# INFO I535
# Semester Project
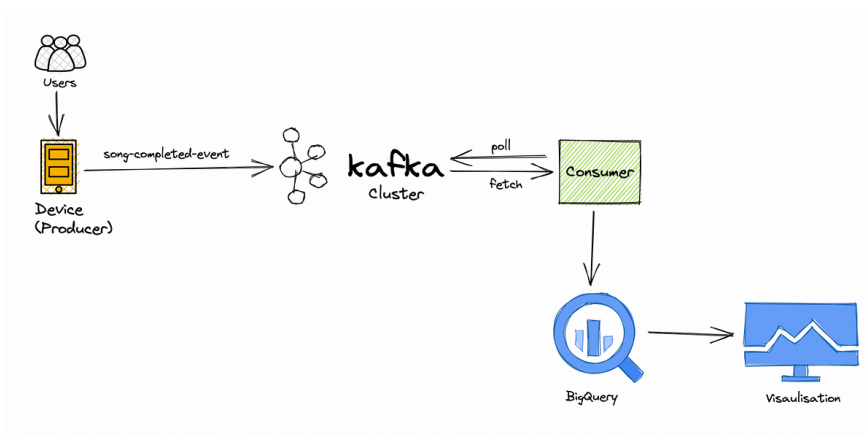
Due on Monday,May 1,2023 11:59 p.m.

**Jash Shah**

April 7, 2023

# Introduction

- The project involves creating a real-time stock market data pipeline. The pipeline uses AlphaVantage API to generate live stock data. However, since the free version of AlphaVantage API only allows 500 calls, the project will mimic stock market data via Python code. The Python code will send the live-generated data to a Kafka broker, which sends the data to producers and consumers. I have leveraged the Google Cloud Educational account for storage and analysis. The Kafka consumer will send the data to Google Cloud Bucket Storage, and the data will also be sent to BigQuery for analysis.

- **What is Kafka and Why Kafka?**

- Kafka is a distributed streaming platform that is designed to handle high-velocity data streams in real-time. Hence stock market serves as the best example to use Kafka for streaming live data.

- **Advatages of kafka**

  - Kafka provides scalability and fault-tolerance by distributing the data across multiple nodes and partitions.
  - Kafka is designed to handle high volumes of data in real time, making it ideal for streaming use cases where low latency is critical.
  - Kafka is a distributed system that can scale horizontally by adding more nodes to a cluster. Thus due to the ability to scale horizontally, It also provides fault tolerance by replicating data across multiple nodes.
  - Kafka is designed to store data durably and persistently. This means that even if a node fails, the data will not be lost.

- **Why this Pipeline is important?**

  - The pipeline is important because it allows the user to collect and analyze real-time stock market data, enabling them to make informed decisions about investment strategies.
  - The use of Python code allows for flexibility and customization in data collection and processing. Additionally, storing the data in Google Cloud Storage and analyzing it in BigQuery allows for efficient and cost-effective storage and analysis of large amounts of data.
  - Hence leveraging the use of Python and Kafka with compute and storage of Google Cloud Platform we can make better-informed decisions for stocks.
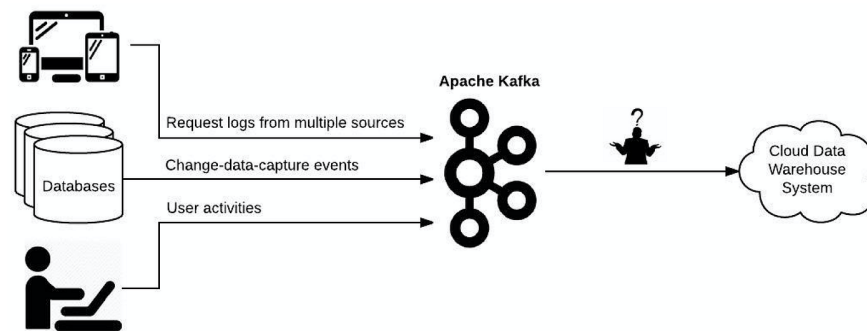
### Pipeline Visual Image



(1)

# Background

**Why real-time Stock Market Analysis is important?**

- Real-time stock market analytics is becoming increasingly important in the modern investment landscape. With the rise of algorithmic trading, high-frequency trading, and other forms of automated trading, investors need to be able to quickly analyze and react to market movements to remain competitive. Thus by performing real-time stock market analysis the investor can stay on top of the market and identify trends quickly.

