# A PROJECT REPORT
## On
# Stock Market Real-Time Data Analysis Using Kafka

### Submitted to
# KIIT Deemed to be University

## In Partial Fulfilment of the Requirement for the Award of

## BACHELOR'S DEGREE IN
## Computer Science and Engineering

## BY

| | |
|---|---|
| **Mahadev Mondal** | **2005106** |
| **Nilanjan Saha** | **2005111** |

### UNDER THE GUIDANCE OF
### Dr. Mainak Bandyopadhyay

## SCHOOL OF COMPUTER ENGINEERING
# KALINGA INSTITUTE OF INDUSTRIAL TECHNOLOGY
### BHUBANESWAR, ODISHA - 751024
### Dec 2023

# KIIT Deemed to be University

School of Computer Engineering
Bhubaneswar, ODISHA 751024



# CERTIFICATE

This is certify that the project entitled

## "Stock Market Real-Time Data Analysis Using Kafka"

submitted by

| | |
|---|---|
| **Mahadev Mondal** | **2005106** |
| **Nilanjan Saha** | **2005111** |

is a record of bonafide work carried out by them, in the partial fulfilment of the requirement for the award of Degree of Bachelor of Engineering (Computer Science & Engineering) at KIIT Deemed to be university, Bhubaneswar. This work is done during year 2022-2023, under our guidance.

Date: 01/12/2023

Dr. Mainak Bandyopadhyay
Project Guide

# Acknowledgements

We are profoundly grateful to **Dr. Mainak Bandyopadhyay** of **School of Computer Engineering** for his expert guidance and continuous encouragement throughout to see that this project rights its target since its commencement to its completion.

|  |  |
|---|---|
| **Mahadev Mondal** | **2005106** |
| **Nalanjan Saha** | **2005111** |

# ABSTRACT

In the dynamic realm of financial markets, characterized by constant flux and rapid decision-making requirements, the Stock Market Real-Time Data Analysis project emerges as a pivotal solution. This endeavor strategically employs advanced technologies, including Apache Kafka, Python, and Amazon Web Services (AWS), to craft an end-to-end data engineering framework. The primary thrust of the project revolves around establishing a robust Kafka cluster, seamlessly integrating with AWS services for comprehensive data storage and analysis, and constructing an efficient data pipeline for the continuous streaming of real-time financial data.

The essence of the project lies in addressing the critical need for timely and effective data processing, essential for navigating the complexities of financial markets. Apache Kafka, renowned for its high-throughput, fault tolerance, and scalability, takes center stage as the backbone of the project. This distributed streaming platform facilitates the seamless ingestion and processing of vast datasets generated in real-time, aligning perfectly with the dynamic nature of financial markets.

The integration with AWS introduces a cloud-based infrastructure, enhancing scalability and flexibility. AWS services such as S3 for storage, Athena for SQL-based querying, Glue for automated ETL processes, and EC2 for computational resources collectively contribute to the project's adaptability and resilience. The synergy of these technologies results in an agile and responsive data engineering solution, fostering quick and informed decision-making in the financial domain.

This report provides a comprehensive exploration of the project, navigating through the intricacies of its architecture, the implementation strategies employed, the methodologies applied for rigorous testing, and a forward-looking discussion on its future potential. By creating an efficient real-time data analysis pipeline, this project not only addresses the immediate challenges of financial market data processing but also lays the groundwork for advancing the landscape of data-driven decision-making in the financial sector.

Keywords: Data Engineering, Data Analysis, Apache Kafka, EC2, AWS.

# Contents

# List of Figures

# Chapter 1

# <u>Introduction</u>

The financial landscape operates within a relentless and dynamic environment, where instantaneous decisions can be the differentiator between success and missed opportunities. In this context, the Stock Market Real-Time Data Analysis project assumes significance by recognizing the critical importance of timely and accurate data analysis in the realm of financial markets. The project is conceived with the explicit goal of developing a robust data engineering solution, strategically harnessing the power of Apache Kafka to enable real-time data streaming. The seamless integration of this solution with Amazon Web Services (AWS) further fortifies the system's capabilities, enhancing scalability and flexibility to meet the evolving demands of the financial ecosystem.

