}]
```

2023-10-04 16:50:22 Message delivered to electronic-store [{"event\_time": "2020-09-28 09:54:04 UTC", "event\_type": "view", "product\_id": "4079305", "category\_id": "2144415922427789416", "category\_code": "computers.components.videocards", "brand": "gigabyte", "price": "499.54", "user\_id": "1515915625473223205", "user\_session": "d07f8cf4-c0a6-4a0b-be8a-554844ff928d"}]

2023-10-04 16:50:22 Message delivered to electronic-store [{"event\_time": "2020-09-28 09:55:52 UTC", "event\_type": "view", "product\_id": "744016", "category\_id": "214415923107266682", "category\_code": "computers.peripherals.printer", "brand": "hp", "price": "569.67", "user\_id": "1515915625520432632", "user\_session": "KjVfbW5itc"}]

2023-10-04 16:50:22 Message delivered to electronic-store [{"event\_time": "2020-09-28 09:56:46 UTC", "event\_type": "view", "product\_id": "3647922", "category\_id": "2144415922427789416", "category\_code": "computers.components.videocards", "brand": "msi", "price": "336.57", "user\_id": "1515915625520475278", "user\_session": "lALG4lJv8u"}]

2023-10-04 16:50:22 Message delivered to electronic-store [{"event\_time": "2020-09-28 09:57:05 UTC", "event\_type": "view", "product\_id": "779865", "category\_id": "214415922427789416", "category\_code": "computers.components.videocards", "brand": "asus", "price": "268.06", "user\_id": "1515915625519044580", "user\_session": "WxN855e18G"}]

2023-10-04 16:50:22 Message delivered to electronic-store [{"event\_time": "2020-09-28 09:59:26 UTC", "event\_type": "view", "product\_id": "4081772", "category\_id": "2144415922427789416", "category\_code": "computers.components.videocards", "brand": "msi", "price": "569.59", "user\_id": "1515915625473223205", "user\_session": "d07f8cf4-c0a6-4a0b-be8a-554844ff928d"}]

2023-10-04 16:50:22 Message delivered to electronic-store [{"event\_time": "2020-09-28 10:00:52 UTC", "event\_type": "view", "product\_id": "3758905", "category\_id": "214415923107266682", "category\_code": "computers.peripherals.printer", "brand": "hp", "price": "467.43", "user\_id": "1515915625520432632", "user\_session": "KjVfbW5itc"}]

# SAMPLE STREAMING DATA TO APACHE KAFKA

The Data is streamed to Apache Kafka from the producers where the Kafka cluster puts it into different brokers and notifies the different consumers that have subscribed to the specific topic.



```
public static void main(String[] args) throws InterruptedException, IOException {
    TableUtils tableUtils = new TableUtils();
    tableUtils.createTable();

    SparkConf sparkConf = new SparkConf().setAppName("JavaNetworkWordCount");

    JavaStreamingContext ssc = new JavaStreamingContext(sparkConf,
        Durations.seconds(2));

    ssc.sparkContext().setLogLevel("ERROR");

