# Improvising Processing of Huge Real Time Data Combining Kafka and Spark Streaming

M.Sc. Research Project

M.Sc. In Cloud Computing

Jeevantika Lingalwar

Student ID: 17156424

School of Computing

National College of Ireland

Supervisor:     Vikas Sahni

# National College of Ireland
## Project Submission Sheet
## School of Computing

| | |
|---|---|
| **Student Name:** | Jeevantika Sanjay Lingalwar |
| **Student ID:** | 17156424 |
| **Programme:** | M.Sc. in Cloud Computing |
| **Year:** | 2018-19 |
| **Module:** | MSc Research Project |
| **Supervisor:** | Vikas Sahni |
| **Submission Due Date:** | 18th April 2019 |
| **Project Title:** | Improvising Processing of Huge real Time Data combining Apache Kafka and Spark Streaming |
| **Word Count:** | 8480 |
| **Page Count:** | 20 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

**ALL** internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| | |
|---|---|
| **Signature:** | |
| **Date:** | 17th April 2019 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | |

Assignments that are submitted to the Program Coordinator office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Improvising Processing of Huge Real Time Data Combining Apache Kafka and Spark Streaming

Jeevantika Lingalwar
17156424

**Abstract**

Cloud Computing era brings lots of new innovations in computing technologies which includes processing of data, storage, network access, internet security, web data etc. Knowledge Discovery plays a key role in inventions and abstractions of new knowledge. World Wide Web also called as WWW has now became a hub of these new inventions. Due to the increase in data WWW became overloaded because of which extracting of data and processing the same got complex. Data must be processed quickly, with the goal that a firm can respond to changing business conditions proactively. Stream processing is the ongoing processing of information simultaneously with its creation. It is a perfect platform for processing information streams such as log data, sensor information etc. There are number of alternatives additionally to do constant preparing over information like Spark, Kafka Stream, Storm etc. In this paper the proposed framework for processing huge real-time data is Apache Spark Streaming combined with Apache Kafka. The execution time is recorded after each experiment to get a clear view about the processing speed. Based on the experiments performed it has been concluded that Spark alone takes less time to processes huge amount of data whereas this framework i.e. Kafka combined with Spark execution time depends on the amount of dataset taken.

**Keywords**: Cloud Computing, MapReduce, Hadoop, Apache Kafka, Apache Spark, Java

# 1   Introduction

Cloud Computing provides a platform for sharing data related to Internet, Software Information and resources and to be provided as a mean of computing equipment. Parallel programing model is the key technology introduced by cloud computing. The cloud computing technology is a part of vitalization technology which is generally used in abstraction and realization of data. The advancement in the development of data mining is referred to as web data mining. web data mining is used to extract useful information from web data, this is a process model which is an incomplete large amount of data which is understandable. These data set's characteristics is novel and effective.

This era is a world of computing and computers. Everything here is dependent on new and intelligent technologies. The computing world is changing from parallel computing to distributed computing to grid and now it has over taken by cloud computing. As the growth in Internet is increasing, the data is also increasing rapidly. Every day and every second huge amount of data is generated which comes from millions of sources whose numbers is increasing constantly. For instance, playing out a purchase where it appears as though we're purchasing only a certain thing-might produce several solicitations that would send and create information. Saving a record in the cloud doesn't mean putting away it on one server, it implies duplicating it over numerous areas for adaption to non-critical failure and accessibility. Additionally, it implies putting away logs and point

by point data about each miniaturized scale venture of the procedure, to have the capacity to recuperate things on the off chance that they return out badly. We log huge amounts of information. We reserve things for quicker access. We imitate information and setup reinforcements.

We spare information for future investigation. All these real-time structures mean specialized prerequisites to structure an information handling framework: where and how the data is created, what is the recurrence of changes and updates in the data, how fast we need to react to the change (Cha and Wachowicz, 2019). A strong framework is produced with the help pf Hadoop and MapReduce for predictive analytics implementation. There is need to build a solution that will get information from an assortment of sources, perform explicit calculation and examination on information on the fly. The implementation sequence is divided into two parts: Data Ingestion and decoupling layer between wellsprings of information and goal of information, Data Processing and response (Georgios Gousio, 2018). This proposed framework in this paper is a solution to the above challenges and issues. Processing of data is done using Apache Kafka combining it with apache Spark Streaming. The result is analysed based on the execution time of each experiment. The motivation behind this research is to provide a new approach for processing huge real time data. We compare the execution time of our approach with the existing one which makes this topic worth to work on.

# 2  Related Work

Since a long time, numerous scientists are learning about the new ways and methods of extraction and handling of enormous measure of information. Different calculations and systems have been proposed to illuminate this issue. Hadoop MapReduce assumed a vital job for all the proposed methodologies. The following is a concise synopsis of some past writing papers which are identified with the theme. The information gives a concise finish of the work done by different analysts additionally a short synopsis of the different calculations and methodologies that are being utilized by them.

## 2.1  Developing a Real-time Data Analytics Framework for Twitter Streaming Data

Twitter being a social networking communication service with in excess of 300 million clients, producing a large amount of data consistently every day. The most critical and important trademark of Twitter is its capacity for managing the clients to tweet about every day events, occasions, emotions and circumstances, suppositions, or not withstanding something which is new progressively. Currently there are various distinctive work processes offering continuous information examination for Twitter, displaying general preparing over gushing information. The examination will endeavour to build up an expository system with the capacity of in-memory preparing to remove and investigate organized and unstructured Twitter information. The Researchers in this paper have proposed a system that will incorporated information ingestion, stream handling, and information perception parts with the Apache Kafka information framework that is utilized to perform information ingestion task.

Moreover, Apache Spark makes it conceivable to perform modern information preparing and AI calculations continuously. The researchers have directed a contextual analysis on tweets about the quake in Japan and the responses of individuals around the world with the investigation on the time and birthplace of the tweets. Analysing information continuously requires data ingestion and preparing of the surge of information before the information stock piling step (Cha and Wachowicz, 2019). A portion of the utilizations of the continuous information examination are reconnaissance, conditions, social insurance, business knowledge, marketing, representation, cybersecurity and social media. This investigation exhibits a continuous information examination structure for breaking down Twitter information. The fundamental distinction between this investigation and different explores is that the proposed structure cannot just play out the fundamental preparing undertakings, be that as it may, makes a framework for performing more refined and muddles examination on the gushing data.

There are various studies carried out on real-time data analysis directed in the field of continuous information investigation. Each one of them makes a few commitments to a classification in day to day life and employment diverse approaches. A portion of these zones of utilization have a high rate of affectability to respond to real information. The securities exchange is a substantial case

of zones, which dependably have had a substantial dependence upon quick and exact examination. Internet based life information then again is generally the kind of information in view of sentiments and emotions. Regardless of this, there are still loads of valuable shrouded data, which might be separated from this kind of information. Dissecting posts on locales, for example, Facebook and Twitter may demonstrate vert valuable for illustration ends and making forecasts about exercises that happen on explicit regions of the worked at specific occasions (A. Bifet, 2018).