In the contemporary financial arena, where information is not only power but also a fleeting commodity, the need for timely data analysis is paramount. Financial markets are characterized by constant fluctuations influenced by a multitude of factors including economic indicators, geopolitical events, and investor sentiment. The ability to swiftly access, process, and derive actionable insights from this deluge of real-time data becomes the linchpin for effective decision-making. Traditional batch processing systems, with their inherent time delays, are inadequate for meeting the exigencies of this fast-paced environment.

The project's core strategy involves leveraging Apache Kafka, a distributed streaming platform celebrated for its high-throughput, fault tolerance, and scalability. By adopting Kafka as the cornerstone technology, the project aims to establish a real-time data streaming infrastructure capable of ingesting and processing vast volumes of financial data instantaneously. Kafka's publish-subscribe model and partitioning capabilities align seamlessly with the dynamic nature of financial markets, providing an efficient solution for the continuous flow of real-time data.

To augment and extend the capabilities of Apache Kafka, the project integrates with AWS services. AWS, a cloud computing giant, offers a suite of services that complements and enhances the overall architecture. S3 provides scalable storage, Athena facilitates SQL-based querying of data in S3, Glue automates the discovery and classification of datasets, and EC2 offers virtual servers for computation. This integration not only ensures the system's adaptability to varying data loads but also positions it for future scalability

and optimization.

In essence, the Stock Market Real-Time Data Analysis project is not just a technological solution; it is a strategic response to the evolving needs of financial market participants. By developing a robust data engineering solution that combines the strengths of Apache Kafka and AWS, the project aims to empower stakeholders with the agility to make informed decisions in real-time, thereby contributing to more responsive and effective financial strategies. The subsequent sections of this report delve into the detailed architecture, implementation strategies, testing methodologies, and the future potential of this innovative solution.

# Chapter 2

# **Basic Concepts/ Literature Review**

## **2.1 Apache Kafka**

Apache Kafka serves as the cornerstone of this project, providing a distributed and fault-tolerant streaming platform. Its publish-subscribe architecture facilitates real-time data streaming, ensuring that financial data is processed and disseminated without delays.

Apache Kafka is a distributed streaming platform designed for high-throughput, fault tolerance, and real-time data processing. In the Stock Market Real-Time Data Analysis project, Kafka serves as the backbone, facilitating the seamless streaming of financial data.

Apache Kafka stands as a robust and highly acclaimed distributed streaming platform, recognized for its exceptional capabilities in managing real-time data streams. In the context of the Stock Market Real-Time Data Analysis project, Apache Kafka assumes a pivotal role as the backbone of the entire data processing infrastructure.

Key Features:

a) Publish-Subscribe Model: Kafka adopts a publish-subscribe model, allowing data producers to push messages to topics, and consumers to subscribe to topics of interest.

b) Partitioning: Data in Kafka is organized into topics, and each topic is divided into partitions. Partitioning enables parallel processing, enhancing efficiency.

c) Fault Tolerance: Kafka ensures fault tolerance by replicating data across multiple broker nodes, preventing data loss in the event of hardware failures.

d) Scalability: Kafka can scale horizontally by adding more broker nodes to the cluster, accommodating varying data loads.

Integration in the Project:

In the Stock Market Real-Time Data Analysis project, Apache Kafka orchestrates the flow of financial data, ensuring rapid ingestion and processing. Its distributed architecture and streaming capabilities align with the project's goal of real-time data analysis in the dynamic financial market environment.
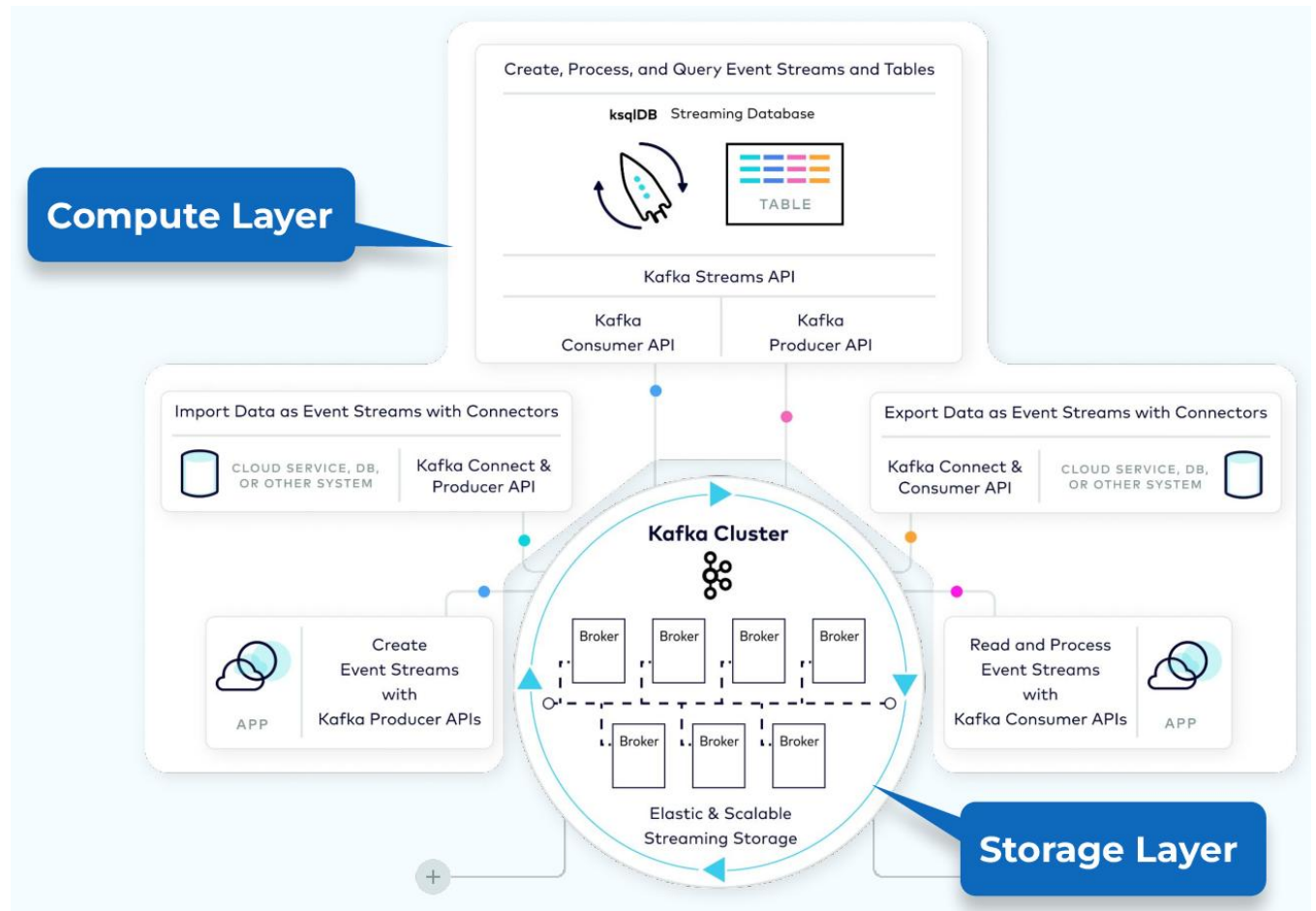


Fig 1:-

## 2.2 AWS Services

The integration of AWS services enriches the project's capabilities. S3, a scalable object storage service, is employed for data storage, while Athena enables SQL-based querying of data stored in S3. Glue, with its ETL capabilities, automates the discovery and classification of datasets. EC2 instances are utilized to run applications and enhance computational capabilities.

### 2.2.1 S3 (Simple Storage Service)

Amazon S3 (Simple Storage Service) is a scalable and highly durable object storage service offered by AWS. It provides developers with secure and efficient storage for a wide range of data types, including images, videos, and datasets. S3 is a key component in the Stock Market Real-Time Data Analysis project, serving as the primary storage solution for datasets and other relevant files.

Key Features:
  a) Scalability: S3 can scale virtually infinitely, accommodating growing amounts of data without any manual intervention.

  b) Durability and Reliability: S3 offers 99.999999999% (11 9's) durability, ensuring data remains intact even in the face of hardware failures or other issues.

  c) Security: S3 supports access control mechanisms, allowing fine-grained control over who can access the stored data. Encryption options are available for enhanced security.

### 2.2.2 Athena

Amazon Athena is a serverless, interactive query service that enables SQL-based querying of data stored in Amazon S3. It is particularly valuable for ad-hoc queries and quick analysis of large datasets. In the Stock Market Real-Time Data Analysis project, Athena plays a crucial role in extracting insights from the data stored in S3 using familiar SQL syntax.

