# Real-time Data Streaming using Apache Spark on Fully Configured Hadoop Cluster

**[1]Kashi Sai Prasad, [2]Dr. S Pasupathy**

[1]Department of Computer Science and Engineering, MLR Institute of Technology, Hyderabad, 500043.

[2]Associate Professor, Department of Computer Science and Engineering, Annamalai University, Annamalainagar - 608002, Tamil Nadu.

## Abstract

*Data plays a major role in today's Internet world.Analyzing historical data became easy due to advancement of analytical tools. Gathering data from social networking websites is a great challenge for today's data scientists. Many advancements and research has been conducted to gather streaming data(data generated every second) .Hadoop has provided a component called Apache Flume to ingest data into HDFS for processing using MapReduce. It has its own benefits,which made many analysis easy for social networking data,but Apache Flume requires a depthknowledge on configuration files and administration.*

*Our work proposes a framework for real-time data streaming of Twitter data. Apache spark which is an enhancement of Hadoop in terms of speed and faster processing provides much more insight than Apache flume.Spark is an in-memory distributed computing engine to increase processing speed over MapReduce, Spark is considered one of the most advanced ecosystem component for Batch and near-real time processing. We in our paper are explaining in detail about data ingestion using Apache Spark and Scala IDE. In our work the data will be directly ingested from Twitter website through tokens and access keys provided,which will be explained in chapter 3,4. Our GUI can also help a user to tweet into Twitter directly without moving on to Twitter website. We have also provided an option to categorize tweet of specific persons using '#' tags.The data thus obtained can be used for statistical analysis and generating reports.*

**Keywords***: Apache Spark;BigData; Flume; Hadoop; MapReduce; Twitter data ingestion.*

## I. Introduction

Data ingestion and processing in parallel is very important because the amount of data is generally inexabytes and above. Currently there are different workflows offering real-time data analysis for Twitter, presenting general processing over streaming data. This study will attempt to develop an analytical framework with the ability of in-memory processing to extract and analyze structured and unstructured Twitter data [III].

Data processing can be classified into two types i.e.,  Batch Processing and Real-time processing. Batch processing dealing with huge amount of data which requires more time for processing.   In Real-time processing is instant, results are generated hardly in milliseconds or seconds.

### Data Processing Types

A. Batch Processing
B. Real-time Processing

### A. Batch Processing

- Core Batch   Hadoop, MapReduce
- Batch Interactive Hive,Impala [VIII]

### B. Real-time Processing

- Near Real-time                Spark Streaming
- Core Real-time                Storm, flink

Apache Spark with many advanced features became the super fast mechanism for data processing. Twitter provides a tweet area with 140 character restricted short messages, where users can express their views and share their opinions, feelings, etc., Companies which capture click stream data of the users like Netflix, YouTube, Amazon etc., gives the scope to gather huge amount of streaming data which can be processed using Apache spark. In twitter   hash tag is the convention of prefixing a word in the tweet with symbol '#' which indicates the keywords of a tweet. These words will be used for  categorization of tweets based on topics and aids in searching. As per the Twitter Inc., statistics more than 500M tweets are generated every day.   So handling such big streaming data using Apache spark is considered as next generation of Bigdata processing.

### Problem in detail:

The streaming data generated through regular activities of users is voluminous to measure, which needs a very huge storage. Continuous triggering of data through an easy storage process can be achieved with a fully distributed cluster which in parallel can store and process the data. Twitter data ingestion through Apache Flume in Hadoop is a very big process. Apache streaming using Scala programming language can reduce the burden of writing the Flume queries through a

simple process. The problem to be addressed is how Apache spark can be a solution for overcoming the process of Apache flume for data ingestion.