The motivation behind this investigation is to build up a structure for investigating Twitter information continuously. This structure has a few qualities which recognize it from conventional information examination approaches. The primary thought here is that there is a requirement for strategies to investigate a huge number of tweets coming each second, in a short measure of time. Additionally, the system ought to be autonomous of imported information volume; this is critical on the grounds that the volume of tweets is developing at a discernible rate. Figure below demonstrates a schematic of an adaptable 330 stream handling. The idea here is to gather occasion streams by various hubs and let numerous handling hubs to break down information in parallel. In this way, the test here is the means by which to oversee spilling information and how to examine it over the groups. Information handling work process often associates figuring assets to mechanize a succession of errands by preparing huge volumes of information, where diverse assets are associated for computerizing diverse assignments. On account of spilling information preparing, an adaptable and conveyed stage is required for consolidating expansive volumes of notable what's more, gushing information in the meantime. The system comprises of three areas: I) data ingestion; ii) data processing and iii) data visualization. The information ingestion area, interfaces straightforwardly to Twitter gushing API and in an adaptable way import information to the system. The information preparing segment with the capacity of gushing preparing over group accesses dispersed imported information, dissects information in memory, and performs handling undertakings on information, lastly sends the outcomes to be observed.



Figure 1: scalable processing structure (Babak Yadranjiaghdam, 2017).

Finally, the researchers have concluded by proposing a structure for Realtime examination of Twitter information. This system is intended to gather, channel, and investigate floods of information and gives us an understanding to what is well known amid a time and condition. The system comprises of three primary advances; information ingestion, stream preparing, and information representation. Information ingestion is performed by Kafka, an incredible message handling framework to import tweets, and to circulate it based on Topics that it characterizes, and to make it accessible over customers' hubs to be prepared by investigative devices. Apache Spark is utilized to get to these customers straightforwardly and examine information by Spark Streaming. The contextual analysis in this examination means to demonstrate the quality what's more, the significance of constant information investigation on social media gushing data. For example, assessment investigation of the tweets from the nation over can be done in presidential discussions. Doing as such, each discussion and the responses of the web-based life clients to it, can be utilized to train calculations and in up and coming discussions the fabricated model can play out a precise spilling handling on the web.

## 2.2 Literature review Performance evaluation of cloud-based log file analysis with Apache Hadoop and Apache Spark

Cloud based log file are generated in many formats by number of devices and in number of formats. When the analysis is carried properly each analysis leads to a useful information regarding each aspects of every system. Cloud computing is the most usable and suitable for such analysis as it can manage huge and high production rate and the huge size and variety in log files. The production rate of these log files can reach to number of TeraBytes or Petabytes. Due to these huge number of datasets reliable database solutions cannot be used for the analysis which is replaced by virtual database which is combined with distributed parallel processing and are considered more appropriate. As per the research done it has been assuming that log analysis is a use case which is of big importance and hence many of the important challenges that are related to cloud-based log file analysis are the big data challenges.

These analysis of log files later come on with other different challenges. Many of these systems are heterogenous and distributed so the logs from various number of other components must be related to each other. In addition, there may be some missing data or duplicate data that can make these log analyses more complicated or sometimes impossible. Hadoop is a framework that is been used by cloud computing in many companies for log analysis. The original design of Hadoop was mainly done for batch processing which provided scalability and fault tolerance, but it was not that useful for fast performance. It helps in running 1000s of nodes in petabytes of data. Apache Hadoop manages huge amount of log by breaking the files into blocks and give them to the Hadoop cluster nodes. It then breaks those jobs into number of small tasks. Apache Hadoop follows a processing model which reads and writes data from the disk. This process affects the performance and makes it non workable and unsuitable for real-time applications. The next approach is using Apache Spark which is more effective and suitable in which is uses more efficient main memory and minimizes the transfer of data from disk and to disk. Spark also provides new set of high-level tools which are used SQL queries, machine learning and graph processing (Mavridis, 2019).

For the performance evaluation of cloud-based log analysis the researchers in this paper have used two cloud computational frameworks, Apache Hadoop and Apache spark. The components of Apache Spark ecosystem used are Spark Core: a general execution engine, Spark SQL: Spark module for structured data processing, Spark Streaming: Enables real time data processing, MLib; Scalable machine learning library, GraphX: Graph computation Engine (Zaharia et al., 2012). Spark's Fault tolerance and parallel data structure i.e. Resilient Distributed Dataset. This is automatically distributed by Spark into the cluster and paralle operations are performed on them. For this analysis the researchers have carried out number of experiments which will evaluate the performance of Hadoop and spark and to do so they have used IAAS to create a private cloud infrastructure and some realistic log analysis operations with real log files. For the experimental purpose virtualized computing resources of Okeanos is been used (Koukis et al., 2013). It is an open source IaaS, which has been developed by Greek technology Network (Koukis & Louridas, 2012). With the help of Okeanos research six virtual machines have been created which are connected to a private virtual local network. In this there is one master node and 5 slave nodes, where the master node carries 8 virtual cores, disk space of 40GB and memory of 8GB. The 5 slave nodes are not that powerful as each of them has 4 virtual cores having a disk space of 40GB and memory of 6GB.

In order to evaluate these log files and investigate with Apache spark and Apache Hadoop the researchers have developed realistic applications in all the frameworks. The applications are developed using Java language with the help of Eclipse IDE and Apache Maven. The log files are stored in HDFS which is accessible by both Apache Hadoop and Apache Spark. In this paper researchers have focused on three main performances that are indicators, execution time, and the number of resources utilized with scalability. They have also compared the power consumption and cost of spark and Hadoop which is based on resource utilization and time measurements in which it is executed. Finally, the researchers have concluded that the analysis of cloud-based log files have been studied with the help of Apache Spark and Apache Hadoop comparing the cost, execution time and

resource utilization with the memory usage. Several experiments have been conducted with number of different slave nodes, having different size of input file and the type of application which will test the best performance of spark as compared to Hadoop. The total execution time of Hadoop and Spark log analysis application has been recorded and the factors affecting the same. [Mavridis and Karatza, 2019]

## 2.3 Literature review on Simplified Data processing on Large clusters: MapReduce

MapReduce is a programming model that is utilized for preparing vast number of datasets which is into number of genuine undertakings. The procedure is conveyed regarding Map and Reduce two unique capacities likewise called as the mapper and reducer work and the runtime condition or the framework parallelizes into groups which is of extensive scales which handles disappointments in machine and between machine interchanges which will utilize different circles and system. MapReduce capacities and employments are executed in the bunches of Google each day which process in excess of 20 petabytes of information every day (Mitali Srivastava, 2017). Pretty much every information is huge and the preparing or the calculation should be dispersed to a huge number of machines to complete the activity in a productive measure of time (Gousios, 2018). Consequently, the issue of parallelizing the calculation and to deal with the disappointments with gigantic arrangement of complex code is required to tackle this issue. Thus, to conquer the previously mentioned issues analysts have planned another deliberation to diminish the multifaceted nature of preparing calculation employments yet it conceals the adaptation to non-critical failure and entangled parallelization subtleties.