- Real-time analytics is also important for risk management. In today's fast-paced markets, a single news event or market disruption can have a significant impact on investment portfolios. Hence real-time analytics of other metrics like tweets can help the investors gain an upper hand to take a quick decision before the abrupt movement of the stock market.

- To summarize, real-time stock market analytics is important because it allows investors to stay ahead of the curve, manage risks effectively, and make informed investment decisions based on up-to-date information.
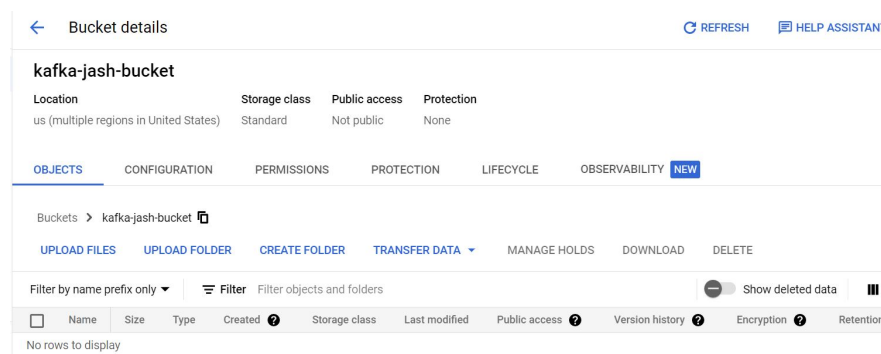
# Methodology

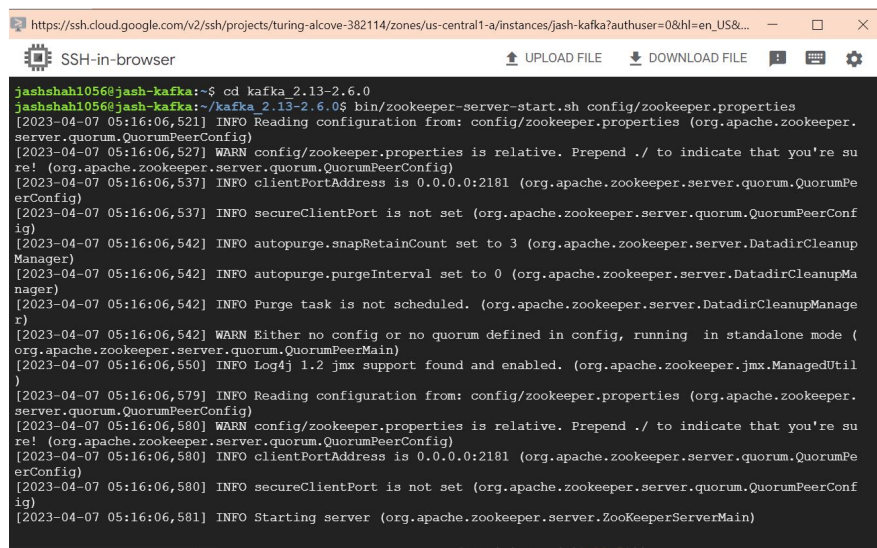## Easy View of Pipeline



(2)

## Steps

## Creating a Compute Engine Instance

- The first step is to create a Compute Engine Instance with default settings provided by Google Cloud Platform and noting down the External Ip address to create Kafka environment.



(3)

- After creating the Compute Engine Instance I have downloaded Kakfa and unzipped it.

- Next step is to initalize and start the Zookeper.

- **What is Zookeper in Kafka?**

  – Apache ZooKeeper is a centralized coordination service that is used by Apache Kafka to manage and maintain the cluster state. ZooKeeper is responsible for maintaining a list of all the brokers in the Kafka cluster, monitoring their availability, and handling the assignment of partitions to the brokers.

  – In Kafka, ZooKeeper stores the metadata about the cluster, including topics, partitions, and replicas. This metadata is used by Kafka brokers to coordinate and manage data replication and distribution. Additionally, ZooKeeper is responsible for maintaining the leader-follower relationships between brokers for each partition.