**Streaming Twitter Data Using Apache Flume**

Prerequisites: Twitter Account

I.     Login to the twitter Account

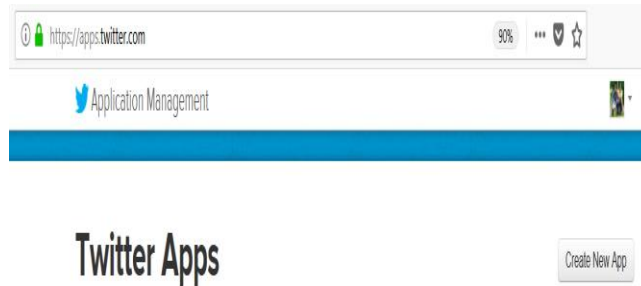II.    Go to the following link and create new App
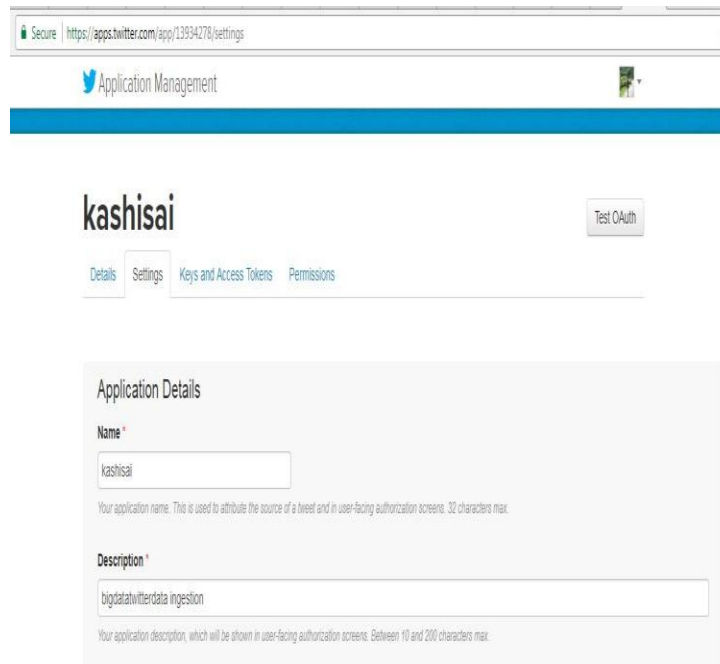


Figure 1

III.   Enter the necessary details.



Figure 2

IV.    Next we should accept the developer agreement and select a button: 'create your Twitter application'.

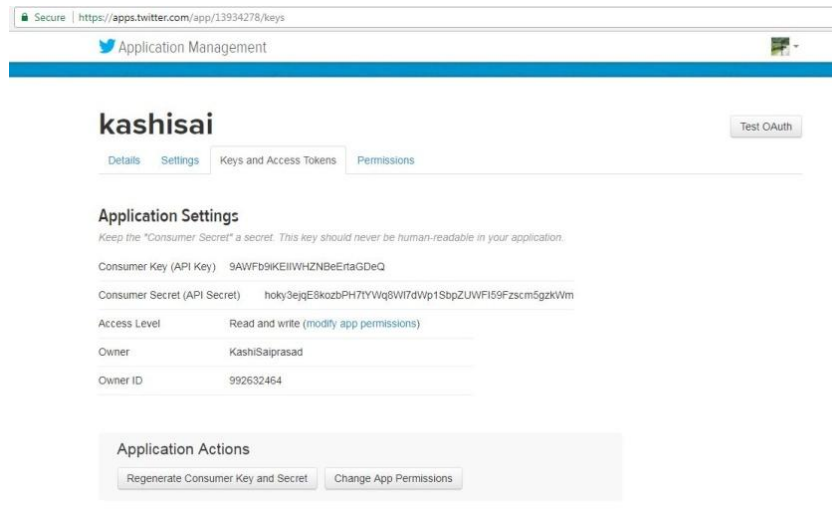V.   Select the "Keys and Access Token" tab [VII].



Figure 3

VI.   Copy the consumer key and the consumer secret code and Copy the Access Token and Access token Secret code.

VII.   We have to decide which Keywords tweet data to be collected from the twitter application. So, you can change the keywords in the TwitterAgent,sources.Twitter.keywordscommand.