The reflection approach depends on Map and Reduce work. the guide work is utilized for all sources of info bringing about a key-esteem sets and later these sets are exchanged to the lessen capacity to produce the yield result by assessing these key-esteem sets (Dean and Ghemawat, 2008). They thought utilizing client determined guide and diminish work alongside useful model permits performing parallelization of expansive handling of calculations and furthermore for adaptation to internal failure utilizing re-execution as the essential required mechanism.(Srivastava et al., 2017) To actualize this methodology utilize straightforward however incredible interface has been made which will open up programmed appropriation of extensive scale calculations with parallelization which will accomplish vast calculation bunches with superior of the equivalent. The fundamental programming model incorporates a lot of sources of info and later creates a yield of key-esteem sets. The Map work takes the info pair which is given by the client and creates set of other key-esteem pair. MapReduce work together procedure the qualities with the middle of the road key 1 and pass it to the diminish work for further preparing. lessen work likewise acknowledges a key-esteem pair which is given by client with middle of the road key 1 and sets some an incentive for a similar key. later these qualities are converged to get a little arrangement of significant worth which is in zero or one (Armbrust, M., 2015). In this paper the scientists have proposed different MapReduce execution techniques which expresses that the correct decision of the strategy is finished by the earth. They clarified this with the assistance of a precedent which says, a usage might be pertinent for a little memory machine which is shared, another for substantial processors or multiprocessors additionally with significantly bigger organized machines.

The most exact clients of MapReduce are a finished improve of the ordering arrangement of the creation which produces the information structure which is utilized by Google web search tool. A contribution of expansive arrangements of reports is taken by information which is assembled by this slithering framework that is put away as GFS document. Utilizing the specially appointed conveyed MapReduce has ended up being helpful in each field. The code managing adaptation to internal failure, parallelization and appropriation are covered up inside the MapReduce library which makes the ordering basic and straightforward.

```
map(String key, String value):
    // key: document name
    // value: document contents
    for each word w in value:
        EmitIntermediate(w, "1");

reduce(String key, Iterator values):
    // key: a word
    // values: a list of counts
    int result = 0;
    for each v in values:
        result += ParseInt(v);
    Emit(AsString(result));
```
Figure 2: Examples of map and reduce function (Anubhav Tarar, 2018)

MapReduce libraries perform calm and in a legitimate manner so the reasonable calculations which are not related anyplace can be kept independently as opposed to blending them up keeping away from additional heap on information.

At long last, the analysts have finished up the system data transfer capacity is a constrained asset. Along these lines, tremendous number of advancements in their framework are modified to diminish the measure of information which must be sent everywhere throughout the system, this enables them to peruse information from the circles which are neighborhood and composing a duplicate to a similar nearby plate to spare system transfer speed. Next, they see that a similar execution can be actualized to facilitate the effect of the machine which are running moderate and to deal with disappointments of machine and loss of information.

# 3    Research Methodology

To achieve the result, four important process need to be performed before any implementation or development which uses Spark streaming and Kafka implementation and they are Transformation, Cleaning, writing the data received from various sources and Cleaning.

## 3.1   Objectives

The main objective of this approach is to keep up with the ever-increasing stream and amount of data that every data the developers are challenged with processing. From the network of data ingestion to processing real time data, this method will make it easier to keep up with the distributed streaming workflow. The main motive here is to reduce the complexity of processing huge real time data which makes it fast, reliable, scalable and pliable data workflow using two frameworks like Apache Kafka and Spark streaming. This system will be used to process real-time data from number of real time sources with some machine learning task and will efficiently experiment with real-time data processing in the real time world.

## 3.2   Proposed Algorithm workflow and framework

**Apache Spark Streaming Pipeline**
In the paper, some standard workflow has been implemented before processing the dataset which consists of 4 main stages which consists of Data Transformation, Data Cleaning, Data Validation and Data writing. For this purpose, various platform can be used like Java, Python and Spark Streaming for some of the workflow. In this paper python, Java and spark is used for the workflow.

   1. **Data Transformation:**
   The process is the first step and the main point for streaming an application. The data transformation can be easy, simple and sometimes complex depending on the required changes in the data that is processed between the source data and the target data. In the Spark Streaming the

data transformation is performed using a built-in function which is filter, map, RDDforeach and many more. Following is an iteration that have been used for the transformation in Spark streaming.

## 2. Data Cleaning:

Cleaning is the second most important stage before processing the data. You don't want any Null values to be disturbing the process for that purpose a clean data is very important. To remove the garbage data, a separate streaming application is executed for cleaning the data and for this purpose some custom built-in libraries in python is used. Following is the iteration used for data cleaning.

## 3. Data Validation:

Data validation is basically validation your data against a certain requirement whether the data produced matches the required fields, string, length etc. For data validation there is a separate spark streaming script written which is as below.

## 4. Data writing:

The final step is writing the data. After the data is processed through above method it is then passed to writing application which will write the data to the HBase for further processing of data. Kafka manages all the communication between all the applications using various topics and various applications.

The job streaming is done as below, and the pseudo code is presented for better understanding of the process (Databricks, 2019).

**Step 1**: **Initializing the Spark Context**

> **spark_context = SparkContext(appName="Application Transformation")**

This is the main step in spark which will allow to connect to spark cluster and create Resilient Distributed Datasets. The streaming application name is appName

**Step 2: Initialize Spark Streaming context**

**streaming_spark_context = StreamingContext(spark_context, 5)**

This is used to create Dstream from another various source. Here the values are taken from Kafka where value 5 indicates the batch interval. Spark streaming make sure that streaming data is divided not batches. Each batch is converted into RDD and these RDD are called as Dstream.

**Step 3: Get data from Kafka topics**

**direct_stream = KafkaUtils.createDirectStream(streaming_spark_context, [list of topics comma separated], {"metadata.broker.list": kafka_host})**

This is used to connect Kafka broker to fetch the messages from Kafka. Here no need of creating multiple input Kafka streams and integrating with DirectStream.

**Step 4: Transformation**

**direct_stream.map(custom_defined_functions)**

Here we have used map methods for the transformation. Note that you can used other transformation method as well as mentioned in this section i.e. filter, foreachRDD etc.

**Step 5: Pipeline Management**

streaming_spark_context.start()

streaming_spark_context.awaitTermination()

**Performance Tunning:**

To track and improve the performance of the pipeline the following configurations can be used.

1. spark.executor.memory: is use to define the memory that is being given per use per executor which shows up in the same format as JVM strings(memory strings).
2. Spark.executor.instances: This gives the number of executors also creates a conflict having spark.dynamicAllocation.enable.
3. Spark.streaming.reciever.maxRate: It gives the number(maximum) of records every sec a receiver will receive.
4. Spark.streaming.kafka.maxRatePerPartition: This gives the number(maximum) of records that will be read from each Kafka stream when using very new Kafka Directstream API. (Opcito Technologies, 2019)

**Apache Kafka**

Kafka is a publish-subscribe messaging system which is known for its real-time data streams ingestion and providing them to the downstream consumers which is managed in parallel and fault-tolerant specification. The Kafka data are stored in the topics which are split into partitions which are the sequence of records that are ordered and can be accessed as structured commit log. Producers writes these logs and the consumers read them as per their own convenience (Databricks, 2019).

