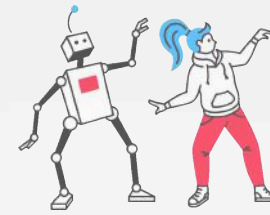# Project Overview

- Introduction to Real-Time Streaming

- Introduction to Apache Kafka

- Data Ingestion & Processing Pipeline

- Set up Kafka for streaming stock data

- Kafka Producer Setup

- Kafka Consumer Setup

- Set up PySpark Streaming in Databricks

- Define schema for stock price data

- Visualization and Streaming Data Sink

- Final Dashboard Visualization

- Libraries and Modules Used

- Conclusion

- Challenges of the Project

- Project Timeline

- References

# Introduction to Real-Time Streaming

**Real-Time Stream Engines:** They processes continuous streams of data as they arrive, enabling immediate insights and actions.

    **Examples:** Amazon Kinesis (Case: I) & Apache Kafka (Case: II)

**Use Case I :** E-Commerce Order Tracking

- **Data Ingestion:** Kinesis streams data from user activities like clicks and purchases.
- **Processing:** Real-time analytics for personalized recommendations and inventory updates.
- **Outcome:** Faster order tracking, instant alerts for delivery issues, and better customer experience.
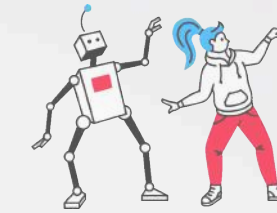
**Use Case II :** Stock Price Monitoring

- **Process:**
    - Kafka receives live stock prices from exchanges.
    - Processes data to calculate moving averages and price volatility.
    - Generates real-time alerts for price fluctuations.
- **Outcome:** Immediate visualization and decision-making for traders.

# Introduction to Apache Kafka

**Apache Kafka** is a distributed streaming platform used for building real-time data pipelines and streaming applications. It handles large volumes of data in real time, offering fault tolerance, scalability, and high throughput.

**Kafka Components:**

- **Producer:** Sends data (messages) to Kafka topics.
- **Consumer:** Reads data from Kafka topics.
- **Broker:** A Kafka server that manages topics and partitions.
- **Topic:** A category for messages, where data is stored.
- **Partition:** A division of a topic to distribute data across multiple servers.

**Example: Real-Time Stock Price Monitoring**

1. **Scenario:** A financial firm tracks real-time stock prices across various markets.
2. **Kafka Flow:**
   - Producer: Sends stock price data to the Kafka topic stock_prices.
   - Kafka Broker: Stores and manages the stock price data in partitions.
   - Consumer: A real-time application consumes the data to calculate metrics like moving averages or trigger alerts if a price threshold is crossed.
3. **Outcome:** Real-time alerts notify traders of price changes, while dashboards visualize the data for further analysis.
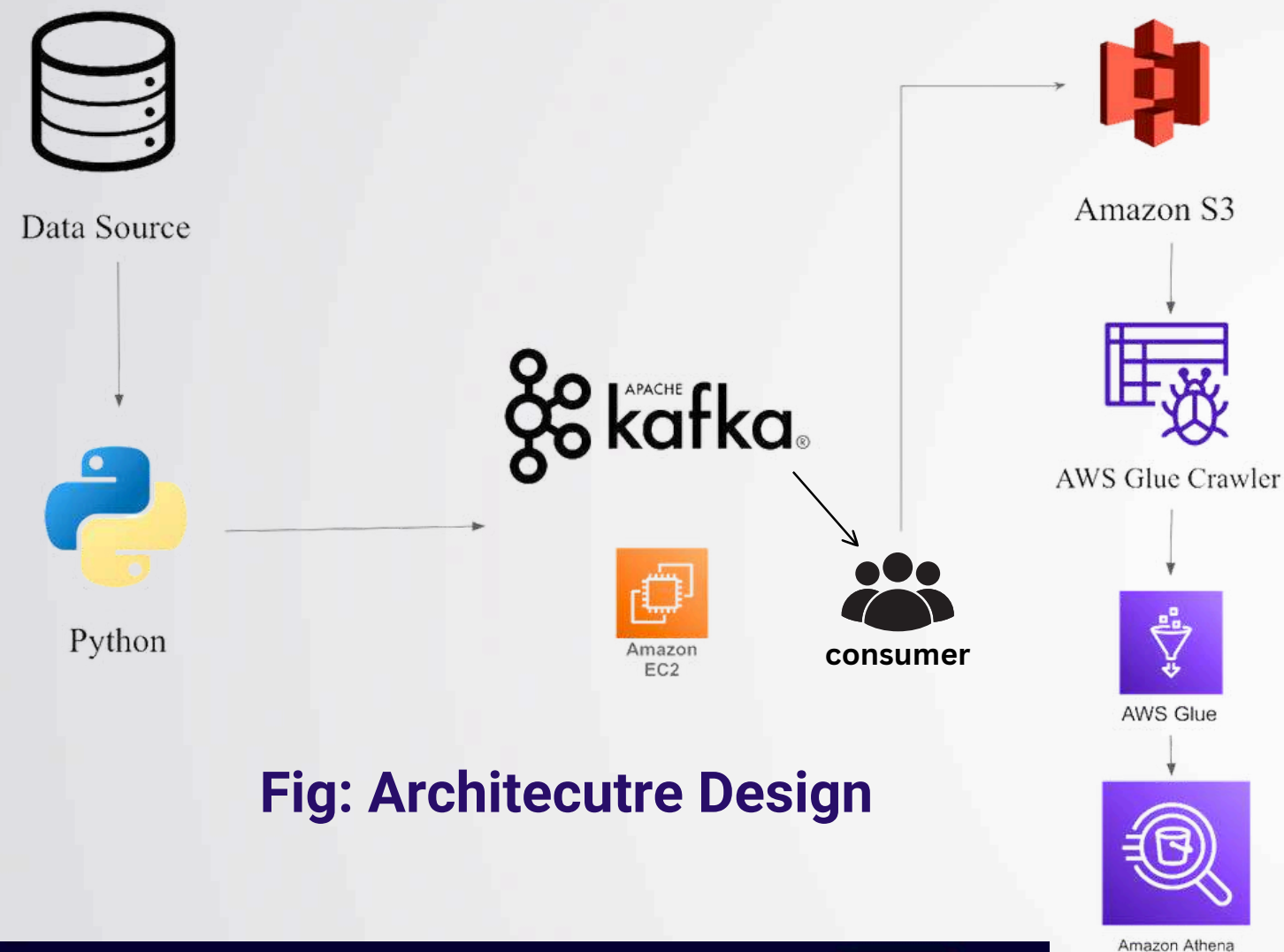
# Data Ingestion & Processing Pipeline