    Map<String, Object> kafkaParams = new HashMap<>();
    kafkaParams.put("bootstrap.servers", "kafka:29092");
    kafkaParams.put("key.deserializer", StringDeserializer.class);
    kafkaParams.put("value.deserializer", StringDeserializer.class);
    kafkaParams.put("group.id", "spark");
    kafkaParams.put("auto.offset.reset", "earliest");
    kafkaParams.put("enable.auto.commit", false);
```

The screenshot shows a terminal window with several tabs. The active tab, labeled 'start-all.sh', displays log output from a Docker container. The logs show two distinct time periods with event counts:

- Time: 1696464014000 ms
  - (cart,2)
  - (purchase,2)
  - (view,16)
- Time: 1696464016000 ms
  - (computers,7)
  - (electronics,9)

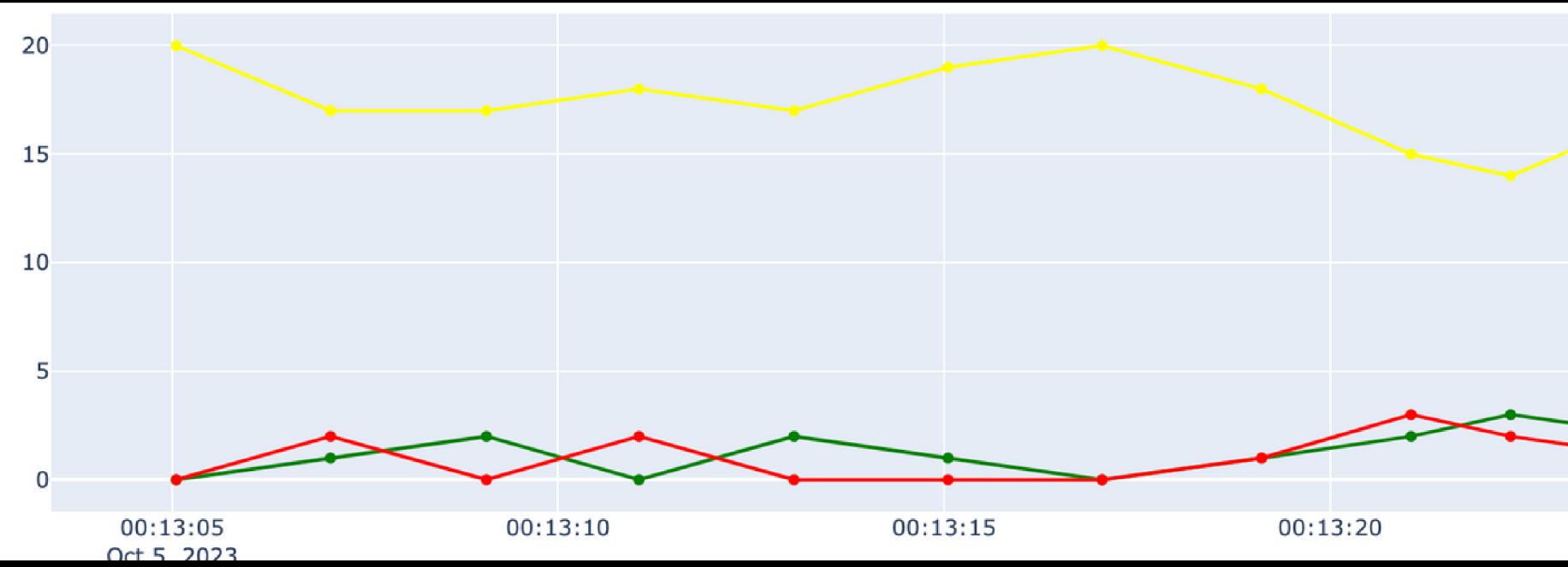
At the top of the terminal window, there is a file tree view showing project files: hbase-env.sh, hbase-site.xml, zoo.cfg, script/start-pseudo.sh, Dockerfile, and src. To the right of the file tree, there are command-line entries:

- 6 docker-compose down
- 7 ./stop-all.sh
- 8
- 9
- 10 docker-compose up -d --build
- 11 ./spark\_sql.sh

# STREAMING WITH SPARK

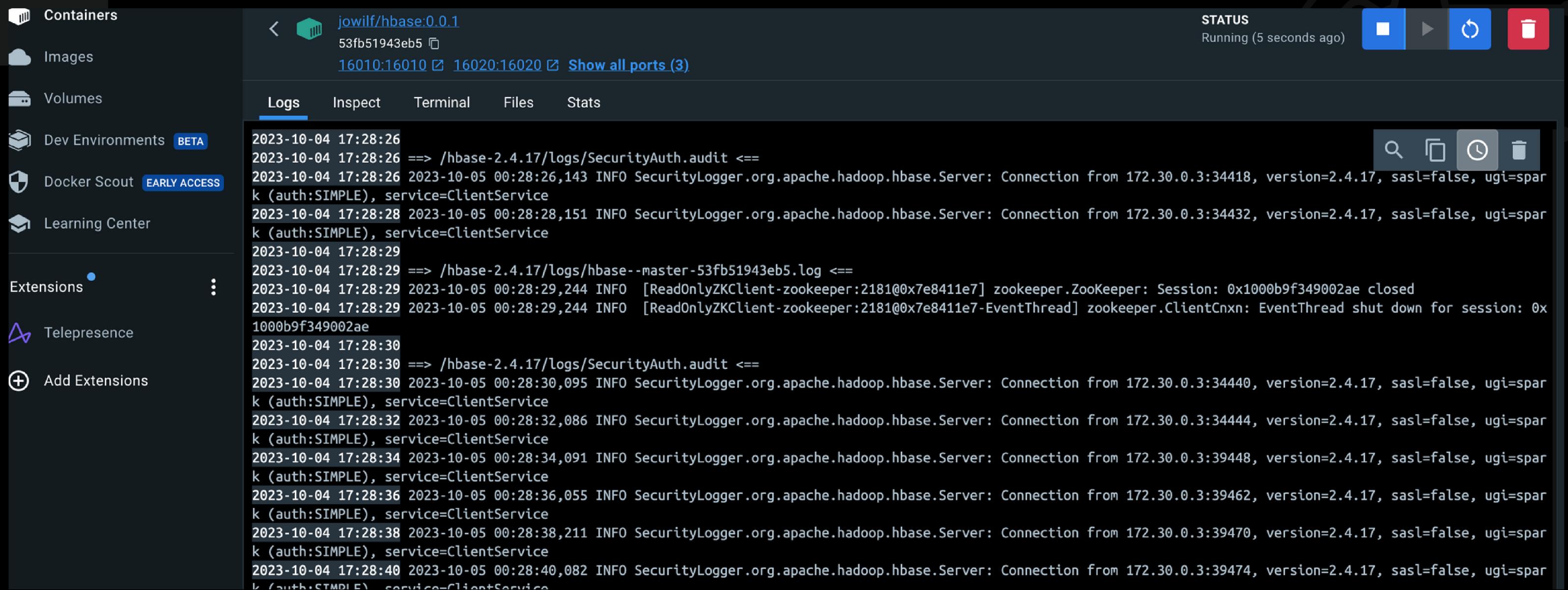


SALFORD & CO.



