

# Real-Time Stock Market Data Extraction with Apache Kafka: An End-to-End Data Engineering Project

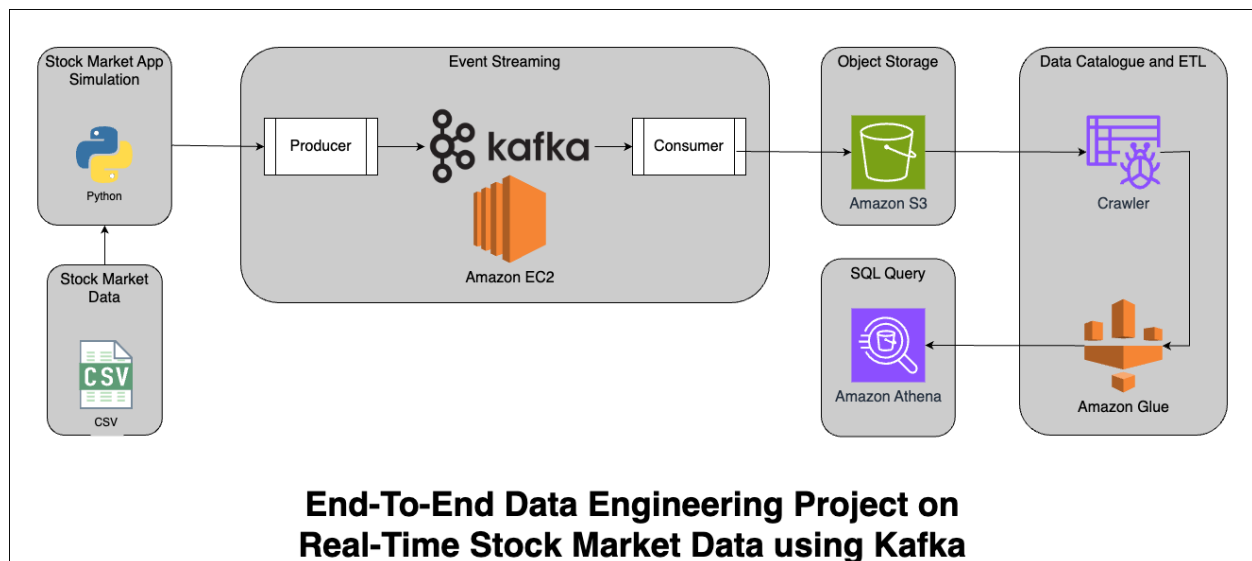
*By Randy A. Velasco*

## I. Introduction:

In the ever-evolving landscape of financial markets, real-time data analysis plays a pivotal role in decision-making for investors and traders. The Real-Time Stock Market Data Analysis project is a comprehensive demonstration of end-to-end data engineering capabilities, showcasing the integration of cutting-edge technologies to process and analyze stock market data in real-time.

## II. Project Overview:

This project leverages the power of Apache Kafka as a distributed streaming platform to ingest and process real-time stock market data. The seamless integration of Python for data manipulation and AWS services such as Glue, Athena, and S3 for data transformation and analysis enhances the project's scalability, efficiency, and accessibility.