**Kafka Topic**



Figure 6: Kafka Topic Partitions (Databricks, 2019)

Kafka subject is endless flow in which the records are taken back for a certain amount of time. The countless nature of this flow it is said that there is a need to first decide which data is to be studied and the place of time it is to be started, when starting a new query. As shown in the figure the three main choice of partition:

- Earliest: In this the older data is deleted which is outdated and the data is been read form the beginning.
- Latest: The data that has been arrived after starting the query is started is being processed.
- Per-partition assign: A unique offset is suggested to start from each partition, which permits fine-grained manipulate over where processing is supposed to start.

**Apache Kafka structured streaming:**

The structured streaming presents a unified batch and streaming API that allows us to see the posted data to Kafka which is as DataFrame. The data when processing in streaming fashion, the same API is used, and the identical information consistency is retrieved which is promised in batch processing.

**Kafka Topics records reading:**

First, the location of the Kafka cluster and the topic that is to be read is specified. Spark is the best framework which allows to read every single topic, set of topics, pattern of topic which in regex or any partitions that is a part of set of topics.

# 4 Design Specification

**Workflow: Processing Huge real Time Data**

The processing of huge real time data is done using Apache Kafka and Apache Spark and the execution time is recorded to compare the difference between both the framework's execution. The following are the steps used to process the huge amount of data.

For coding Apache NetBeans have been used which is an integrated development environment (IDE) for Java. To run and compile Linux applications on windows operating system Cygwin is used. Cygwin is an open source platform that allows Linux application to run on windows. For a user to use windows specific service on a different platform winutils is important. It is a collection of useful TcI commands which takes part of Win32 API. For processing the data, Apache Spark and Apache Kafka framework have been used. Apache Spark is a consolidated analytics engine for large-scale or huge amount of data processing. On the other hand, Apache Kafka is a distributed data streaming platform. Apache Kafka is used to process streams of records as they occur. It is also used for building real-time streaming applications that react to streams of data. In this project Apache Kafka will be used to process huge amount real -time data and compare with the processing time of Apache Spark.

The following are the steps involved for the processing of data:

| Steps No | Steps |
|----------|-------|
| Step 1 | For this project we need Windows 7 32-bit OS or more update version |
| Step 2 | The tools required for this project are Cygwin, Winutils, Apache Spark, Apache Kafka, NetBeans. |
| Step 3 | Download and Install NetBeans IDE 8.0.2 Java SE edition. |
| Step 4 | Download and Install Cygwin. |
| Step 5 | Download and Install Apache kafka and Apache Spark. |
| Step 6 | Download and Install Zookeeper-3.4.14. |
| Step 7 | Create a new Java project naming it ImprovisingProcessingFramework having three classes including one main class, MainFrame and Receiver. The MainFrame consists of Data Ingestion, Processing data using Spark and Kafka view. This is will give the complete view Loading the dataset and ingesting it into Spark and Kafka. |
| Step 8 | : Create a new project naming it ApacheSpark. This project will contain Java Source Package which will have three java classes one Main class, one Spark Frame, and Spark Receiver. The Dependencies folder will contain all the Spark and Java Dependencies required to process the data. The MainFrame consists of Datasets that is ingested from Processing Framework MainFrame which needs to be processed using Spark. |
| Step 9 | Create another new project naming it ApacheKafka. This project folder will contain Java Source Package which will contain three java classes, Main class, KafkaFrame class, KafkaReceiver Class. The Dependencies folder contains all the spark, Kafka |

| | |
|---|---|
| | and Java dependencies that is required for processing the data using Kafka. The Kafka MainFrame consists of dataset that is ingested from Processing Framework MainFrame. |
| **Step 10** | Click on the button Data Processing using Apache Spark in the Spark Main Frame. The process is started and when finished it will show a popup saying Data processing completed. Click on Ok. |
| **Step 11** | Click on the button Data Processing using Apache Kafka in the Kafka Main Frame. The process is started and when finished it will show a popup saying Data processing completed. Click on Ok. |
| **Step 12** | Go to the Improvising Processing Framework MainFrame and check the result of both Apache Spark and Apache Kafka where the processed data s recorded. |
| **Step 13** | Click on Comparison Tab where you can see the processing time comparison between apache Spark and apache kafka. |
| **Step 14** | Perform the above experiment 5 times and record the time of each execution to compare the exact difference |

Table 1: Steps for Processing Huge Real Time Data

**Flow Diagram**

The following is the flow diagram of the complete design specs.



Figure 10: Design Specification flow

# 5   Implementation

In this section, the implementation of the project has been discussed in detail. The processing of huge real-time data is done using Kafka Frame, Spark Frame and Processing Data main frame that will

ingest the dataset in Kafka and Spark frame for processing. The Frames are created using Java libraries. The detail implementation is explained as below,

## 5.1   Dataset – Market Basket Analysis

As per the definition, MapReduce is a programming model used for distributed programming based on Java. The Map function takes a set of data and converts it into another, and the reducer break these individual elements into smaller set tuple called as key/value pairs.



Figure 11: MapReduce Workflow (Talend Real-Time Open Source Data Integration Software, 2019)

Based on the above MapReduce workflow, Market basket analysis dataset is taken for this project. In each transaction many items are purchased by customers.
For Example,

**Transaction     Items**

t1        {T-shirt, Trousers, Belt}
t2        {T-shirt, Jacket}
t3        {Jacket, Gloves}
t4        {T-shirt, Trousers, Jacket}
t5        {T-shirt, Trousers, Sneakers, Jacket, Belt}
t6        {Trousers, Sneakers, Belt}
t7        {Trousers, Belt, Sneakers}

## 5.2   Project Execution

### Processing Real-time Data Framework

The Processing Data framework consists of three Java classes
1.   Main Class
2.   MainFrame Class
3.   MainFrameReceiver Class

The Main Class consists of function that is needed to call the processing data main frame class and the main frame receiver class that will receive the data which is used for processing the data. The main () method is a starting point for the Java application. The Main Frame is the UI for presenting the processing of real-time data. For reading and writing the file Java FileInputStream library is used. To build the UI java free UI Application frame library is used. To connect the application to any available port on local host DatagramSocket is used. The UI consist of loading the huge real-time data, ingesting the data to Apache Spark Frame and Apache Kafka Frame, comparing the executing time for Apache Spark an Apache Kafka. The MainFrameReceiver includes the function that will return the processed data in the Spark processed result frame and Kafka processed result frame which

is used for comparison of execution time. Below is the class diagram for this framework which gives a brief detail.



Figure 12: Class Diagram for Processing Data Framework

## Apache Spark framework

The Apache Spark framework consists of three Java classes
1. Main Class
2. SparkFrame Class
3. SparkReceiver Class

The Main class consists of function that is used to call the Spark Main Frame class and the Spark Main Frame Receiver class which will receive the data used for processing. The SparkFrame consists of the function that will contain the data that was ingested from the processing data main frame. The org.apache.spark.api.java.JavaRDD is a local handle class for Resilient Distributed Dataset. It is used to write function that is required to transform the information and the data to collect it locally. The org.apache.spark.api.java.JavaPairRDD is like JavaRDD which is also a local handle for a distributed collection of data [Dei.unipd.it, 2019]. Here maptoPair () method is used to count the frequency of each word and reduceByKey () method counts the repetitions of word in the text file. JavaSparkContext is Java friendly and a part of SparkContext which returns back JavaRDDs and implemented with collections in Java instead of the one in scala.