Figure 4

VIII.     Create a new directory inside HDFS path, where the Twitter tweet data should be stored. hadoop fs –mkdir –p /usr/flume/tweet

  IX.     For fetching data from Twitter, Use the below command to fetch the twitter tweet data into HDFS cluster path.

Flume-ng agent –n TwitterAgent –f <location of created /edited conf file>

   X.     The above command will start fetching data from Twitter and streams in into the HDFS given path.

  XI.     Once, the tweet data started streaming it into the given HDS path we can use 'Ctrl+c' command to stop the streaming process.

As Flume has Complex structure, configuration and maintenance will be difficult. In Flume throughput depends on the backing store of the channel so scalability and reliability in not up to the mark.

**Fully Configured Hadoop Cluster**

We have access to a fully configured Hadoop cluster with 9 Nodes (3 Master nodes, 5 data nodes and an edge node). The data ingested from Twitter directly comes into HDFS which can be accessed by an user by providing the credentials.



Fig 5: Multi Node Hadoop Cluster Architecture

## II.     Literature Survey

Hadoop though provides the best solution for the Bigdata requirements, there are much more requirements due to vast usage of networking sites. Apache Spark was built to prevail over the restrictions of Hadoop MapReduce cluster computing paradigm. Apache spark is built using the existing MapReduce paradigm but the computing power and processing strategy is completely different [IV].As the computing power of Spark is very fast it is said that spark is hundred times

faster than Hadoop MapReduce because of the in memory processing of data the time spent in moving data in and out is not required [V].

As per Microsoft statistics by 2020 there will be 30B connected devices, each individual will generate 1.5GB of data per day, Smart homes will generate approximately 50 GB of data per day, Smart buildings with smart devices can generate 150GB data per day.

The limitations which are identified in terms of data processing by MapReduce are below:

1. MapReduce uses only Java for building an application.

2. Hadoop MapReduce uses disk-based processing and not suitable for performing streaming data.

3. Not suitable for Real-Time streaming data processing

4. High latency

5. Reading and writing data from disk (No in memory processing)

6. Reusability

7. Fault Tolerance

8. Complex to use

9. MapReduce needs an external job scheduler like Oozie to schedule complex jobs.

10. MapReduce can't cache the data in memory for future requirements.

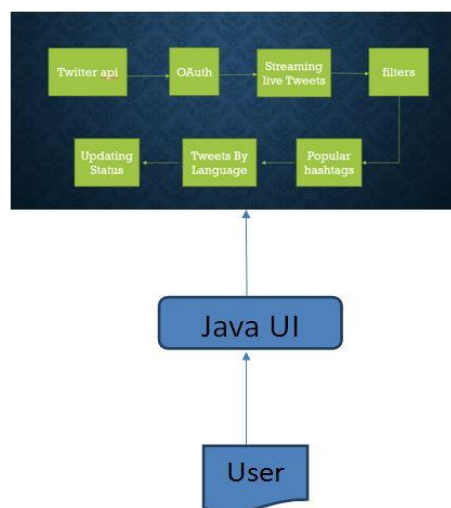## III. System Architecture
**Dataflow Diagram**



Figure 6

Apache spark is a general purpose engine to solve variety of data processing problems and also it can be used for traditional batch processing operations with additional support for Real-time

Processing. The primary purpose of Spark is to efficiently handle data of a large magnitude [II]. It's a fast In-Memory Data Processing engine. There are specific modules for several use cases like Spark Core, Spark SQL, Spark Streaming module for real-time processing[V].

**Workflow:**

An user uses the UI developed using Java and in the background we have used the Twitter API to connect to Twitter account and using OAuth standard our application provides a secure access to twitter data.

OAuth is a standard that application can use to provide client application with secured delegate access. OAuth can help the developer to create and debug applications which use Twitter's Application program Interface.If we want to use OAuth our application need to obtain the access tokens of an user. OAuth is a protocol that supports authorization workflows.

The streaming live tweets can be ingested into HDFS in a single click until we close the UI. Later filters are applied to get the Tweets using hash (#) tags, Tweets are also obtained from various languages using filters. The data thus obtained will be stored in the Hadoop cluster and using Scala we can process the data for insights.

