

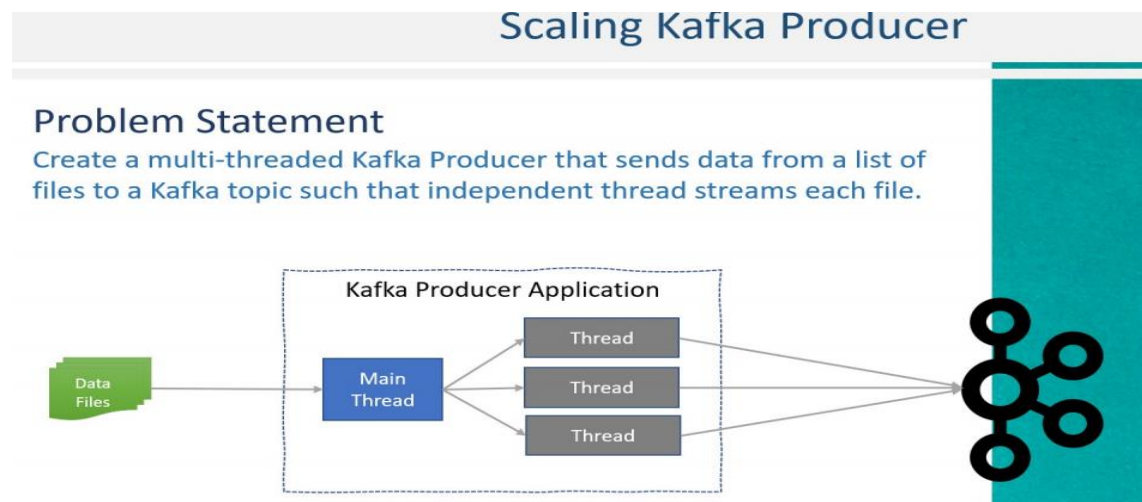
In data acquisition there were three steps:

1. Collecting the data from different sources
2. Using Kafka to read the data file into Kafka cluster.
3. Integrated Kafka cluster to Spark.

Here step 2 and 3 are described:

The task is to read Data from CSV file, collected from different sources, one of them is Example: [Hstdata.com](https://hstdata.com).

We will be using Kafka to take this data and create a producer to read the data from the csv file.



The output will be printed in the Kafka consumer console, We wanted to integrate it with spark, but were unsuccessful in doing so. The documentation we followed given [here](#).

Step 1:

Installing and configuring Kafka and spark. My focus will be on Kafka.

1. Download Kafka scala latest version from [here](#), you will be directed [here](#), where you can get the tgz file.
2. Now unzip it and set KAFKA_HOME as your Environment Variable, also specify the path for Kafka bin directory.
3. You will be needing the [Java JDK](#) and [Maven](#) as we will be using that as well. For the IDE we will be using [IntelliJ](#).
4. After installing all of these we are good to go. Any problem? Refer [here](#).

Step 2:

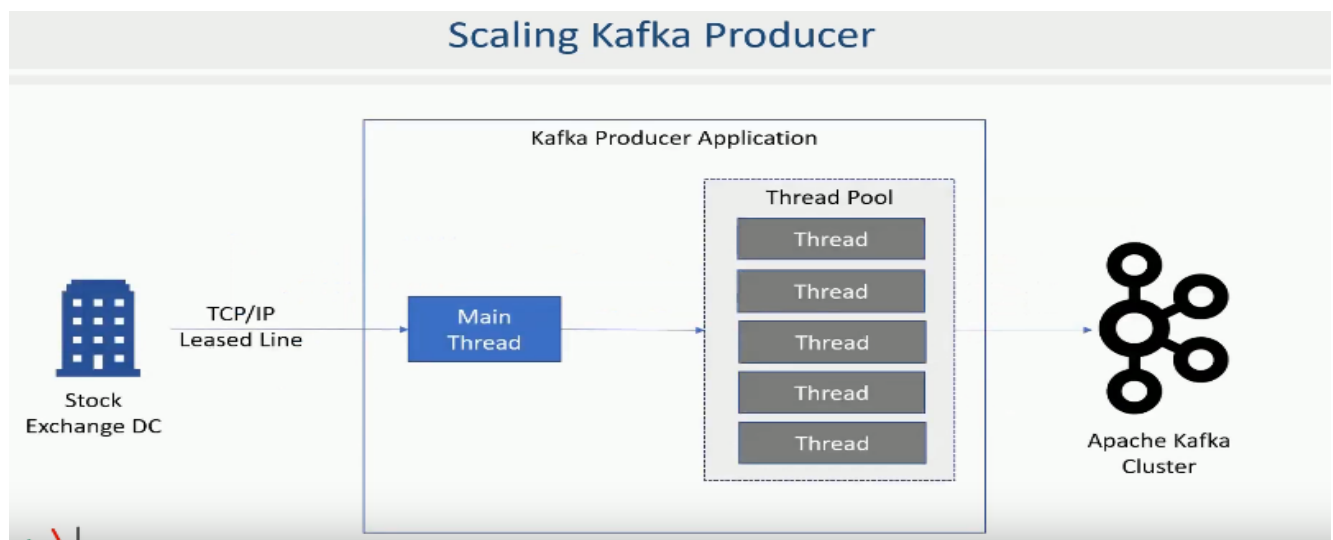
Now let us go ahead and work on our Kafka project. We learn about scaling the application (Kafka Producer).