**Fig: Architecutre Design**


**Fig: Confluent Platform**

1. **Data Source:** The raw data originates from a data source, which could be a database, logs, or other systems.
2. **Python:** A Python application processes or formats the raw data before sending it to the streaming platform.
3. **Apache Kafka:** Acts as a distributed messaging system for streaming data in real time. Kafka ensures scalability and durability for event-driven data.
4. **Amazon EC2:** Hosts services to manage Kafka producers, consumers, or other processing tasks.
5. **Amazon S3:** Data is streamed and stored in Amazon S3, a scalable object storage service.
6. **AWS Glue Crawler:** Automatically detects and catalogs the schema of the data in S3.
7. **AWS Glue:** Performs ETL (Extract, Transform, Load) processes to prepare data for querying.
8. **Amazon Athena:** Provides serverless SQL query capabilities on the prepared data for insights and analytics.
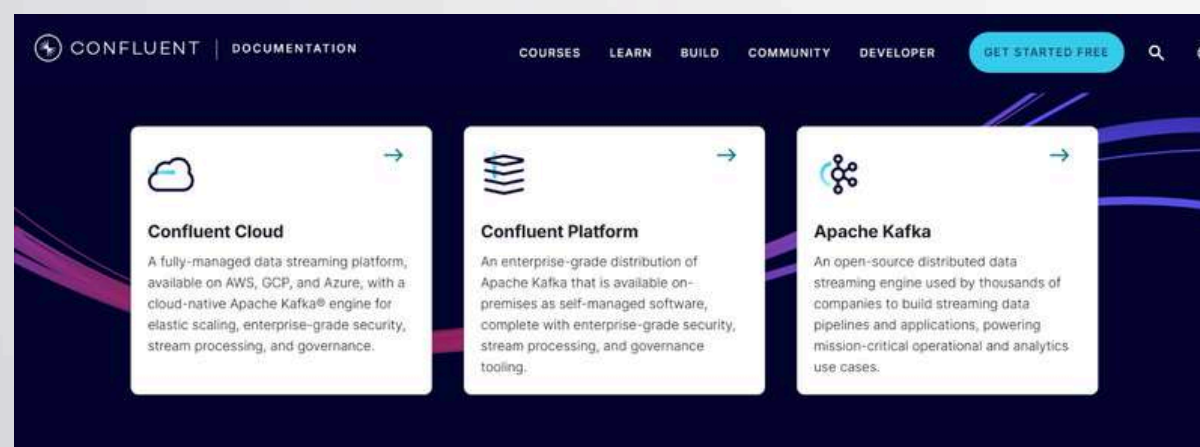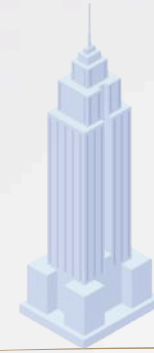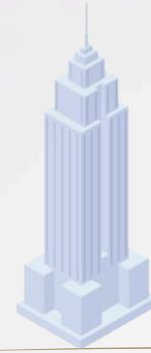
# Set up Kafka for streaming stock data