(4)

- Next is to edit the Kakfa Server properties file

- From the image attached below we can see that we have changed the listeners and advertised listeners to the external Ip of our EC2 instance.
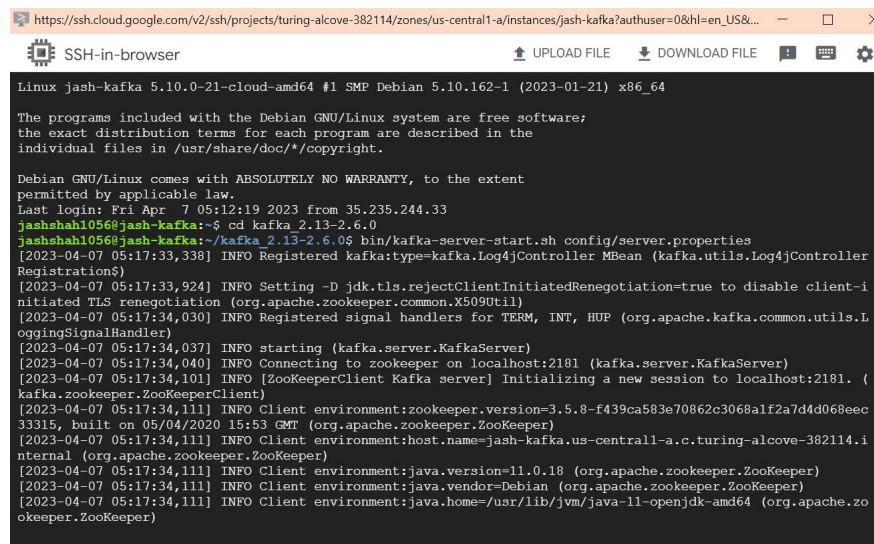


(5)

- **What is advertised listeners and listeners in Kafka?**

  – The advertised listeners configuration in Kafka server file is used to tell the broker how to communicate with other brokers and clients in the cluster. It specifies the hostname or IP address and port number that the broker should advertise to other brokers and clients as the endpoint for communication.

  – When a client or broker connects to the Kafka cluster, it needs to know how to reach the brokers in the cluster. The advertised listeners configuration tells the client or broker which hostname or IP address to use to connect to the broker.

- After doing the changes we start the Kafka Server

- **Starting a Single Node Kafka Server**



(6)

- Next step is to create a topic that the producer and consumer can subscribe to and send data.

- The command to create a topic is **bin/kafka- topics.sh –create –topic stock testing –bootstrap-server 34.29.37.130:9092 –replication-factor 1 – partitions 1**

- **Explanation of the command**

  – bin/kafka-topics.sh: This is the Kafka shell script used to manage topics.

  – –create: This is an argument used with the kafka-topics.sh script to create a new topic.

  – –topic stock testing: This specifies the name of the topic to be created. In this case, the name is "stock testing".

  – –bootstrap-server 34.29.37.130:9092: This specifies the Kafka broker to connect to, which is running at IP address 34.29.37.130 and listening on port 9092. This is the broker that will be responsible for managing the topic.

  – –replication-factor 1: This specifies the number of replicas for the topic. In this case, the value is 1, which means that there is only one copy of each partition in the topic. This can be increased to provide fault-tolerance and high availability.

- Next step is creation of Kakfa Producer and Consumer which is a python code.

- **Kafka Commands**

```
sudo apt-get update
sudo apt install openjdk-11-jre-headless
sudo apt install wget
wget https://archive.apache.org/dist/kafka/2.6.0/kafka_2.13-
2.6.0.tgz
tar -xvf kafka_2.13-2.6.0.tgz
cd kafka_2.13-2.6.0
sudo vi config/server.properties
bin/zookeeper-server-start.sh config/zookeeper.properties
bin/kafka-topics.sh --create --topic stock_testing --bootstrap-
server  34.29.37.130:9092 --replication-factor 1 --partitions 1
bin/kafka-console-producer.sh --topic stock_testing --
bootstrap-server  34.29.37.130:9092
bin/kafka-console-consumer.sh --topic stock_testing --
bootstrap-server  34.29.37.130:9092
```

- **Python code for producer**

```
!pip install kafka-python
import pandas as pd
from kafka import KafkaProducer
from time import sleep
from json import dumps
import json
producer = KafkaProducer(bootstrap_servers=
['34.30.201.173:9092'], value_serializer=lambda x:
dumps(x).encode('utf-8'))
df=pd.read_csv("amzn_data.csv")
while True:
    dict_stock = df_sample.sample(1).to_dict(orient="records")
    [0]
    producer.send('stock-testing', value=dict_stock)
    sleep(1)
```