```
def fig(self):  
    return go.Figure(  
        data=[  
            go.Scatter(  
                x=self.dates,  
                y=self.views,  
                line=dict(color='yellow'), # Set line color to green  
                legendgroup="view",  
                name="View",  
                legendgrouptitle=Legendgrouptitle(text="View"),  
            ),  
            go.Scatter(  
                x=self.dates,  
                y=self.carts,  
                line=dict(color='green'), # Set line color to green  
                name="Cart",  
                legendgroup="cart",  
                legendgrouptitle=Legendgrouptitle(text="Cart"),  
            ),  
            go.Scatter(  
                x=self.dates,  
                y=self.purchases,  
                line=dict(color='red'), # Set line color to red  
                name="Purchase",  
                legendgroup="purchase",  
                legendgrouptitle=Legendgrouptitle(text="Purchase"),  
            )  
        ]  
    )
```

# REAL-TIME DATA VISUALIZATION OR PRESENTATION



# SPARKSTREAM - HBASE

Integrating Apache Spark Streaming with Apache HBase enables real-time data processing and storage capabilities. Spark Streaming allows the processing of live data streams, while HBase serves as a scalable, distributed NoSQL database. By combining these technologies, organizations can analyze and process real-time data streams using Spark Streaming and then store the results in HBase for efficient retrieval and analysis. This integration facilitates seamless handling of large volumes of real-time data, making it ideal for applications requiring instant data processing, analysis, and storage.



```
./stop-all.sh
#docker-compose up -d --build
./spark_sql.sh
#docker-compose restart dashboard
./spark.sh

docker-compose down
./stop-all.sh

docker-compose up -d --build
./spark_sql.sh
# build the jar file
mvn package
# copy our dataset inside the spark docker container
docker cp stream/dataset.csv spark:/dataset.csv
# copy the generated jar file inside the spark docker container
docker cp final-project-1.0-SNAPSHOT.jar spark:/app_sql.jar
# Run spark submit inside the docker container
docker exec -it spark spark-submit --master=local --conf spark.sql.shuffle.partitions=1 --class com.cs523.Spa
docker-compose restart dashboard
./spark.sh
```

# HOW TO RUN THE APPLICATION

To run the application you have to start all shell scripts and it will stop all the running containers. Then it will start again all the containers in the same network by creating new images(if not available) otherwise it will run the existing ones as docker containers.or you can explicitly run by docker-compose command.



SALFORD & CO.

cs523-finalproject – HBaseWriter.java

```
import org.apache.hadoop.hbase.util.Bytes;
import org.apache.hadoop.hbase.util.Tuple2;
import java.io.IOException;
import java.util.ArrayList;
import java.util.List;

2 usages  ↳ Biniam Gheb
public class HBaseWriter {
    5 usages
    private Table table;
    3 usages
    private Connection connection;
    1 usage
    TableUtils tableUtils = new TableUtils();