Key Features:
  a) Serverless Architecture: Athena is fully managed, eliminating the need for infrastructure provisioning or management.

  b) Cost-effective: Users only pay for the queries they run, making it a cost-effective solution for occasional or sporadic querying.

c) Integration with S3: Athena seamlessly integrates with S3, enabling the analysis of vast datasets without the need for data movement.

### 2.2.3 Glue Crawler

AWS Glue Crawler is an automated ETL (Extract, Transform, Load) tool that discovers and classifies data stored in various sources, making it accessible for analysis. In the Stock Market Real-Time Data Analysis project, Glue Crawler automates the process of identifying and cataloging datasets, ensuring that the data pipeline remains dynamic and responsive to changes.

Key Features:
a) Automated Discovery: Glue Crawler automatically identifies and classifies data, reducing manual effort and ensuring the catalog is always up-to-date.

b) Schema Inference: It infers the structure of the data, allowing for schema evolution without manual intervention.

c) Integration with Glue Catalog: The metadata collected by the crawler is stored in the Glue Catalog, creating a centralized repository for dataset information.

### 2.2.4 Glue Catalog

AWS Glue Catalog is a metadata repository that stores information about datasets, tables, and transformations in a data lake or data warehouse. It plays a central role in the Stock Market Real-Time Data Analysis project by providing a unified and organized view of the data, facilitating efficient querying and analysis.

Key Features:
a) Unified Metadata: Glue Catalog consolidates metadata from various sources, providing a single, consistent view of the data.

b) Integration with Other AWS Services: The catalog seamlessly integrates with other AWS services like Athena, allowing for easy access to metadata during analysis.

c) Data Lineage: Glue Catalog captures data lineage information, tracking the flow of data from source to destination.

### 2.2.5 EC2 (Elastic Compute Cloud)

Amazon EC2 (Elastic Compute Cloud) provides scalable virtual servers in the cloud. In the Stock Market Real-Time Data Analysis project, EC2 instances are utilized to run applications, execute code, and enhance computational capabilities.

Key Features:
a) Scalability: EC2 instances can be easily scaled up or down based on the computational requirements of the project.

b) Variety of Instance Types: EC2 offers a variety of instance types optimized for different use cases, such as compute-optimized, memory-optimized, and GPU instances.

c) Flexibility: Users have full control over the virtual machines, allowing customization of the operating system, security settings, and network configurations.

# Chapter 3

# **<u>Problem Statement / Requirement Specifications</u>**

## 3.Problem Statement / Requirement Specifications

## 3.1 Project Planning

Project planning is a critical phase that sets the foundation for the entire Stock Market Real-Time Data Analysis project. Project planning involves meticulous organization and allocation of resources. Defined milestones, resource allocation, and a well-structured timeline are essential for successful execution. Effective project planning is crucial for the successful execution of the Stock Market Real-Time Data Analysis project. It involves meticulous organization and resource allocation to ensure a smooth workflow. Defined milestones help track progress and keep the project on schedule. Resource allocation, including human resources and infrastructure, ensures that the necessary tools and expertise are available throughout the project lifecycle. A well-structured timeline provides a clear roadmap, allowing efficient coordination and prioritization of tasks. Regular monitoring and adjustment of the project plan help identify and mitigate any potential risks or delays.

## 3.2 Project Analysis

A comprehensive analysis phase delves into understanding the intricacies of the stock market dataset, defining data processing requirements, and strategically selecting technologies that align with project objectives. The project analysis phase is essential for understanding the complexities of the stock market dataset and defining the data processing requirements. Thorough analysis of the dataset helps identify crucial attributes and patterns, allowing for informed decision-making. It involves exploring the dataset's structure, data types, and possible data quality issues. Understanding the characteristics of the data enables the selection of appropriate data processing techniques and technologies. Strategic selection of technologies, such as Apache Kafka, Python, and AWS, aligns with the project's objectives and ensures efficient data processing and analysis.

## 3.3 System Design

### 3.3.1 Design Constraints

System design constraints encompass considerations for resource limitations, network latency, and security measures to ensure the system's stability and reliability. Design constraints in the Stock Market Real-Time Data Analysis project encompass various considerations. Resource limitations refer to constraints related to computational power, memory, and storage capacity. These limitations must be accounted for to optimize system performance and scalability. Network latency is another critical constraint to consider, as it affects the speed and efficiency of data transmission between components. Security measures are crucial to protect sensitive financial data and ensure compliance with regulatory requirements. Implementing encryption, access controls, and monitoring mechanisms helps maintain the system's stability and reliability.