Below is the class diagram for this framework that gives a brief detail about it.



Figure 13: Class Diagram for Spark Framework

## Apache Kafka Framework

The Apache Kafka framework consists of three Java classes
1. Main Class

2. KafkaFrame Class
3. KafkaReceiver Class

The Main class consists of function that is used to call the Kafka Main Frame class and the Kafka Main Frame Receiver class which will receive the data used for processing. The kafkaFrame consists of the function that will contain the data that was ingested from the processing data main frame. JavaDStream is an interface to DStream, the essential abstraction in Apache Spark Streaming which speaks with the constant stream of data. DStreams can either be made from live information for e.g. Kafka, TCP sockets etc. The JavaPairDStream is an interface of DStream which consists of key-value pairs which generates more methods like join and reduceByKey. The Kafka Frame UI consists of two buttons, Processing data using Kafka and forwarding that data to result in the main frame. On click of button 1, a function is called that will process the data in Kafka and on the click of second button the processed data is sent to the main frame. Below is the class diagram for this framework that gives a brief detail about it.



Figure 14: class Diagram for Kafka Framework

## 5.3 Project Setup

**System Requirement**: Windows 64-bit with access to 32-bit application run

For the implementation purpose NetBeans IDE 8.0.2 and Java code deployment has been used.

The Market Basket dataset is taken from Kaggle which is used for processing using two frameworks and execution time is recorded for the both and compared in end Specifications and evaluations are explained in the section below. Figure shows datasets details that is used for processing in Kafka and Spark.

# 6 Evaluation

This section gives the details of the evaluation scenarios that is implemented and processed.

Figure 15: Dataset Market Basket Analysis

## Processing the Data using Apache Kafka

Apache Kafka being an open source streaming system is used for building real time application and streaming huge amount of data. To process the data, in the project Kafka topics are created. Kafka runs as a cluster on more than one server called as brokers. Kafka returns stream of records.


Figure 16: Kafka Topics

Each of these partitions are ordered and are sequence of record that is updated in the log. The Apache Kafka producers and Consumers play a very important role. Producer write the data in the topic. Consumers receives a message as soon as they subscribe to it. Consumer can be a part of a group or can be independent. Not only huge amount of data in ingested by Kafka but also works well with small data providing flexibility and scalability which makes is fast. For processing this huge real-time data there is a need for a framework of stream processing. There are number of streaming platforms like Kafka, Spark, Storm etc. In this paper Kafka and spark has been used. Kafka Stream is fast and lightweight that works well if the data ingestion is proper.


Figure 17: Kafka producers and Consumers (Lenadroid.github.io, 2019)

**Processing Data using Apache Spark**

Apache Spark being an open source framework for large scale computing environment. In this project, Apache spark is used for processing huge-real time data. Although spark is same as MapReduce, but it is more fast and powerful. For the implementation latest release of spark is used which supports continuous processing and execution as well as micro-batch processing. Apache Spark consists of workers which are the physical nodes. There is a driver that looks after all the workers. The execution of the tasks is managed by the executor where the tasks are the small unit of execution which are independent.



Figure 18: Spark working structure (Lenadroid.github.io, 2019)

As discussed above Apache Kafka is scalable when it comes to ingesting the streams of data produced by consumers. Whereas, Apache Spark is great for processing huge real-time or large amount of data. When combining both Apache Kafka and Apache Spark the processing of real-time streaming data gets easier and faster then before. There are number of experiments performed based on the above and following are the results generated for Kafka and Spark as per the execution time.

## 6.1 Experiment 1

The below figure shows the output for experiment 1 performed. The result generated shows the time comparison in seconds for the processing huge-real-time data using Kafka and spark. The Data is ingested via main frame into Spark frame and Kafka frame. The processed result is forwarded to the main frame result section and the comparison graph is generated as below.



Figure 19: Observation 1 Execution time apache spark and Apache Kafka

The above observation clearly specifies that the execution Time for Apache Kafka and Apache Spark differs.

**Execution Time Recorded for Kafka:** 6.501 sec
**Execution Time Recorded for Spark:** 8.109 sec

These observations were noted 5 times and final readings have been taken.

According to the above observation it can be said that processing huge real time data using Kafka and streaming platform takes less time than processing data using Spark alone.

## 6.2  Experiment 2

The below figure shows the output for experiment 2 performed.



Figure 20: Observation 2 Execution time apache spark and Apache Kafka

The above observation clearly specifies that the execution Time for Apache Kafka and Apache Spark differs.

**Execution Time Recorded for Kafka:** 6.111 sec
**Execution Time Recorded for Spark:** 4.960 sec

 These observations were noted 5 times and final readings have been taken.

According to the above observation it can be said that processing huge real time data using Kafka and streaming platform takes more time than processing data using Spark alone.

## 6.3  Experiment 3

The below figure shows the output for experiment 3 performed.



Figure 21: Observation 3 Execution time apache spark and Apache Kafka

The above observation clearly specifies that the execution Time for Apache Kafka and Apache Spark differs.

**Execution Time Recorded for Kafka:** 5.897 sec
**Execution Time Recorded for Spark:** 4.990 sec
These observations were noted 5 times and final readings have been taken.
According to the above observation it can be said that processing huge real time data using Kafka and streaming platform takes less time than processing data using Spark alone.

## 6.4 Experiment 4

The below figure shows the output for experiment 4 performed.



Figure 22: Observation 4 Execution time apache spark and Apache Kafka

The above observation clearly specifies that the execution Time for Apache Kafka and Apache Spark differs.
**Execution Time Recorded for Kafka:** 7.516 sec
**Execution Time Recorded for Spark:** 7.789 sec
These observations were noted 5 times and final readings have been taken.
According to the above observation it can be said that processing huge real time data using Kafka and streaming platform takes less time than processing data using Spark alone.

## 6.5 Experiment 5

The below figure shows the output for experiment 5 performed.



Figure 23: Observation 5 Execution time apache spark and Apache Kafka

The above observation clearly specifies that the execution Time for Apache Kafka and Apache Spark differs.

**Execution Time Recorded for Kafka:** 7.102 sec
**Execution Time Recorded for Spark:** 5.112 sec
These observations were noted 5 times and final readings have been taken.
According to the above observation it can be said that processing huge real time data using Kafka and streaming platform takes more time than processing data using Spark alone.

## 6.6  Discussion

There are 5 experiments performed with the same dataset and all observations having different execution times. As all observation's execution time differs from each other, the aggregate has been taken to get the exact result of the experiment. The result varies as per the performance. It has been noted Apache Spark take less time for processing real-time data whereas Kafka takes a bit long. But when combined both the execution time again varies depending on the amount of data. According to the analysis done and result generated it can be suggested that result various as per the amount of data taken. Apache Kafka is best worked in Kafka > Kafka context whereas Spark Streaming is best used in Kafka > Database kind of context.