    2 usages  ↳ Biniam Gheb
    public HBaseWriter() throws IOException {
        connection = tableUtils.newConnection();
    }
}

1 usage  ↳ Biniam Gheb
```

Project pom.xml (final-project) main.py hbase.py SparkSQLHive.java spark\_sql.sh HBaseWriter.java Maven

hbase config script src main java com.cs523 HBaseWriter KafkaStream KafkaWriter SparkSQLAnalyze TableUtils stream app.py dataset.csv Dockerfile requirements.txt .gitignore docker-compose.yml

final-project Lifecycle clean validate compile test package verify install site deploy Plugins Dependencies

# SAMPLE JAVA CODE

## HBASE WRITER



MAHARISHI INTERNATIONAL UNIVERSITY

```
sudo cp /etc/hive/conf/hive-site.xml /usr/lib/spark/conf/
```



```
spark-submit --class "cs523.finalProject.SparkSQLHive"  
--master yarn /home/cloudera/workspace/finalProject/target/finalProject-0.0.1-SNAPSHOT.jar
```

```
SELECT brand, COUNT(event_type) AS viewed FROM electronics WHERE event_type='view'  
GROUP BY brand ORDER BY viewed DESC
```

```
SELECT brand, MAX(price) AS highestPrice FROM electronics GROUP BY brand ORDER BY brand LIMIT 7;
```

```
SELECT brand, SUM(price) AS highestPrice FROM electronics GROUP BY brand ORDER BY brand LIMIT 7;
```

# SPARKSQL COMMANDS

```
3  
4 import ...  
5  
6  
7  
8  
9 public class SparkSQLHive {  
10     public static void main(String[] args) {  
11         JavaSparkContext sc = new JavaSparkContext(new SparkConf().setAppName("electronics").setMaster("local").set("spark.sql.warehousePath", "file:///tmp/electronics"));  
12         HiveContext hiveContext = new HiveContext(sc);  
13         //total sale per brand  
14         DataFrame mostPurchased = hiveContext  
15             .sql("SELECT brand, SUM(cast(price as double)) AS totalSale FROM electronics WHERE event_type='purchase' GROUP BY brand");  
16         mostPurchased.show();  
17     }  
18 }  
19  
20 }  
21 }
```

The screenshot shows a Jupyter Notebook cell containing Java code. The code defines a class named SparkSQLHive with a main method. It sets up a JavaSparkContext and a HiveContext. Then, it runs a SQL query to calculate the total sale per brand from the electronics table where the event type is 'purchase'. The results are grouped by brand and displayed using the show() method.



# WHOLENESS OF PROJECT

In Spark Streaming, real-time data is processed in micro-batches, enabling rapid analysis and response to live data streams. In this project we have been actively involved with different technologies, we have observed the power of individual technologies but more than that we have emphasized how all those technologies can solve realtime processing of huge amounts of data, monitor them using different visualization tools, and how this amount of data can be stored into a scalable Databases for future use.

Similarly, the Science of Consciousness delves into the continuous, ever-changing nature of human awareness and perception. Both concepts emphasize the importance of understanding and adapting to ongoing, evolving information. In a metaphorical sense, Spark Streaming can be seen as analyzing the constant flow of data in the digital world, mirroring the intricate processes underlying the ever-changing landscape of human consciousness explored in the Science of Consciousness. Both highlight the significance of real-time adaptation and analysis in their respective domains, emphasizing the dynamic nature of knowledge and understanding and this can be achieved by diving deep into pure consciousness through Transcendental meditation.

INTEGRATING SCIENCE TECHNOLOGY AND CONSCIOUSNESS



```
CREATE EXTERNAL TABLE electronics (event_time STRING, event_type STRING, product_id int, category_id int, category_code STRING, brand STRING, price FLOAT, user_id int, user_session STRING ) ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde' LOCATION '/home/cloudera/cs523/inputProject' TBLPROPERTIES ('skip.header.line.count'=1);
```

```
//Load command  
LOAD DATA INPATH '/user/cloudera/cs523/inputProject/popular_brand_only.csv' OVERWRITE INTO TABLE employee;
```

## HIVE TABLE



MAHARISHI INTERNATIONAL UNIVERSITY

# THANK YOU

BIG DATA TECHNOLOGY WITH  
PROF. MRUDULA MUKADAM