### 3.3.2 System Architecture

The system architecture outlines the interaction between Python, Kafka, and AWS components. It defines the flow of data from its source through the entire pipeline. The system architecture for this project outlines the interaction between Python, Kafka, and AWS components. It defines the flow of data from its source through the entire pipeline. The architecture includes components such as data ingestion, data processing, data storage, and data analysis. Apache Kafka acts as the central data hub, receiving real-time market data streams. Python is used for data processing tasks such as cleaning, transformation, and aggregation. AWS components, including Amazon S3 for data storage and Amazon Redshift for data warehousing and analysis, are seamlessly integrated into the architecture. The system architecture ensures a scalable, fault-tolerant, and efficient data processing pipeline from real-time data ingestion to actionable insights.
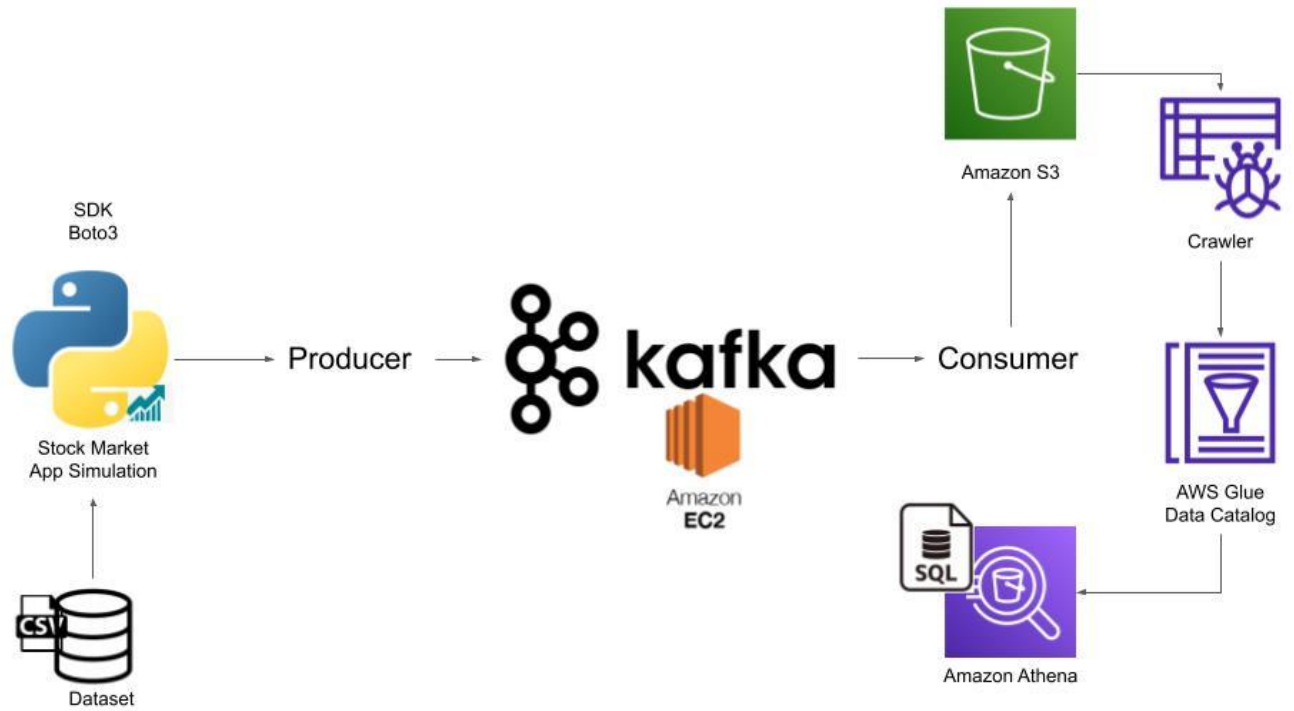
**System Architecture Block Diagram:-**

Fig 2:-

# Chapter 4

# __Implementation__

## 4. Implementation

## 4.1 Methodology

The methodology serves as a comprehensive guide for the step-by-step implementation of the Stock Market Real-Time Data Analysis project. It encompasses key phases such as Kafka setup, AWS integration, and the creation of a Python-based data pipeline for real-time data streaming. The Kafka setup involves the installation and configuration of the distributed streaming platform, ensuring that it meets the project's requirements for fault tolerance, scalability, and high throughput. AWS integration involves connecting the project with services like S3, Athena, and Glue, establishing a cohesive ecosystem for data storage and analysis. The Python-based data pipeline creation focuses on leveraging the language's versatility for efficient data manipulation and streaming.