# 7    Conclusion

As the demand of stream processing is increasing every day lots of new technologies and frameworks are emerging. The reason behind this increasing demand is the huge amount of data generating each day. The implementation of this project has highlighted some pros and cons of using Apache Kafka and Apache Spark for processing huge real-time data. The streaming process, the execution time, complexity level according to data set has been noted. From processing data using Spark and Kafka alone to combining both, the result has been discussed in detail in above section. Apache Spark is a general structure for huge scale information handling that bolsters bunces of various programming dialects and ideas for e.g. MapReduce, in-memory preparing, stream handling, chart handling or AI. Spark streaming gives and abnormal state reflection called DStream, which speaks to a nonstop stream of data. DStreams cab be made either from data streams from other sources like Kafka, Flume etc. The DStream is referred to as a grouping RDDs. Kafka Streams is a customer library for preparing and breaking down information put away in Kafka and either compose the subsequent information back to Kafka or send the last yield to an outer framework.
Kafka Streams straightforwardly addresses a ton of the troublesome issues in stream preparing: Occasion at once preparing (not micro batch) with millisecond idleness. Stateful preparing including conveyed joins and collections. A helpful DSL. Windowing without-of-request information utilizing a Dataflow-like model. Appropriated preparing and adaptation to internal failure with quick failover. No-vacation moving arrangements. According to the observations retrieved from the implementation it can be said that additional plugins can also be utilized in Kafka because the data that is once sent is retained and can be used many times. Apache Spark can be utilized with Kafka to stream the information yet if you are sending a Spark group for the sole reason for this new application, that is unquestionably a major intricacy hit. So, to reduce the complexity, we can utilize undeniable stream handling structure and after that Kafka streams comes into picture with the accompanying objective.

## 7.1  Future Work

For those who like to utilize cloud conditions for enormous information handling, this may intrigue. Azure currently has an option in contrast to running Kafka on HDInsight. You can leave your current Kafka applications as is and use Event Hubs as a backend through Kafka API. Occasion Hubs is an administration for gushing information on Azure, thoughtfully fundamentally the same as Kafka. As such, Event Hubs for Kafka biological systems gives a Kafka endpoint that can be utilized by your current Kafka based applications as an option in contrast to running your own Kafka bunch. It bolsters Apache Kafka 1.0 and more up to date customer forms, and works with existing Kafka applications, including Mirror Maker - you should simply change the association string and begin gushing occasions from your applications that utilization the Kafka convention into Event Hubs. Practically, obviously, Event Hubs and Kafka are two distinct things. Be that as it may, this element can be valuable on the off chance that you as of now have administrations written to work with Kafka, and

you'd like to not deal with any framework and attempt Event Hubs as a backend without changing your code.

## Acknowledgements

# References

URL- https://github.com/Jeevantika/ImproviseProcess

Koukis, E. and Louridas, P., 2012, December. ~ okeanos IaaS. In *EGI Community Forum 2012/EMI Second Technical Conference* (Vol. 162, p. 007). SISSA Medialab.

Armbrust, M., Xin, R.S., Lian, C., Huai, Y., Liu, D., Bradley, J.K., Meng, X., Kaftan, T., Franklin, M.J., Ghodsi, A. and Zaharia, M., 2015, May. Spark sql: Relational data processing in spark. In *Proceedings of the 2015 ACM SIGMOD international conference on management of data* (pp. 1383-1394). ACM.

Koukis, E. and Louridas, P., 2012, December. ~ okeanos IaaS. In *EGI Community Forum 2012/EMI Second Technical Conference* (Vol. 162, p. 007) SISSA Medialab.

Mavridis, I. and Karatza, H. (2019). Performance evaluation of cloud-based log file analysis with Apache Hadoop and Apache Spark.
Opcito Technologies. (2019). Building a real-time data pipeline using Spark Streaming and Kafka - Opcito Technologies. [online] Available at: https://www.opcito.com/blogs/building-a-real-time-data-pipeline-using-spark-streaming-and-kafka/ [Accessed 4 Apr. 2019].

Databricks. (2019). *Processing Data in Apache Kafka with Structured Streaming in Apache Spark 2.2.* [online] Available at: https://databricks.com/blog/2017/04/26/processing-data-in-apache-kafka-with-structured-streaming-in-apache-spark-2-2.html [Accessed 5 Apr. 2019].

Jacinto Arias, Jose A Gamez, and Jose M Puerta. Learning distributed discrete bayesian network classifiers under mapreduce with apache spark. Knowledge-Based Systems, 117:16–26, 2017.

Jeffrey Dean and Sanjay Ghemawat. Mapreduce: simplified data processing on large clusters. Communications of the ACM, 51(1):107–113, 2008.

Georgios Gousios. Big data software analytics with apache spark. In Proceedings of the 40th International Conference on Software Engineering: Companion Proceeedings, pages 542–543. ACM, 2018.

S. Cha and M. Wachowicz. Developing a real-time data analytics framework using Hadoop. 2017 IEEE International Congress on Big Data, pages 657–660, June 2017

A. Bifet, "Mining Big Data in real time," Informatica, 37(1), 2018, Pages 15 - 20.

Babak Yadranjiaghdam, Seyedfaraz Yasrobi, and Nasseh Tabrizi. In *Big Data (BigData Congress), 2017 IEEE International Congress on*, pages 329–336. IEEE, 2017.

URL- https://www.kaggle.com/xvivancos/market-basket-analysis

Talend Real-Time Open Source Data Integration Software. (2019). *MapReduce 101: What It Is & How to Get Started - Talend*. [online] Available at: https://www.talend.com/resources/what-is-mapreduce/ [Accessed 14 Apr. 2019].

Dei.unipd.it. (2019). *Most useful Spark methods — Big Data course documentation*. [online] Available at: http://www.dei.unipd.it/~capri/BDC/HOMEWORKS/api.html [Accessed 14 Apr. 2019].

Díaz, M., Martín, C. and Rubio, B. (2019). *State-of-the-art, challenges, and open issues in the integration of Internet of things and cloud computing*. (Díaz, Martín and Rubio, 2019)

Lenadroid.github.io. (2019). *Processing streams of data with Apache Kafka and Spark: ingestion, processing, reaction, examples*. [online] Available at: https://lenadroid.github.io/posts/distributed-data-streaming-action.html [Accessed 15 Apr. 2019].

dzone.com. (2019). *Spark Streaming vs. Kafka Streaming - DZone Big Data*. [online] Available at: https://dzone.com/articles/spark-streaming-vs-kafka-stream-1 [Accessed 15 Apr. 2019].

Ken Hess. Hadoop vs. spark: The new age of big data, Feb 2016. URL
.
K Jayamalini and M Ponnavaikko. Research on web data mining concepts, techniques and applications. In *Algorithms, Methodology, Models and Applications in Emerging Technologies (ICAMMAET), 2017 International Conference on*, pages 1–5. IEEE, 2017.

Kyong-Ha Lee, Yoon-Joon Lee, Hyunsik Choi, Yon Dohn Chung, and Bongki Moon. Parallel data processing with mapreduce: a survey. *AcM sIGMoD Record*, 40(4): 11–20, 2012.