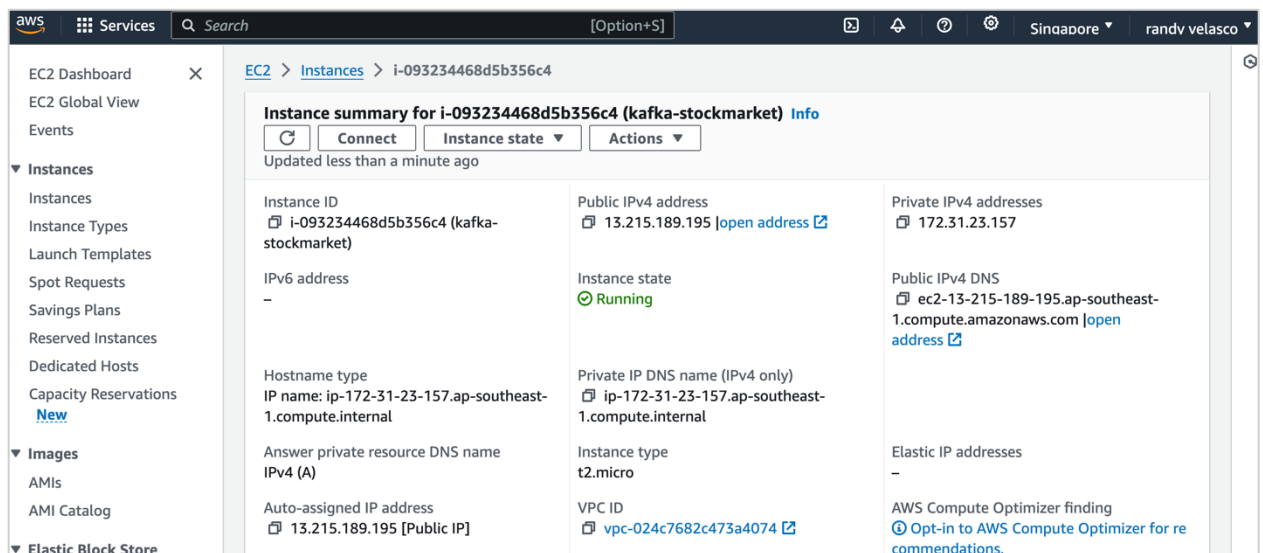
The workflow commences with the generation of a stock market data stream to simulate API behavior. This process utilizes a CSV file containing stock market data. Subsequently, a Python script is executed every second to randomly capture a single data row from the CSV. The captured data is then transmitted to Apache Kafka through a broker and producer, ultimately being published to a specific topic.

The data is consumed by a consumer and subsequently stored in JSON format on Amazon S3. Amazon Glue undertakes data cataloging tasks using a crawler, organizing and cataloging the data before saving it in a database. Amazon Athena is then employed to execute SQL queries, ensuring the accuracy of the entire workflow process.

### III. Actual Steps Taken:

#### 1. Configured Apache Kafka in Amazon EC2:

##### a) Launched Amazon EC2 Instance:



##### b) Connect to Amazon EC2 instance via SSH, download and install Apache Kafka and Java Runtime Environment

Apache Kafka is a distributed event streaming platform used for building real-time data pipelines and streaming applications. It will be used to capture streaming data generated by Python simulating a stock market API source. Kafka captures streaming data through producers that publish records to the “KSMP” topic. These records are then consumed by the consumer and saved in Amazon S3.

Java plays a central role in Apache Kafka, as Kafka is implemented in Java and is designed to run on the Java Virtual Machine (JVM).

Code:

```
> ssh -i "kafka_stockmarket.pem" ec2-user@ec2-13-215-189-195.ap-southeast-1.compute.amazonaws.com
> wget https://downloads.apache.org/kafka/3.6.1/kafka_2.13-3.6.1.tgz
> tar -xvf kafka_2.13-3.6.1.tgz
> sudo yum install java-1.8.0-openjdk
```

### c) Run Apache Zookeeper via SSH

The primary function of ZooKeeper in Kafka is to manage and coordinate the distributed nature of Kafka brokers and to maintain metadata and configuration information.

Code:

```
> bin/zookeeper-server-start.sh config/zookeeper.properties
```

```

er.server.auth.DigestAuthenticationProvider)
[2024-02-07 00:10:25,714] INFO zookeeper.DigestAuthenticationProvider.enabled =
true (org.apache.zookeeper.server.auth.DigestAuthenticationProvider)
[2024-02-07 00:10:25,722] INFO zookeeper.snapshot.trust.empty : false (org.apach
e.zookeeper.server.persistence.FileTxnSnapLog)
[2024-02-07 00:10:25,751] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,751] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,751] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,752] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,753] INFO (org.apache.zookeeper.server.ZooKeeperServer)
[2024-02-07 00:10:25,759] INFO Server environment: zookeeper.version=3.8.3-6ad6d3
64c70bcfb0de452d54ebefa3058098ab56, built on 2023-10-05 10:34 UTC (org.apache.zoo
keeper.server.ZooKeeperServer)

```

### d) Run Apache Kafka Server

Apache Kafka server functions as a distributed event streaming platform. It manages the ingestion, storage, and real-time streaming of data. Kafka enables producers to publish records to topics, and consumers subscribe to these topics, allowing scalable, fault-tolerant data streaming across multiple nodes. The server maintains a distributed commit log, ensuring durability and sequential storage.