- **Prerequisites:** Ensure Java is installed and the system meets Kafka's requirements.

- **Download Kafka:** Get the Kafka binary from the official website and extract it.

- **Start Zookeeper:** Launch Zookeeper, which Kafka uses to manage metadata.

- **Start Kafka Broker:** Run the Kafka server to handle message distribution.

- **Configure Kafka:** Adjust the server.properties file for broker ID, log directories, and networking.

- **Create Topics:** Set up logical channels for message publishing and consumption.

- **Test the Setup:** Use Kafka's console tools to verify producer and consumer functionality.

- **Monitor and Scale:** Track performance with tools and add brokers to grow the cluster.

# Kafka Producer Setup

```python
def delivery_report(err, msg):
    """Callback for delivery reports."""
    if err is not None:
        print(f"Delivery failed for record {msg.key()}: {err}")
    else:
        print(f"Record {msg.key()} successfully produced to {msg.topic()} [{msg.partition()}]")

def send_stock_data():
    try:
        while True:
            # Randomly select a stock symbol from the list
            stock_symbol = random.choice(stock_symbols)

            # Generate random stock data with the selected stock symbol
            stock_message = {
                'Index': stock_symbol,  # Use the stock symbol here
                'Date': datetime.now().strftime("%Y-%m-%d %H:%M:%S"),
                'Open': round(random.uniform(100, 1500), 2),
                'High': round(random.uniform(1500, 2000), 2),
                'Low': round(random.uniform(50, 1499), 2),
                'Close': round(random.uniform(100, 1500), 2),
                'Adj_Close': round(random.uniform(100, 1500), 2),
                'Volume': random.randint(10000, 1000000),
                'CloseUSD': round(random.uniform(100, 1500), 2)
            }
```

```python
            # Produce message to Kafka
            kafka_producer.produce(
                'stocks_analysis',
                key=stock_symbol,  # The stock symbol serves as the key
                value=json.dumps(stock_message),
                callback=delivery_report
            )
            print(f"Sent: {stock_message}")
            kafka_producer.poll(0)
            time.sleep(3)
    except KeyboardInterrupt:
        print("Stopping data production.")
    finally:
        kafka_producer.flush()  # Ensure all messages are delivered

send_stock_data()
```

The producer generates stock price data and sends it to a Kafka topic.

**Setup and Configuration:**
- Imports required libraries (Producer from confluent_kafka, json, random, etc.).
- Configures the Kafka producer with necessary parameters like bootstrap servers.

**Simulated Stock Data:**
- Generates random stock data for symbols like MSFT, AMZN, etc.
- Fields include Index, Date, Open, High, Low, Close, Adj_Close, Volume, and CloseUSD.

**Publishing to Kafka:**
- Each record is serialized into JSON format.
- Published to the Kafka topic stocks_analysis using the produce method.
- A delivery callback (delivery_report) is used to handle success or error responses.

**Continuous Streaming:** The producer continuously sends new stock data every 3 seconds (time.sleep(3)). Terminates gracefully on a keyboard interrupt.

- *snippest from the code*

# Kafka Consumer Setup

### Integration with Kafka:

- **Reading Data from Kafka in PySpark:** PySpark can act as a Kafka consumer to ingest real-time streaming data, ensuring low latency for smooth processing.