Wang Lei and Liu Chong. Implementation and application of web data mining based on cloud computing. In *2015 International Conference on Intelligent Transportation, Big Data & Smart City (ICITBS)*, pages 470–473. IEEE, 2015.

Anand Singh Rajawat and Akhilesh R Upadhya. Big web data mining for predicting usage behaviour using fusion map reduce model. 2018.

R Shyam, Bharathi Ganesh HB, Sachin Kumar, Prabaharan Poornachandran, and KP Soman. Apache spark a big data analytics platform for smart grid. *Procedia Technology*, 21:171–178, 2015.

Mitali Srivastava, RAKHI Garg, and PK Mishra. A mapreduce-based parallel data cleaning algorithm in web usage mining. *International Journal of Computer Science and Applications*, 14(2), 2017.

Anubhav Tarar. How does spark use mapreduce? - dzone big data, Jan 2018. URL https://dzone.com/articles/how-does-spark-use-mapreduce.

Wenzheng Zhu and Changhoon Lee. A new approach to web data mining based on cloud computing. *Journal of Computing Science and Engineering*, 8(4):181–186, 2014.

# 8 Appendix - Configuration Manual

Configuration manual gives a detail knowledge about software used, how to install and Implement. Following sections give a brief information about executing the project and the observations. For evaluation, downloading and installing the software is needed.

# 1. Downloading Software

## 1.1 Java Development Kit

The Java SE Development Kit is a development environment that is used for building applications, applets, and other components which makes use of Java Programming Language. The Java Development Kit also includes some tools that are used for development and testing of programs which are written in Java Programming Language and that runs on the Java platform. The JDK can be downloaded from the link provided in references (*Java SE Development Kit 8 - Downloads*; n.d.) Choose the setup as per the configuration of the systems where the software to be installed.



Figure 1: Java SE Development Kit

## 1.2 Java SE Runtime Environment 8

When there is a need to run a Java program, but there is no development required for this purpose we need to download Java Runtime environment which is also Called as JRE. If there is a need to develop applications for Java, Java Development Kit needs to be downloaded also called as JDK. This JDK includes JRE due to which you there is no need to download them separately. The JRE can be downloaded with the link provided in the references.

## Java SE Runtime Environment 8u191

You must accept the Oracle Binary Code License Agreement for Java SE to download this software.

○ Accept License Agreement  ◉ Decline License Agreement

| Product / File Description | File Size | Download |
|---|---|---|
| Linux x86 | 68 MB | jre-8u191-linux-i586.rpm |
| Linux x86 | 83.72 MB | jre-8u191-linux-i586.tar.gz |
| Linux x64 | 64.85 MB | jre-8u191-linux-x64.rpm |
| Linux x64 | 80.67 MB | jre-8u191-linux-x64.tar.gz |
| Mac OS X x64 | 76.15 MB | jre-8u191-macosx-x64.dmg |
| Mac OS X x64 | 67.75 MB | jre-8u191-macosx-x64.tar.gz |
| Solaris SPARC 64-bit | 52.25 MB | jre-8u191-solaris-sparcv9.tar.gz |
| Solaris x64 | 50.14 MB | jre-8u191-solaris-x64.tar.gz |
| Windows x86 Online | 1.8 MB | jre-8u191-windows-i586-iftw.exe |
| Windows x86 Offline | 63.17 MB | jre-8u191-windows-i586.exe |
| Windows x86 | 66.52 MB | jre-8u191-windows-i586.tar.gz |
| Windows x64 | 71.16 MB | jre-8u191-windows-x64.exe |
| Windows x64 | 71.32 MB | jre-8u191-windows-x64.tar.gz |

Figure 2: Java SE Runtime Environment 8

### 1.3 7-ZIP Download

7-Zip is a file archiver having a high compression ratio. It is a free software which is open source. Most part o the code comes under GNU LGPL license, whereas the remaining code comes under the BSD clause 3. 7-Zip can be used on any computer including the commercial organization's computers as well. It's completely free. The link to download 7-Zip can be found in the references (*7-zip.org*).

## Download

**Download 7-Zip 19.00 (2019-02-21) for Windows:**

| Link | Type | Windows | Description |
|---|---|---|---|
| Download | .exe | 32-bit x86 | 7-Zip for 32-bit Windows |
| Download | .exe | 64-bit x64 | 7-Zip for 64-bit Windows x64 (Intel 64 or AMD64) |
| Download | .7z | x86 / x64 | 7-Zip Extra: standalone console version, 7z DLL, Plugin for Far Manager |
| Download | .7z | Any | 7-Zip Source code |
| Download | .7z | Any / x86 / x64 | LZMA SDK: (C, C++, C#, Java) |
| Download | .msi | 32-bit x86 | (alternative MSI installer) 7-Zip for 32-bit Windows |
| Download | .msi | 64-bit x64 | (alternative MSI installer) 7-Zip for 64-bit Windows x64 (Intel 64 or AMD64) |

Figure 3: Download 7-Zip

### 1.4   Download Apache Zookeeper New Releases.

Apache Zookeeper is used to develop open-source and highly maintained server which provides access to highly reliable distributed coordination. Zookeeper is a centralized medium which is used for maintaining the configuration information, providing distributed sync and group service. These services are used by distributed applications in other forms. The link to download apache Zookeeper is provided in the references (*http://zookeeper.apache.org/releases.html*).



Figure 4: Step 1: Apache Zookeeper Download



Figure 5: Step 2 Download newest release

## 1.5 Download Apache Kafka

Apache Kafka is a distributed streaming platform. It is used for building real-time streaming data pipelines that has data in between applications and systems. Kafka is also used for building real-time streaming applications which change itself to the streams of the data. The link to download Apache Kafka is provided in the references (http://kafka.apache.org/downloads).

Figure 6: Download Apache Kafka

## 1.6 Download Apache Spark

Apache Spark is an exceptionally fast analytics engine that is used for Machine Learining, AI and Big Data.Apache Spark is 100% open source, facilitated at the vendor autonomous Apache Software Foundation. Fast, adaptable, and developer friendly, Apache Spark is the main stage for vast scale SQL, data and stream processing, stream handling, and AI.



Figure 7: Download Apache Spark

## 1.7 Download Cygwin

Cygwin is a vast accumulation of GNU and Open Source apparatuses which give usefulness like a Linux appropriation on Windows. A DLL (cygwin1.dll) which gives generous POSIX API usefulness.

Figure 8: Download Cygwin

## 1.8 Download Apache NetBeans

Apache NetBeans is significantly more than a editor. It features source code semantically, lets you effectively refactor code, with a scope of convenient and incredible assets. Apache NetBeans gives editors, wizards, and formats to enable you to make applications in Java, PHP and numerous different languages.



Figure 9: Download Apache Netbeans.

## 2    Installation

In this section step by step installation procedures are explained for each software.