Code:

```
> bin/kafka-server-start.sh config/server.properties
```

```

[ec2-user@ip-172-31-23-157 ~]$ bin/kafka-server-start.sh config/server.properties
OpenJDK 64-Bit Server VM warning: If the number of processors is expected to increase from one, then you should configure the number of parallel GC threads appropriately using -XX:ParallelGCThreads=N
[2024-02-07 00:12:56,610] INFO Registered kafka:type=kafka.Log4jController MBean (kafka.utils.Log4jControllerRegistration$)
[2024-02-07 00:12:56,321] INFO Setting -D jdk.tls.rejectClientInitiatedRenegotiation=true to disable client-initiated TLS renegotiation (org.apache.zookeeper.common.X509Util)
[2024-02-07 00:12:56,534] INFO Registered signal handlers for TERM, INT, HUP (org.apache.kafka.common.utils.LoggingSignalHandler)
[2024-02-07 00:12:56,543] INFO Starting (kafka.server.KafkaServer)
[2024-02-07 00:12:56,544] INFO Connecting to zookeeper on localhost:2181 (kafka.server.KafkaServer)
[2024-02-07 00:12:56,586] INFO [ZooKeeperClient Kafka server] Initializing a new session to localhost:2181. (kafka.zookeeper.ZooKeeperClient)
[2024-02-07 00:12:56,933] INFO Client environment:zookeeper.version=3.8.3-6ad6d364c7e0cb0fde452d54ebef3a585898ab56, built on 2023-10-05 10:34 UTC (org.apache.zookeeper.ZooKeeper)
[2024-02-07 00:12:56,593] INFO Client environment:host.name=ip-172-31-23-157.ap-southeast-1.compute.internal (org.apache.zookeeper.ZooKeeper)
[2024-02-07 00:12:56,594] INFO Client environment:java.version=1.8.0_392 (org.apache.zookeeper.ZooKeeper)
[2024-02-07 00:12:56,594] INFO Client environment:java.vendor=Red Hat, Inc. (org.apache.zookeeper.ZooKeeper)
[2024-02-07 00:12:56,594] INFO Client environment:java.home=/usr/lib/jvm/java-1.8.0-openjdk-1.8.0_392-b08-2.amzn2.0.i.x86_64/jre (org.apache.zookeeper.ZooKeeper)

```

e) Create Kafka topic and run Kafka Producer

In Kafka, a topic is a category or feed name to which records (messages) are published by producers and from which records are consumed by consumers. Topics serve as a way to organize and categorize the messages within the Kafka messaging system. Multiple producers can publish records to a topic, and multiple consumers can subscribe to and consume records from the same topic. Topics help in organizing and managing the flow of data within the Kafka cluster. For this project, the topic is named “KSMP”.

Code:

```
> bin/kafka-topics.sh --create --topic ksmp --bootstrap-server 13.215.189.195:9092 --replication-factor 1 --partitions 1
> bin/kafka-console-producer.sh --topic ksmp --bootstrap-server 13.215.189.195:9092
```

```
(base) randyvelasco@Randy-Air Kafka Stockmarket Project % ssh -i "kafka_stockma
rket.pem" ec2-user@ec2-13-215-189-195.ap-southeast-1.compute.amazonaws.com
Last login: Wed Feb  7 00:12:11 2024 from 112.201.188.203

_#_
~\####_ Amazon Linux 2
~~\#####\
~~\###| AL2 End of Life is 2025-06-30.
~~\#/
~~V~!!-->

~~~~/
~~./_/_/_/
_/_/m/_/

A newer version of Amazon Linux is available!

Amazon Linux 2023, GA and supported until 2028-03-15.
https://aws.amazon.com/linux/amazon-linux-2023/

1 package(s) needed for security, out of 1 available
Run "sudo yum update" to apply all updates.
-bash: warning: setlocale: LC_CTYPE: cannot change locale (UTF-8): No such file
or directory
[ec2-user@ip-172-31-23-157 ~]$ cd kafka_2.13-3.6.1
[ec2-user@ip-172-31-23-157 kafka_2.13-3.6.1]$ bin/kafka-topics.sh --create --topi
c ksmpp_demo --bootstrap-server 13.215.189.195:9092 --replication-factor 1 --par
titions 1
OpenJDK 64-Bit Server VM warning: If the number of processors is expected to inc
rease from one, then you should configure the number of parallel GC threads app
ropriately using -XX:ParallelGCThreads=N
```

### f) Run Kafka Consumer

A Kafka consumer is responsible for subscribing to Kafka topics, retrieving and processing records/messages produced by Kafka producers. For this project, the consumer retrieves data from the “ksmp” topic and save it to Amazon S3.

Code:

```
> bin/kafka-console-consumer.sh --topic ksmp --bootstrap-server 13.215.189.195:9092
```



The screenshot shows the AWS Glue console with the 'ksmp-crawler' configuration page. The left sidebar lists navigation options under 'AWS Glue', including 'Data Catalog' and 'Crawlers'. The main content area displays the crawler's properties in a table format.

Crawler properties			
Name	ksmp-crawler	IAM role	glue-admin-role
Description	-	Security configuration	-
Database	ksmp-database	Lake Formation configuration	-
State	READY	Table prefix	-
Maximum table threshold	-		

Buttons at the top right include 'Run crawler', 'Edit', and 'Delete'. The 'Last updated (UTC)' timestamp is February 7, 2024 at 03:16:36. Below the properties table is an 'Advanced settings' section. At the bottom, there are tabs for 'Crawler runs', 'Schedule', 'Data sources', 'Classifiers', and 'Tags'.

The screenshot shows the AWS Glue console with the 'Tables' page. The left sidebar lists navigation options under 'AWS Glue', including 'Data Catalog' and 'Tables'. The main content area displays a list of tables with a search bar and a table of results.

Tables (1) Last updated (UTC) February 7, 2024 at 03:15:22

View and manage all available tables.

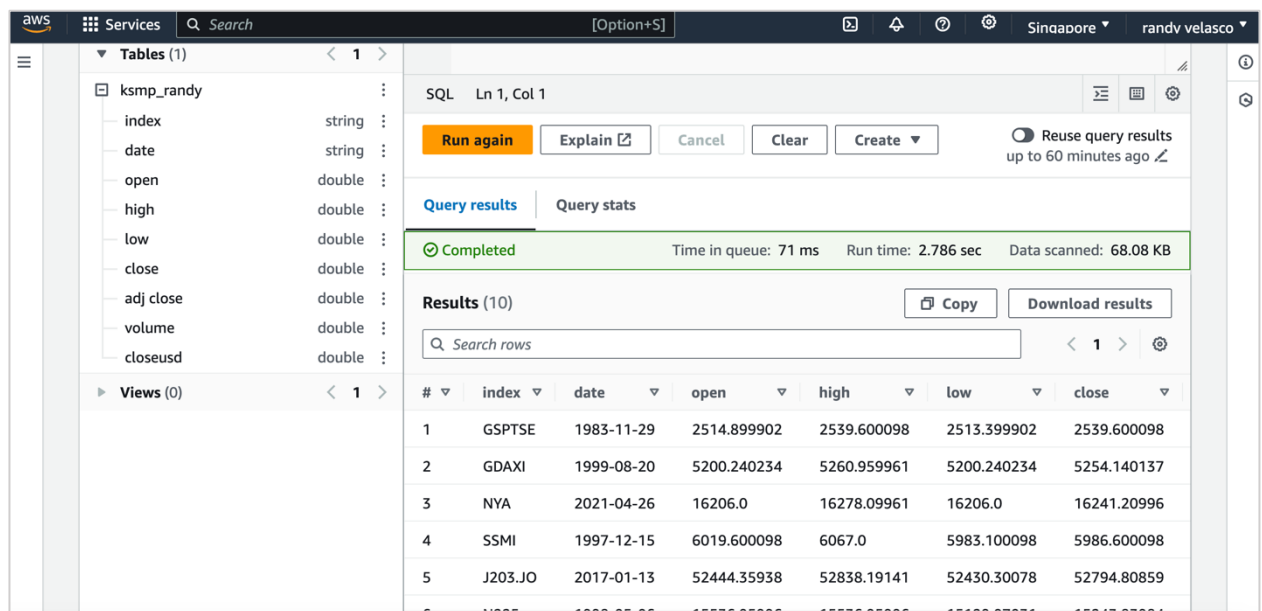
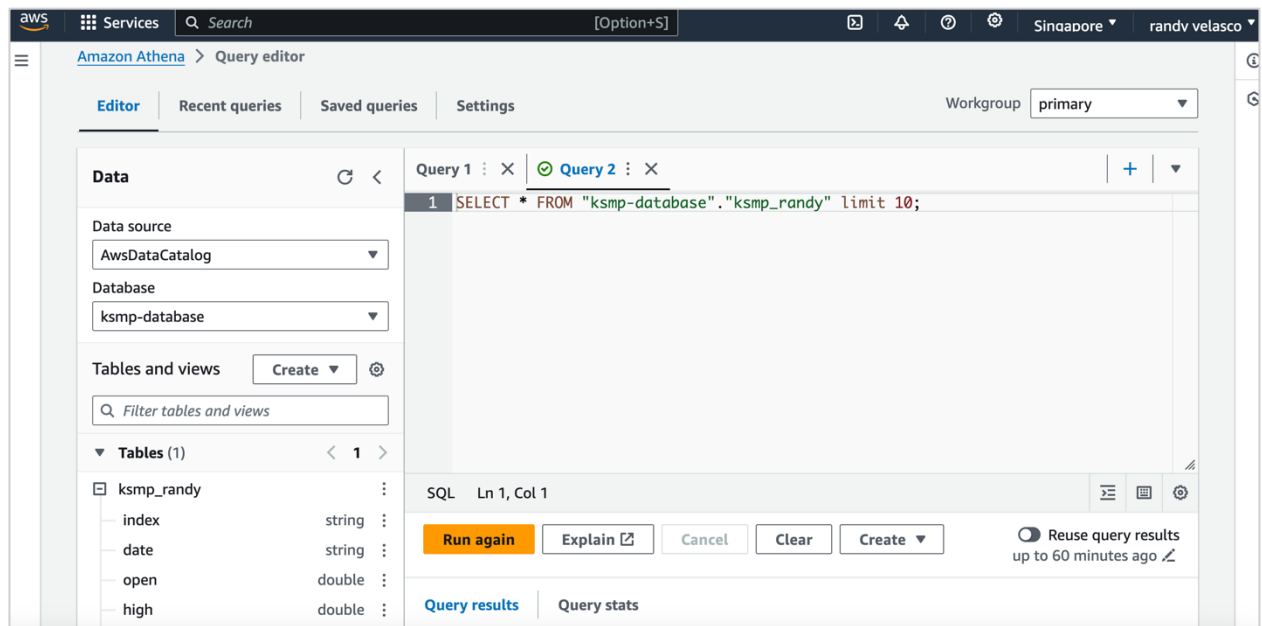
Filter tables

<input type="checkbox"/>	Name	Database	Location	Classific...	Depreca...	View data	Data quality
<input type="checkbox"/>	ksmp_randy	ksmp-database	s3://ksmp-randy	JSON	-	Table data	View data qual...

Buttons at the top right include 'Delete', 'Add tables using crawler', and 'Add table'.