1. We can send data from multiple CSV file, here event and it is the Kafka broker who receives the messages and acknowledge the successful receipt of the message.

So, if you are send 10000 of messages you can increase the number of Kafka brokers in the Kafka cluster.

As single Kafka broker can handle 100 or maybe 1000 of message per second. Meanwhile you can send more messages to the Kafka producer. This provides Linear scalability.

2. Here we will be using Multi-threaded Kafka producer.
 - a. The application receives the tick by tick data packets from stock exchange over a TCP/IP socket. The packets arrive at high frequency we create a multi threaded data handler.
 - b. The main thread listens to the socket and read the data packet as they arrives and immediately hands over the data packet to the thread pool to be passed on to Kafka cluster.
 - c. The main thread then goes back to reading the data again.
 - d. The threads in the thread pool are responsible for uncompressing the data packets and reading the content in the packets, validating the message and sending it to kafka broker.
 - e. So if we are passing three data files it must create 3 threads for each data file. Each thread is responsible for reading the records from each file and sending it to Kafka broker. Here we don't want to create multiple producer instances.



Step 3:

1. In this step we will set up the pom.xml file with all the dependencies.

```
<build>
  <plugins>
    <!-- Maven Compiler Plugin-->
    <plugin>
      <groupId>org.apache.maven.plugins</groupId>
      <artifactId>maven-compiler-plugin</artifactId>
      <version>3.8.0</version>
      <configuration>
        <source>${java.version}</source>
        <target>${java.version}</target>
      </configuration>
    </plugin>
  </plugins>
</build>

<dependencies>
  <!-- Apache Kafka Clients-->
  <dependency>
    <groupId>org.apache.kafka</groupId>
    <artifactId>kafka-clients</artifactId>
    <version>${kafka.version}</version>
  </dependency>

  <!-- Apache Log4J2 binding for SLF4J -->
  <dependency>
    <groupId>org.apache.logging.log4j</groupId>
    <artifactId>log4j-slf4j-impl</artifactId>
    <version>2.11.0</version>
  </dependency>
</dependencies>
```

2. The [Log4J](#) for logging results.

```
<Configuration status="ERROR">
  <Appenders>
    <Console name="stdout" target="SYSTEM_OUT">
      <PatternLayout pattern="%d] (%c) %p - %m %n"/>
    </Console>
  </Appenders>
  <Loggers>
    <Root level="error">
      <AppenderRef ref="stdout"/>
    </Root>
    <Logger name="org.apache.kafka.clients" level="warn" additivity="false">
      <AppenderRef ref="stdout"/>
    </Logger>
    <Logger name="guru.learningjournal.kafka.examples" level="trace">
```

```
additivity="false">
    <AppenderRef ref="stdout"/>
</Logger>
</Loggers>
</Configuration>
```

3. Now let us set up the Kafka scripts so the we don't have to run the cmd prompt again and again and we can use the IDE for all of the steps.