### 2.1  Java Development Kit Setup

1.  Start the JRE Installation setup and click the" Change destination folder" and install.
2.  Change the path of the directory without any spaces and click next.
3.  Next, open the system environment variables by clicking on control panel > system > Advanced System settings > Environment variables.
4.  Click New User Variable button in the section of user variables and type JAVA_HOME in the section of variable name and the path where you have your JRE installed in the variable value.

Figure 10:  Set Path for JDk- 1

5. Click OK. Now search for the variable in the section name "System Variable" in the section called "Environment Variables" and open it.
6. Edit the path and type the following text "; %JAVA_HOME%\bin" and the very end where already entered text ends as shown in the image below



Figure 11: Set Path for JDK- 2

7. Check if the installation has been done properly and to do so, open command prompt and type "java -version" if everything is installed correctly you will see a the current installed version of your java.

## 2.2    Apache Zookeeper Installation

1. Go to the directory where you have downloaded Zookeeper and go to its config.
2. Once the Config file is found rename your "zoo_sample.cfg" to "zoo.cfg".
3. Edit the zoo.cgf in notepad++ or any other text editor.
4. Search for dataDir=/tmp/zookeeper and replace it with: \zookeeper- (your version) \data.
5. Set the path in system environment variable similarly done for Java.
6. Add ZOOKEEPER_HOME = C:\zookeeper-3.4.7 to the system variable.
7. Run the zookeeper using the command prompt by typing zkserver.
8. If everything is installed properly you will see the output as the image below.

Figure 12: Zookeeper running on cmd.

## 2.3 Apache Spark Installation

1. Download Spark from https://spark.apache.org/downloads.html and choose "Pre-built for Apache Hadoop 2.7 or more later revision."
2. Unpack downloaded spark tgz file in a directory.
3. Download Hadoop 2.7's winutils.exe and put it in the location C:\Installations\Hadoop\bin
4. Now set HADOOP_HOME = C:\loaction to your Hadoop file
5. After Installation start windows shell ignore the warnings and move



Figure 13: Install Apache Spark

## 2.4 Install Apache NetBeans

1. Click on .exe installer file and run
2. If you downloaded the **All** or **Java EE** bundle, follow the following steps. Click Customize, make your choice and click ok.
3. Click Next
4. Accept the License
5. Chose the destination to install your IDE or select the default. Chose default JDK installation or chose your own.
6. Update your JDK to the newer version.
7. Click begin to start the installation
8. Click Finish

Figure 14: Apache NetBeans Install

## 2.5 GitHub Configuration

GitHub is an online hosting service for adaptation control utilizing Git. It is for the most part utilized for computational code. It offers most of the Distributed version control and source code of Git including its own characteristics.



Figure 15: Creating new Git Repository

## 3 Implementation

This section describes about implementation, importing procedures for any project and executing the codes to get the expected output. Figures below shows the clear explanation of creating, implementation and executing a project with their respective outputs.

## 3.1 Transforming Data



Figure 16: Spark Streaming script for Data Transformation

## 3.2 Data Cleaning



Figure 17: Spark Streaming python script for cleaning Data.

## 3.3 Data Validation



Figure 18: Spark Streaming script for data validation

## 3.4 Kafka Structured Streaming

1. pseudo code java for reading data from Kafka (single topic)

Figure 19: Subscribe to 1 topic

2. Pseudo code in java for multiple topics



Figure 20: Subscribe to multiple topics

3. Pseudo code in Java for a pattern



Figure 21: Subscribe to a pattern

## 3.5 Creating a New Java Project in NetBeans



Figure 22: Create New Project



Figure 23: Required Main classes

## 3.6 Executing the created project Apache ImprovingProcessing Framework



Figure 24: Improvising Processing Framework



Figure 25: Executing Port

## 3.7 Executing Apache Spark Framework



Figure 26: Apache Spark Frame

## 3.8 Executing Apache Kafka Framework



Figure 27: Apache Kafka Frame

## 3.9 Loading Data in Main Frame



Figure 28: Data Loaded

## 3.10  Data Ingestion



Figure 29: Data Ingested to Spark

Figure 30: Spark Frame with data



Figure 31: Data Ingested to Apache Kafka



Figure 32: Spark Frame with data

## 3.11     Data Processing Using Spark



Figure 33: Processing Data- Spark



Figure 34: Successfully processed Data- Spark



Figure 35: Processed Data in Spark

Figure 36: Forwarding Processed Data to MainFrame- Spark


Figure 37: Spark Processed Data in Main Frame

### 3.12 Data Processing Using Apache Kafka


Figure 38: Processing Data- Kafka

Figure 39: Successfully processed Data- Kafka



Figure 40: Processed Data in Spark



Figure 41: Forwarding Processed Data to MainFrame- Kafka

Figure 42: Spark Processed Data in Main Frame

## 3.13    Comparing the Execution Time



Figure 43: Comparing Execution Time



Figure 44: Comparison Graph

38

## 3.14    Dataset

Market basket analysis might tell a retailer that customers often purchase shampoo and conditioner together, so putting both items on promotion at the same time would not create a significant increase in revenue, while a promotion involving just one of the items would likely drive sales of the other. Market basket analysis may provide the retailer with information to understand the purchase behavior of a buyer.
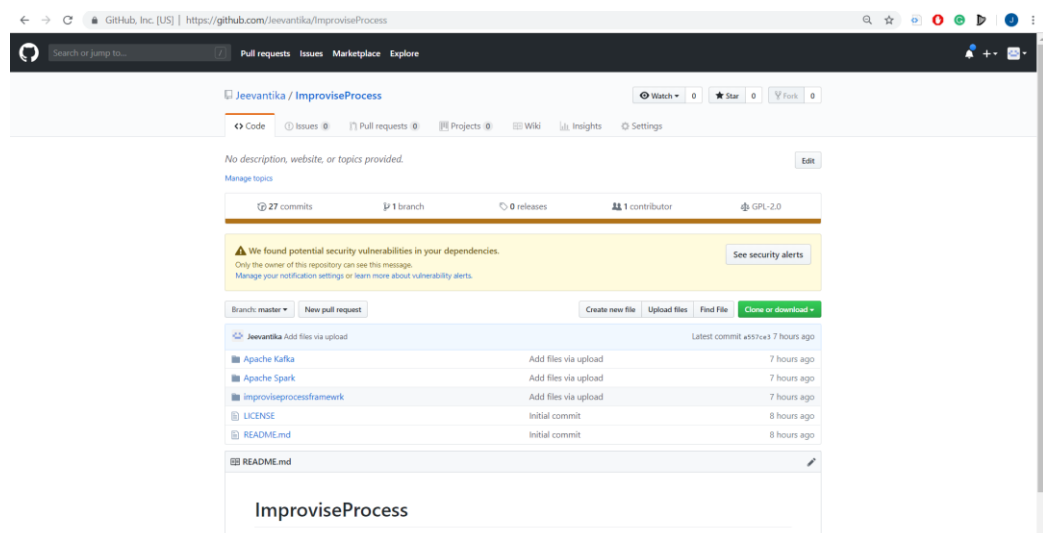


Figure 45: Dataset

## 3.15    GitHub



Figure 46: Project on GitHub

The Project is uploaded on GitHub
*URL-https://github.com/Jeevantika/ImproviseProcess*