#### 4. Perform SQL query using Amazon Athena

Use Amazon Athena to query tables and data cataloged by AWS Glue in Amazon S3 using standard SQL. This is to ensure the accuracy of the entire workflow process.



#### IV. Skills Showcased

- Apache Kafka: Showcase proficiency in setting up and configuring Kafka for real-time data streaming, ensuring reliable and scalable data ingestion.
- Python: Demonstrate skills in data manipulation and preprocessing using Python, ensuring the seamless integration of real-time data into the Kafka pipeline.
- Amazon Web Services (AWS):
  - AWS Glue: Illustrate expertise in automated ETL (Extract, Transform, Load) processes, facilitating the transformation of raw data into a structured and quarriable format.

- Amazon Athena: Showcase the ability to use serverless SQL queries for interactive analysis of data stored in Amazon S3, providing a cost-effective and efficient solution for ad-hoc querying.
- SQL: Demonstrate proficiency in writing SQL queries to perform complex data analysis tasks, showcasing the ability to derive valuable insights from real-time stock market data.

## **V. Conclusion:**

The Real-Time Stock Market Data Extraction project provides a hands-on demonstration of end-to-end data engineering skills, combining the power of Apache Kafka, Python, and AWS services to enable real-time data analysis in the dynamic and fast-paced realm of stock markets. This project not only highlights technical proficiency but also emphasizes the ability to derive meaningful insights from real-time data for informed decision-making in financial markets.

----- ( E N D ) -----