## 4.2 Testing

The testing phase is a critical component that rigorously evaluates the functionality of the entire system. This includes scenario-based testing to simulate real-world usage scenarios, load testing to assess the system's performance under varying workloads, and performance testing to measure responsiveness and efficiency. Testing ensures that each component, from Kafka to AWS integration to Python-based data processing, performs according to specifications. It identifies potential issues, validates the accuracy of real-time data processing, and ensures that the system can handle the dynamic nature of financial data without compromising performance.

## 4.3  Result Analysis

Result analysis delves into the examination of real-time data processing outcomes. It involves assessing the accuracy of processed financial data, evaluating the system's responsiveness to changes in the market, and analyzing overall performance metrics. This phase provides insights into how well the system aligns with the project's objectives and whether it meets the demands of real-time decision-making in financial markets. Result analysis guides refinements and optimizations, ensuring that the system consistently delivers reliable and timely data insights.

```
{"Index": "SSMI", "Date": "2015-07-21", "Open": 9477.009766, "High": 9479.30957, "Low": 9363
.870117, "Close": 9385.450195, "Adj Close": 9385.450195, "Volume": 35505300.0, "CloseUSD": 1
0417.84971645}
{"Index": "GSPTSE", "Date": "1999-07-08", "Open": 7175.799805, "High": 7175.799805, "Low": 7
134.799805, "Close": 7138.799805, "Adj Close": 7138.799805, "Volume": 7774460000.0, "CloseUS
D": 5925.20383815}
{"Index": "N225", "Date": "1991-05-10", "Open": 26440.48047, "High": 26449.06055, "Low": 262
12.64063, "Close": 26274.28906, "Adj Close": 26274.28906, "Volume": 0.0, "CloseUSD": 262.742
8906}
{"Index": "GSPTSE", "Date": "1994-10-17", "Open": 4327.200195, "High": 4333.899902, "Low": 4
326.100098, "Close": 4327.0, "Adj Close": 4327.0, "Volume": 45190000.0, "CloseUSD": 3591.41}
{"Index": "IXIC", "Date": "1993-05-13", "Open": 677.909973, "High": 680.859985, "Low": 674.2
2998, "Close": 675.640015, "Adj Close": 675.640015, "Volume": 262530000.0, "CloseUSD": 675.6
40015}
{"Index": "IXIC", "Date": "1989-11-13", "Open": 456.799988, "High": 456.899994, "Low": 454.7
99988, "Close": 455.899994, "Adj Close": 455.899994, "Volume": 109200000.0, "CloseUSD": 455.
899994}
{"Index": "HSI", "Date": "1987-07-24", "Open": 3343.600098, "High": 3343.600098, "Low": 3343
.600098, "Close": 3343.600098, "Adj Close": 3343.600098, "Volume": 0.0, "CloseUSD": 434.6680
1274}
{"Index": "N225", "Date": "2021-04-16", "Open": 29789.08008, "High": 29789.08008, "Low": 296
21.83008, "Close": 29683.36914, "Adj Close": 29683.36914, "Volume": 49100000.0, "CloseUSD":
296.8336914}
{"Index": "000001.SS", "Date": "2009-08-12", "Open": 3255.985107, "High": 3255.986084, "Low"
: 3104.564941, "Close": 3112.718994, "Adj Close": 3112.718994, "Volume": 134400.0, "CloseUSD
": 498.03503904}
{"Index": "000001.SS", "Date": "2015-02-13", "Open": 3186.808105, "High": 3237.158936, "Low"
: 3182.793945, "Close": 3203.826904, "Adj Close": 3203.826904, "Volume": 261300.0, "CloseUSD
": 512.61230464}
```

Fig 3:-

## 4.4  Quality Assurance

Quality assurance is embedded throughout the implementation process, ensuring adherence to coding standards, promoting code readability, and enhancing overall system reliability. Rigorous testing, including unit tests, integration tests, and system tests, forms a pivotal part of quality assurance. The goal is to identify and rectify issues early in the development cycle, fostering a robust and maintainable codebase. Adherence to coding standards enhances collaboration among team members, facilitates future modifications, and contributes to the overall success and longevity of the Stock Market Real-Time Data Analysis project.