### Parsing JSON Data:

- **Convert JSON Strings:** Use PySpark's from_json function to parse JSON strings into structured data, based on a predefined schema.
- **Transform into Tabular Format:** Convert the raw JSON data into a structured DataFrame for more advanced analysis and real-time processing.

```python
# Function to handle consuming messages and processing them
def process_kafka_messages():
    try:
        while True:
            message = kafka_consumer.poll(1.0)  # Poll for new messages

            if message is not None and not message.error():
                try:
                    # Decode key and value safely and log for debugging
                    message_value = message.value().decode("utf-8") if message.value() else "{}"
                    print(f"Decoded Kafka message: {message_value}")  # Debugging: Log the decoded message

                    # Parse message value as JSON
                    try:
                        record = json.loads(message_value)
                    except json.JSONDecodeError as e:
                        print(f"Error decoding JSON: {e}")
                        continue

                    # Extract fields from the record with defaults to ensure robustness
                    index = record.get("Index", "UnknownIndex")  # Correct field name
                    date = record.get("Date", "UnknownDate")
                    open_price = float(record.get("Open", 0.0))
                    high_price = float(record.get("High", 0.0))
                    low_price = float(record.get("Low", 0.0))
                    close_price = float(record.get("Close", 0.0))
                    adj_close = float(record.get("Adj_Close", 0.0))
                    volume = int(record.get("Volume", 0))
                    close_usd = float(record.get("CloseUSD", 0.0))
```

```python
                    # Create a DataFrame for the single record
                    data_frame = spark_session.createDataFrame(
                        [(index, date, open_price, high_price, low_price, close_price, adj_close, volume, close_usd)],
                        schema=data_schema
                    )

                    # Display the DataFrame in Databricks (optional for debugging/validation)
                    data_frame.show()

                    # Write the data to a Delta table in append mode
                    try:
                        data_frame.write.format("delta").mode("append").save(delta_table_location)
                        print("Data written to Delta table successfully.")
                    except Exception as e:
                        print(f"Error writing to Delta table: {e}")

                except Exception as e:
                    print(f"Error processing Kafka message: {e}")
            elif message is not None and message.error():
                print(f"Kafka error: {message.error()}")
```

- *snippest from the code*

# PySpark Streaming in Databricks

## Why Databricks? databricks

- **Unified Platform:** Databricks integrates PySpark with collaborative tools for streamlined real-time data streaming and analysis.
- **Scalable Performance:** It handles large streaming datasets efficiently with automatic resource scaling.
- **Ease of Use:** Simplifies cluster management and monitoring, making streaming pipelines easier to build and maintain.

## Steps Involved databricks

- **Import Libraries:** Load PySpark and Kafka-related libraries for handling streaming data and schemas.
- **Configure Spark Session:** Set up a Spark session with the Kafka broker details (bootstrap.servers).
- **Subscribe to Topic:** Read data from the stock_prices Kafka topic using PySpark's structured streaming API.
- **Define Schema:** Create a schema to parse the incoming JSON data for structured processing.
- **Process Stream:** Transform and display real-time data using PySpark's streaming query capabilities.

# Define schema for stock price data

## Why Define a Schema?

- **Simplifies Parsing:** Makes data easier to process.
- **Ensures Consistency:** Guarantees uniform data structure.

## Schema Fields:

- **Index:** *String type, nullable (True).*
- **Date:** *String type, nullable (True).*
- **Open:** *Float type, nullable (True).*
- **High:** *Float type, nullable (True).*
- **Low:** *Float type, nullable (True).*
- **Close:** *Float type, nullable (True).*
- **Adj_Close:** *Float type, nullable (True).*
- **Volume:** *Integer type, nullable (True).*
- **CloseUSD:** *Float type, nullable (True).*

```python
data_schema = StructType([
    StructField("Index", StringType(), True),
    StructField("Date", StringType(), True),
    StructField("Open", FloatType(), True),
    StructField("High", FloatType(), True),
    StructField("Low", FloatType(), True),
    StructField("Close", FloatType(), True),
    StructField("Adj_Close", FloatType(), True),
    StructField("Volume", IntegerType(), True),
    StructField("CloseUSD", FloatType(), True)
])
```

*Defining schema*

# Visualization and Streaming Data Sink

- *Streaming Data Sink (Top Left)*
- *Stock Symbol Filtering (Top Middle)*
- *Real-Time Aggregation (Moving Averages & Volatility) (Top Right)*
- *Daily Summary Calculation (Bottom Left)*
- *Price Change and Volatility Calculation (Bottom Middle)*
- *Stock Price Alerts (Bottom Right)*

# Final Dashboard Visualization

- *Bar Chart (Top Left): Shows the total trading volume for stocks (AAPL, AMZN, GOOG, MSFT), identifying the most actively traded stock.*
- *Pie Chart (Top Right): Displays the percentage contribution of each stock to overall volatility.*
- *Line Chart (Middle Left): Compares maximum moving averages (price trends) and maximum volatilities for stocks.*
- *Bar chart (Bottom Left): Correlates volatility with moving averages to analyze stock price stability and trends.*
- *Summary Table (Middle Right): Provides key metrics like average closing price for each stock by date.*
- *Window Data Table (Bottom Right): Shows real-time streaming batches with processing windows and metrics.*
- *Raw Data Table (Bottom Right Corner): Displays unprocessed stock data like Open, High, and Volume for analysis.*