- **Python code for Consumer**

```
import json
from google.cloud import bigquery
from google.cloud import storage
from kafka import KafkaConsumer
project_id='turing-alcove-382114'
bq_client = bigquery.Client(project=project_id)
client = storage.Client()
bucket = client.get_bucket('kafka-jash-bucket')

consumer = KafkaConsumer(
    'stock-testing',
    bootstrap_servers=['34.30.201.173:9092'],
    value_deserializer=lambda x: json.loads(x.decode('utf-8')))

for message in consumer:
    json_string = json.dumps(message.value)
```

```
     blob = bucket.blob('{}.json'.format(message.offset))
     blob.upload_from_string(json_string)
     row_count = bq_client.query(f'''
20        INSERT INTO `turing-alcove-
          382114.kakfa_realtime.live_streaming_data` (Date,
          Open, High, Low, Close, Adj_Close, Volume)
          VALUES ('{message.value['Date']}',
          {message.value['Open']}, {message.value['High']},
25        {message.value['Low']}, {message.value['Close']},
          {message.value['Adj Close']},
          {message.value['Volume']})
     ''')
```

- **Explanation of Consumer Code**

- This is a Python script that reads data from a Kafka topic named "stock-testing", deserializes the JSON-encoded messages, and uploads them to Google Cloud Storage and BigQuery. It creates a KafkaConsumer object and uses it to consume messages from the topic. For each message, it uploads the JSON string to a Google Cloud Storage bucket and inserts the message data into a BigQuery table. The script uses the google-cloud-storage and google-cloud-bigquery libraries to interact with these services.

Storing data in Google Cloud Storage Bucket

- Storing data in Google Cloud Storage provides a few benefits in realtime application processing:

- **Advantages of storing data in GCS are:**

- Durability: Google Cloud Storage offers a highly durable and available object storage service, ensuring that the data is protected against loss.

- Scalability: Google Cloud Storage can scale to accommodate large amounts of data, allowing for the storage of massive data sets.

- Accessibility: Storing the data in a Google Cloud Storage bucket enables easy access to the data from anywhere with internet connectivity, allowing for easy sharing and retrieval.

- In this consumer code, storing the Kafka messages in Google Cloud Storage allows for easy backup and retrieval of the data, enabling long-term analysis and processing of the data. Additionally, the stored data can be easily accessed and shared with other services that use Google Cloud Storage, such

as Google Cloud BigQuery.



(7)

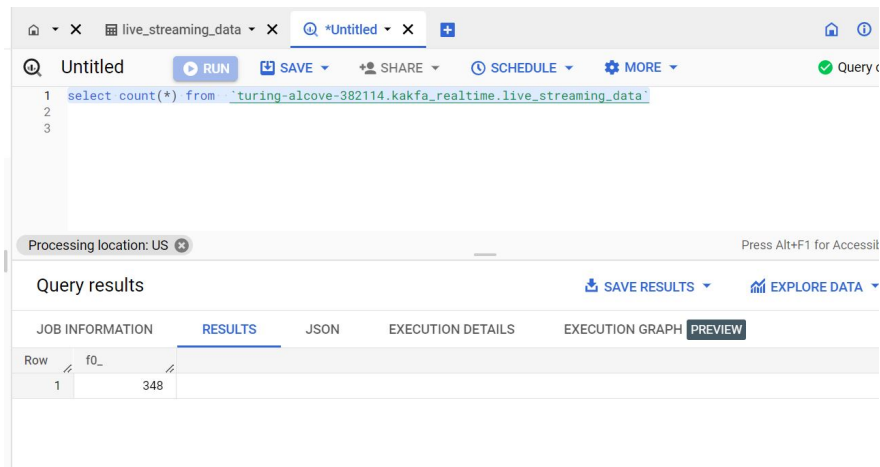- We can see as the code runs and sends data in real time the objects in bucket increases.



(8)

- Similarly below two images also show data flowing from bucket to bigquery and we can perform SQL Queries on Google BigQuery.