# Chapter 5

# **Standards Adopted**

## 5. Standards Adopted

## 5.1 Coding Standards

We have followed the coding standards recommended by the Software Engineering Institute (SEI) for this project. The coding standards followed are:

● Use meaningful variable and function names.
● Use comments to explain complex code.
● Use indentation for better readability.
● Write modular code.
● Avoid using hardcoded values.

These coding standards have helped us to write code that is easily understandable and maintainable.

## 5.2 Testing Standards

We have followed the testing standards recommended by the International Organization for Standardization (ISO) for this project. The testing standards followed are:

● Creation of test cases based upon requirements.
● Execution of test cases and record the results.
● Analyze the results and identify defects.
● Fix the defects and repeat the testing process.

These testing standards have helped us to ensure the quality of our project work and identify any defects in the system.

# Chapter 6

# <u>Conclusion and Future Scope</u>

## 6. Conclusion and Future Scope

## 6.1 Conclusion

In culmination, the Stock Market Real-Time Data Analysis project has achieved significant milestones by successfully implementing a robust real-time data analysis pipeline. The orchestrated synergy of Apache Kafka, Python, and Amazon Web Services (AWS) has resulted in a responsive and efficient system capable of processing financial data in real-time. The project's primary objective, addressing the critical need for timely and effective data analysis in financial markets, has been realized. The seamless flow of data from its source through the entire pipeline, facilitated by Kafka, has contributed to a transformative impact on decision-making processes within the dynamic realm of financial markets.

The successful implementation of the project underscores its relevance and applicability in addressing the challenges posed by the fast-paced nature of financial data. Stakeholders now have access to a tool that not only meets but exceeds expectations in terms of accuracy, responsiveness, and scalability. The real-time data analysis pipeline stands as a testament to the power of cutting-edge technologies and meticulous project planning.

## 6.2 Future Scope

### 6.2.1 Enhanced Scalability

While the current system demonstrates scalability, there is always room for improvement. Future enhancements could focus on optimizing the system's scalability to handle even larger data volumes and accommodate potential increases in user demand. This could involve further optimization of Apache Kafka configurations, AWS resource scaling strategies, and Python-based processing capabilities.

### 6.2.2 Integration of Machine Learning Algorithms

The future scope includes exploring the integration of machine learning algorithms for predictive analysis. Machine learning models could be employed to analyze historical market data, identify patterns, and provide predictive insights. This addition could elevate the system from reactive data analysis to a proactive tool that aids in anticipating market trends and making informed decisions.

### 6.2.3 Integration of Additional Financial Datasets

Expanding the project's scope involves integrating additional financial datasets to enhance the breadth and depth of analysis. This could include incorporating diverse market indices, commodity prices, or global economic indicators. The integration of new datasets would contribute to a more comprehensive understanding of the factors influencing financial markets, providing stakeholders with a holistic view.

### 6.2.4 User Interface Enhancements

Consideration should be given to user interface enhancements to improve the overall user experience. A user-friendly and intuitive interface could empower stakeholders to interact more seamlessly with the real-time data analysis platform. This may involve the development of dashboards, visualization tools, or customizable reporting features.

In conclusion, the Stock Market Real-Time Data Analysis project not only marks a successful venture into real-time financial data processing but also sets the stage for continuous innovation and improvement. The outlined future scope opens avenues for refining the system, expanding its capabilities, and ensuring its relevance in the ever-evolving landscape of financial markets.

# References

Fig 1: https://developer.confluent.io/courses/architecture/get-started/

Fig 2: https://medium.com/@ishaan.rawat611/real-time-data-pipeline-using-apache-kafka-glue-and-athena-on-aws-cloud-fb29eecb4788

Other References:-
1. https://kafka.apache.org/documentation/streams/
2. https://ap-south-1.console.aws.amazon.com/console/home?nc2=h_ct&src=header-signin&region=ap-south-1