# Conclusion

- **Confluent Kafka Integration:** Real-time stock price streaming through Kafka's message queuing.
- **PySpark Data Processing:** Utilizes PySpark for efficient and scalable data processing and analytics.
- **Interactive Dashboard:** Provides visualizations (bar charts, pie charts, line graphs) for tracking stock performance and trends.
- **Stock Symbol Coverage:** Currently limited, but can be expanded for broader market analysis.
- **Scalability Challenges:** Potential issues with scaling, addressed in future enhancements.
- **Future Enhancements:** Expanding stock symbols, adding advanced analytics, and improving scalability for better performance.

# Future Work

**Enhanced Data Analytics**

Integrate machine learning for price predictions, sentiment analysis from news or social media, and anomaly detection to provide smarter insights for traders.

**Real-Time Alerts**

Implement an alert system that notifies users of significant market changes, such as price shifts or volatility spikes, enabling quick decision-making.

**Cloud Deployment**

Deploy the system on cloud platforms to handle larger datasets, improve scalability, and provide easier access for a broader user base.

# Challenges of the Project

**Data Quality and Consistency**

Validating and cleaning real-time data for accuracy.

**Scalability**

Configuring Kafka and PySpark to handle large data volumes.

**Latency**

Minimizing delays to ensure real-time processing.

**Fault Tolerance and Reliability**

Ensuring system recovery with Kafka replication and PySpark checkpoints.

**Complex Event Processing (CEP)**

Detecting real-time patterns and anomalies efficiently.

**Real-Time Visualization**

Displaying live data and insights through dashboards.

# Project Timeline

Project Start

Project Completion

**DAY 1**

**DAY 2**

**DAY 3**

**DAY 4**

**DAY 5**

**DAY 6**

**DAY 7**

**DAY 9**

**Kafka Setup and Configuration**

**Kafka Producer and Data Streaming**

**PySpark Environment Setup**

**Schema Definition and Data Parsing**

**Data Consumption and Real–Time Processing**

**Visualization and Fault Tolerance**

**Final Testing, Optimization, and Reporting**

**Final Files Delivery**

- Installed Apache Kafka locally, configured Kafka on cloud (if applicable), and created the stock_prices topic for streaming data.

- Developed a Kafka producer script to simulate real–time stock prices. The producer publishes stock data at regular intervals to the stock_prices topic.

- Set up PySpark on Databricks (or locally). Verified successful integration between PySpark and Kafka for real–time streaming.

- Defined the schema for stock price data, including fields like stock symbol, price, and timestamp. Implemented PySpark logic to parse incoming Kafka data into a structured format.

- Configured PySpark streaming to consume data from the Kafka stock_prices topic. Implemented real–time metrics like moving averages and price fluctuations.

- Developed a dashboard to visualize live stock price data, moving averages, and alerts.
- Set up Delta tables for persistent storage and configured Kafka replication and PySpark checkpointing for fault tolerance.

- **Dec 9:** Submit project
- **Dec 9:** Presentation

**Legend:**

🔴 Bharat Dhungana

🔵 Samir Bhandari

🟡 Satish Kandel

🔴 Bikash Thapa

🟢 Aravind

# References

- Confluent Documentation. (n.d.). Confluent documentation. Retrieved December 8, 2024, from *https://docs.confluent.io/*

- Apache Spark. (n.d.). Cluster mode overview. Retrieved December 8, 2024, from *https://spark.apache.org/docs/latest/cluster-overview.html*

- Confluent. (n.d.). confluent_kafka API. Retrieved December 9, 2024, from *https://docs.confluent.io/platform/current/clients/confluent-kafka-python/html/index.html*

- Iwuchukwu, C. (2024, April 23). Analyzing stock prices using PySpark. Medium. Retrieved December 7, 2024, from *https://medium.com/@ceejayiwufitness/analyzing-stock-prices-using-pyspark-acdaef8a5511*

- ResearchGate. (n.d.). Real-time streaming data analysis using Spark. Retrieved December 9, 2024, from *https://www.researchgate.net/publication/322674233_Real time_Streaming_Data_Analysis_using_Spark*

Team Everest

# Thank You
## For Your Attention

## Any Queries?

### We are open for discussion.

Access this ppt via: **bit.ly/aml2054**

**BDM 3603: Big Data Framework | Real-Time Stock Price Analysis : PySpark and Kafka Streaming**