# Step-by-Step Workflow

July 24, 2025

**End-to-End Data Streaming Pipeline with Kafka, AWS S3, Glue, and Athena** Prepared by Your Name

# Introduction

This document outlines the step-by-step process of building an end-to-end data streaming pipeline using Apache Kafka on AWS. The pipeline streams data from a local machine to a Kafka broker on an EC2 instance, stores it in Amazon S3, catalogs it with AWS Glue, and enables querying with Amazon Athena. This professional workflow is designed to be clear, reproducible, and scalable for enterprise use.

# Step-by-Step Workflow

### Set Up an EC2 Instance for Kafka Broker

○ **Launch EC2 Instance**: Create an Amazon EC2 instance using Amazon Linux or Ubuntu with sufficient resources (e.g., t2.micro for testing).

○ **Install Java**: Install Java, a prerequisite for Kafka and Zookeeper:

```
1  sudo yum install java-1.8.0-openjdk -y  % For Amazon Linux
2
```

○ **Install Kafka**: Download and set up Apache Kafka:

```
1  wget https://archive.apache.org/dist/kafka/3.2.0/kafka_2.13-3.2.0.tgz
2  tar -xzf kafka_2.13-3.2.0.tgz
3  cd kafka_2.13-3.2.0
4
```

○ **Configure Security Group**: Allow traffic on:
  - Port 9092 (Kafka broker)
  - Port 2181 (Zookeeper)

○ **Start Zookeeper and Kafka**:

```
1  bin/zookeeper-server-start.sh config/zookeeper.properties &
2  bin/kafka-server-start.sh config/server.properties &
3
```

### Connect to EC2 Instance via SSH

○ **Access EC2**: Use the provided SSH key (.pem file) to connect from your local machine:

```
1  ssh -i "your-key.pem" ec2-user@your-ec2-public-ip
2
```

○ **Verify Connection**: Ensure you can access the EC2 instance terminal.

### Create a Kafka Topic

○ **Create Topic**: On the EC2 instance, create a Kafka topic named `demo_test`:

```
1  bin/kafka-topics.sh --create --topic demo_test --bootstrap-server localhost:9092 --
     partitions 1 --replication-factor 1
2
```

○ **Verify Topic**: Check the topic creation:

```
1  bin/kafka-topics.sh --list --bootstrap-server localhost:9092
2
```

### Set Up a Python Producer on Local Machine

○ **Install Libraries**: Install required Python libraries locally:

```
1  pip install kafka-python pandas
2
```

○ **Create Producer Script**: Write a Python script to read from `indexProcessed.csv` and send data to the Kafka topic:

```
1  from kafka import KafkaProducer
2  from json import dumps
3  from time import sleep
4  import pandas as pd
5
```

```python
6  # Load CSV data
7  df = pd.read_csv('indexProcessed.csv')
8
9  # Initialize Kafka Producer
10 producer = KafkaProducer(
11     bootstrap_servers=['your-ec2-public-ip:9092'],
12     value_serializer=lambda x: dumps(x).encode('utf-8')
13 )
14
15 # Send data to Kafka
16 while True:
17     data = df.sample(1).to_dict(orient="records")[0]
18     producer.send('demo_test', value=data)
19     sleep(1)
20
```

○ **Run Producer**: Execute the script to stream data to Kafka.

### Set Up a Kafka Consumer on EC2

○ **Run Console Consumer**: On the EC2 instance, verify data reception:

```
1  bin/kafka-console-consumer.sh --topic demo_test --bootstrap-server localhost:9092 --
       from-beginning
2
```

○ **Confirm Data Flow**: Ensure messages from the local producer are received.

### Store Consumed Data in Amazon S3

○ **Create Consumer Script**: Write a Python script to consume Kafka messages and upload to S3:

```python
1  from kafka import KafkaConsumer
2  import boto3
3  from json import loads
4
5  # Initialize Kafka Consumer
6  consumer = KafkaConsumer(
7      'demo_test',
8      bootstrap_servers=['your-ec2-public-ip:9092'],
9      auto_offset_reset='earliest',
10     value_deserializer=lambda x: loads(x.decode('utf-8'))
11 )
12
13 # Initialize S3 client
14 s3_client = boto3.client('s3')
15
16 # Consume and upload to S3
17 for message in consumer:
18     data = message.value
19     s3_client.put_object(
20         Bucket='your-bucket-name',
21         Key=f'data/message_{message.offset}.json',
22         Body=str(data)
23     )
24
```

○ **Configure AWS Credentials**: Set up AWS CLI or environment variables for S3 access.
○ **Run Consumer**: Execute the script to store data in S3.

### Set Up AWS Glue Crawler

○ **Create Crawler**: In AWS Glue, configure a crawler to scan the S3 bucket (e.g., s3://your-bucket-name/data/).
○ **Run Crawler**: Execute the crawler to create a table in the AWS Glue Data Catalog.
○ **Verify Table**: Confirm the table is available in the Glue Data Catalog.

### Query Data Using Amazon Athena

○ **Access Athena**: In the AWS Management Console, navigate to Amazon Athena.
○ **Select Database**: Choose the database created by the Glue crawler.

3

○ **Query Data**: Run SQL queries to analyze the data:

```sql
SELECT * FROM your_table_name LIMIT 10;

```

○ **Verify Results**: Confirm the streamed data is queryable.

## Architecture Overview

The pipeline follows this flow:
○ Local Python producer streams CSV data to a Kafka topic (`demo_test`) on EC2.
○ A Kafka consumer reads messages and uploads them to an S3 bucket.
○ AWS Glue Crawler catalogs the S3 data into a table.
○ Amazon Athena enables SQL-based querying of the cataloged data.

## Conclusion

This pipeline demonstrates a robust, end-to-end data streaming solution using Kafka and AWS services. It is scalable and can be extended with automation or additional integrations as needed.