Here are the codes:

```
Start zookeeper: zookeeper-server-start.bat %KAFKA_HOME%\config\zookeeper.properties
Kafka Server 0: kafka-server-start.bat %KAFKA_HOME%\config\server-0.properties
```

```
Kafka Server 1: kafka-server-start.bat %KAFKA_HOME%\config\server-1.properties
```

```
Kafka Server 2: kafka-server-start.bat %KAFKA_HOME%\config\server-2.properties
```

```
The Kafka Topic: kafka-topics.bat --create --zookeeper localhost:2181 --topic
the-finance-data-topic --partitions 5 --replication-factor 3
```

And finally Kafka consumer to print the output in console;

```
Kafka Consumer: kafka-console-consumer.bat --bootstrap-server localhost:9092 --topic
the-finance-data-topic --from-beginning
```

4. Now we need to set up the App Config:

```
class AppConfigs {
    final static String applicationID = "Multi-Threaded-Producer";
    final static String topicName = "the-finance-3-data-topic";
    final static String kafkaConfigFileLocation = "kafka.properties";
    final static String[] eventFiles = {"data/Raw_data.csv", "data/Raw_data.csv"};
}
```

5. Now we should keep some producer level config outside the source code, here broker coordinates. This arrangement allows us to deploy the application and connect to any Kafka cluster by changing the connection details in the properties file.

```
bootstrap.servers=localhost:9092,localhost:9093
```

6. Now let us create the Dispacter class which run the runnable interface, this allows us to use the instance of this class as a separate thread.
 - a. This takes Kafka producer instance, Topic Name and the File Location.

```

public class Dispatcher implements Runnable {
    private static final Logger logger = LogManager.getLogger();
    private String fileLocation;
    private String topicName;
    private KafkaProducer<Integer, String> producer;

    Dispatcher(KafkaProducer<Integer, String> producer, String topicName, String
fileLocation) {
        this.producer = producer;
        this.topicName = topicName;
        this.fileLocation = fileLocation;
    }

    @Override
    public void run() {
        logger.info("Start Processing " + fileLocation);
        File file = new File(fileLocation);
        int counter = 0;

        try (Scanner scanner = new Scanner(file)) {
            while (scanner.hasNextLine()) {
                String line = scanner.nextLine();
                producer.send(new ProducerRecord<>(topicName, null, line));
                counter++;
            }
            logger.info("Finished Sending " + counter + " messages from " + fileLocation);
        } catch (FileNotFoundException e) {
            throw new RuntimeException(e);
        }
    }
}

```

7. Another class, DispatcherDemo class allows us to create thread and start them.
 - a. This is a main thread, which will produce a single instance of Kafka producer.
 - b. It then creates two threads using the runnable dispatcher and pass the producer instance, along with the Topic name and the File Location to the Dispatcher thread.
 - c. Each Dispatcher thread will send all the line in the file to a given topic.
 - d. Main thread waits for all the thread to be processed and close the Producer and Terminate.

```

public class DispatcherDemo {
    private static final Logger logger = LogManager.getLogger();
    public static void main(String[] args){

        Properties props = new Properties();
        try{
            InputStream inputStream = new
FileInputStream(AppConfigs.kafkaConfigFileLocation);
            props.load(inputStream);

```

```

        props.put(ProducerConfig.CLIENT_ID_CONFIG, AppConfigs.applicationID);
        props.put(ProducerConfig.KEY_SERIALIZER_CLASS_CONFIG,
IntegerSerializer.class);
        props.put(ProducerConfig.VALUE_SERIALIZER_CLASS_CONFIG,
StringSerializer.class);
    }catch (IOException e){
        throw new RuntimeException(e);
    }

    KafkaProducer<Integer,String> producer = new KafkaProducer<>(props);
    Thread[] dispatchers = new Thread[AppConfigs.eventFiles.length];
    logger.info("Starting Dispatcher threads...");
    for(int i=0;i<AppConfigs.eventFiles.length;i++){
        dispatchers[i]=new Thread(new
Dispatcher(producer,AppConfigs.topicName,AppConfigs.eventFiles[i]));
        dispatchers[i].start();
    }

    try {
        for (Thread t : dispatchers) t.join();
    }catch (InterruptedException e){
        logger.error("Main Thread Interrupted");
    }finally {
        producer.close();
        logger.info("Finished Dispatcher Demo");
    }
}
}

```

Step 4:

Run the Application:

1. Start the zookeeper script
2. Start all Kafka brokers, server 0,1,2.
3. Create a topic.
4. Execute the producer application.
5. Excute the Kafka consumer script to see the output of the console.