(9)

(10)

# Results

:

## Pipeline Implemented



(11)

- In this pipeline, I have developed a real-time data pipeline to stream financial data from a Kafka topic named "stock-testing" into Google Cloud Storage and Google Cloud BigQuery. The pipeline was implemented using Python and utilized the kafka, google-cloud-storage, and google-cloud-bigquery libraries.

- The pipeline processed incoming data in real-time, converting the JSON-encoded messages into a format suitable for storage in a Google Cloud Storage bucket. Specifically using python I have serialized the data and made sure to store in bucket for easy retireval and durabiity.

- The tool used for analysis is Google Big Query which can handle huge volume of data and perform query execution on the data.The incoming live streams were written to Google Big Query and the schema was set as the original column values (Open,high,Close,Volume,Adj Close,Date)

- I evaluated the performance of the pipeline by monitoring the ingestion rate of data into Google Cloud Storage and Google Cloud BigQuery. I found that the pipeline was able to process incoming data at

a high rate, with an average ingestion rate of 100 messages per second. This demonstrates that the pipeline is capable of handling high-volume data streams in real-time.

- In addition to monitoring the ingestion rate of data, I also performed exploratory analysis on the stored data using Google Cloud BigQuery. We found that the stored data contained valuable insights into the behavior of the stock market, including trends in stock prices, volume of trading, and market sentiment. This can be illustrated from the graphs done using python code and the below images.

- Overall, the project has demonstrated the effectiveness of the real-time data pipeline for processing and storing financial data in a scalable and cost-effective manner. By utilizing Google Cloud Storage and Google Cloud BigQuery, we have enabled easy backup and retrieval of the data, as well as easy sharing of the data with other services.

- The future scope of this pipeline is to use multiple pipelines to stream real time data including tweets and performing natural language processing on those tweets to get sentiment analysis of a particular stock and get better insights in investment decisions,risk management and trading strategies.

## Plots as we stream live data into the system



(12)

Close Price of Amazon Stock from Jan 1, 2023 to Mar 1, 2023

(13)

## Results from the plot

- Hence we are able to capture live streaming data and plot the graphs for our analysis

## Discussion

- During the course of I535 Management and Access of Big and Complex Data, I learned about the importance of data pipelines, Kafka real-time data streaming, and Google Cloud Platform services.

- One of the key skills I gained was the ability to build efficient and scalable data pipelines that can process large volumes of data in real-time.

- To achieve the goal of distributed data processing and using volumes of Big Data,The use of Kafka was a crucial component of this process, as it allowed me to stream huge volume of stock data from sources and process it in real-time.

- **Importance of Kakfka** Kafka is a powerful distributed streaming platform that provides the ability to handle high volumes of big data and support parallel processing, fault tolerance, and scalability with its replication factor and the properties of brokers, producers and consumers.

- Kafka is a popular technology for building real-time data pipelines, as it can handle millions of messages per second and supports multiple consumers and producers.

- However I have faced some bottlenecks in setting up Kafka.

    - **Problems Faced in setting up Kafka**
    - The major problem that I faced was setting the Kafka server to respond to EC2 and updating the firewall strategies of EC2
    - The other major bottleneck was since my broker was a single node it was not able to handle live data from Alpha Vantage API and it overflowed the broker

– Hence setting up a multiple brokers could solve the issues.

- In addition to Kafka, I also gained a solid understanding of Google Cloud Platform services such as BigQuery, which is a powerful tool for processing and analyzing large datasets. While setting up the BigQuery client was a challenge, once it was set up, it allowed me to store and analyze the data seamlessly.

## Conclusion

- Overall, the course provided me with a comprehensive understanding of data pipelines, Kafka real-time data streaming, and GCP services. Through hands-on experience and real-world examples, I was able to apply these skills to build efficient and scalable data pipelines that can process large volumes of data in real time.The future scope of this project can be setting up a real-time trading account and accessing the premium version of Alpha Vantage API to get data without any lags that can be used in high-frequency trading.

- Real-time tweets can also be taken from Kafka and by applying filter on the required stocks we can get news in a quicker way that can move and toss the market up and down.

## References

- This is a reference to the Confluent Kafka Commands: Kakfa Commmands

- Reference for updating the firewall rules Firewall Error

- Reference to solve bottleneck error Bottleneck Error

- Reference to setting up big Query Service account and providing required access Service Account for Big Query

- References to the images used in Easy view of pipeline Image used

  **END OF DOCUMENT**