The screenshot shows an IDE with a project structure on the left. The project is named 'guru.learningjournal.kafka.examples' and contains a 'src' directory with 'main' and 'resources' subdirectories. The 'main' directory contains 'AppConfigs', 'Dispatcher', and 'DispatcherDemo'. The 'resources' directory contains 'log4j2.xml'. The 'DispatcherDemo' class is selected, and its code is visible in the editor. The code imports 'org.apache.logging.log4j.LogManager' and 'org.apache.logging.log4j.Logger', and 'java.io.*' and 'java.util.*'. It defines a 'DispatcherDemo' class with a 'main' method. The 'main' method logs the start of the Dispatcher threads, processes data from 'data/Raw_data.csv', and logs the completion of sending 1008 messages. The logs in the bottom pane show the execution of the 'DispatcherDemo' class, with timestamps and log levels (INFO) indicating the start and end of the process. The process finished with exit code 0.

```
import org.apache.logging.log4j.LogManager;
import org.apache.logging.log4j.Logger;

import java.io.*;
import java.util.*;

public class DispatcherDemo {
    private static final Logger logger = LogManager.getLogger();
    public static void main(String[] args){
        Properties props = new Properties();
        DispatcherDemo
    }
}
```

Run: zookeeper-start x 0-kafka-server-start x 1-kafka-server-start x 2-kafka-server-start x DispatcherDemo x

Process finished with exit code 0

The screenshot shows the same IDE as the first image, but with the 'consumer-start.cmd' file selected in the project structure. The code in the editor is the same as the first image. The logs in the bottom pane show the execution of the 'consumer-start' command, with timestamps and log levels (INFO) indicating the start and end of the process. The logs show a list of data points, each consisting of a date, time, and a series of numbers.

```
import org.apache.logging.log4j.LogManager;
import org.apache.logging.log4j.Logger;

import java.io.*;
import java.util.*;

public class DispatcherDemo {
    private static final Logger logger = LogManager.getLogger();
    public static void main(String[] args){
        Properties props = new Properties();
        DispatcherDemo
    }
}
```

Run: zookeeper-start x 0-kafka-server-start x 1-kafka-server-start x 2-kafka-server-start x consumer-start x

GOOG,2014-06-27,577.18,579.87,573.8,577.24,2230800.0,0.0,1.0,577.18,579.87,573.8,577.24,2230800.0

GOOG,2014-06-20,556.85,557.58,550.3915,556.36,4496000.0,0.0,1.0,556.85,557.58,550.3915,556.36,4496000.0

GOOG,2014-06-13,552.26,552.3,545.56,551.76,1217200.0,0.0,1.0,552.26,552.3,545.56,551.76,1217200.0

GOOG,2014-06-06,558.06,558.06,548.93,556.33,1732000.0,0.0,1.0,558.06,558.06,548.93,556.33,1732000.0

GOOG,2014-05-30,560.8,561.35,555.91,559.89,1766300.0,0.0,1.0,560.8,561.35,555.91,559.89,1766300.0

GOOG,2014-05-22,541.13,547.5999,540.78,545.06,1611400.0,0.0,1.0,541.13,547.5999,540.78,545.06,1611400.0

GOOG,2014-05-15,525.7,525.87,517.42,519.98,1699700.0,0.0,1.0,525.7,525.87,517.42,519.98,1699700.0

GOOG,2014-05-08,508.46,517.23,506.45,511.0,2015800.0,0.0,1.0,508.46,517.23,506.45,511.0,2015800.0

GOOG,2014-05-01,527.11,532.93,523.88,531.35,1900300.0,0.0,1.0,527.11,532.93,523.88,531.35,1900300.0

GOOG,2014-04-24,530.07,531.65,522.12,525.16,1878000.0,0.0,1.0,530.07,531.65,522.12,525.16,1878000.0

GOOG,2014-04-16,543.0,557.0,540.0,556.54,4879900.0,0.0,1.0,543.0,557.0,540.0,556.54,4879900.0

GOOG,2014-04-09,559.62,565.37,552.95,564.14,3321700.0,0.0,1.0,559.62,565.37,552.95,564.14,3321700.0

GOOG,2014-04-02,565.106,604.83,562.19,567.0,146700.0,0.0,1.0,565.106,604.83,562.19,567.0,146700.0

The step 3: Integrated Kafka cluster to Spark, was not achieved in this project, given time till 31st Jan, that can be done as well. 😊