## IV. Implementation

### A) Configuration

Adding Scala nature to already existing Java project:



Figure 7

Setting Scala verion:

Figure 8

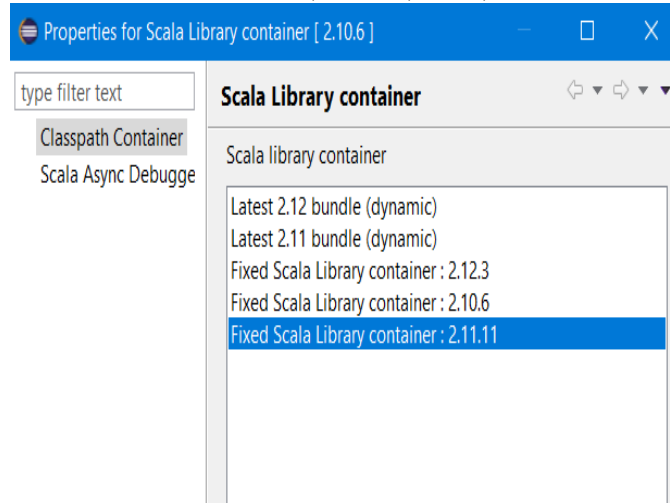Scala runs primarily on the Java Virtual Machine (JVM). As Scala and the JVM improve independently over time, Scala may drop compatibility with older JVM versions, in order to better take advantage of new JVM features.

**B) Configuring build path for adding spark jars:**

Download spark jars from
http://spark.apache.org/downloads.html
Add all jars from \spark-2.2.1-bin-hadoop2.7\jars to project build path.



Figure 9

The developed project consists of both Java and Scala IDE and we need to choose compatible JVM and Scala versions.

## V. Results

## A) Screenshots



Figure 10

```java
private void tweetsActionPerformed(java.awt.event.ActionEvent evt) {
    // TODO add your handling code here:
    TweetsConsole.main(new String[] {});
}

private void hashTagsActionPerformed(java.awt.event.ActionEvent evt) {
    // TODO add your handling code here:
    HashTagConsole.main(new String[] {});
}

private void statsActionPerformed(java.awt.event.ActionEvent evt) {
    // TODO add your handling code here:
    Twitterstats.main(new String[] {});
}

private void postTweetActionPerformed(java.awt.event.ActionEvent evt) {
    // TODO add your handling code here:
    PostTweet.main(new String[] {});
}
```

Figure 11

Streaming Tweets

Retreiving Tweets

```
-----------------------------------------
RT @ImRaina: This is a good habit to wake up to every morning! Another?Congratulations #SathishSivalingam  on our third #GC2018Weightl
RT @soompi: #EXO_CBX Announces Special Daily Live Broadcasts For "Blooming Days" Comeback
https://t.co/9ceSXRmGU0 https://t.co/ljjArgbMt8S
RT @ZachsCubClinic: Today in one of our med classes we talked about LGBT patients and the professor isn't taking anyone's shit

Student: Wh†
RT @angelicola: ugh why do i keep talking to this guy ? he¿s so fucking annoying but i feel bad for him ugh
RT @pennystockalpha: YOU HAVE NEVER EXPERIENCED A RUNNER UNTIL YOU ARE A PART OF WHATS COMING YOUR WAY !!!

CRYPTOPENNYSTOCK@GMAIL.COM

#Pe†
@Sith9890 @TheQuantumCat Its not that hard on hard but extreme is aids
Britain sees the Commonwealth as its trading empire. It is sadly deluded https://t.co/Qfpj1tgct5 https://t.co/dHXZwjv0A1
Shoutout to Lauryn Hill
RT @aligatie: I think I¿m in love with @sabrinaclaudio https://t.co/WAv1rq8rxX
@ForgetAmnesia Why have I never thought of this? THIS so much
...


-----------------------------------------
Time: 1523079754000 ms
-----------------------------------------
RT @geoff032: @iman_marshall8 @tswag03 We have the best wide receiver core in the nation period...@tswag03 @MikePitt_Jr @TrevonSid
RT @Srithar_VJ07: Reached 485K Tweets !!
```
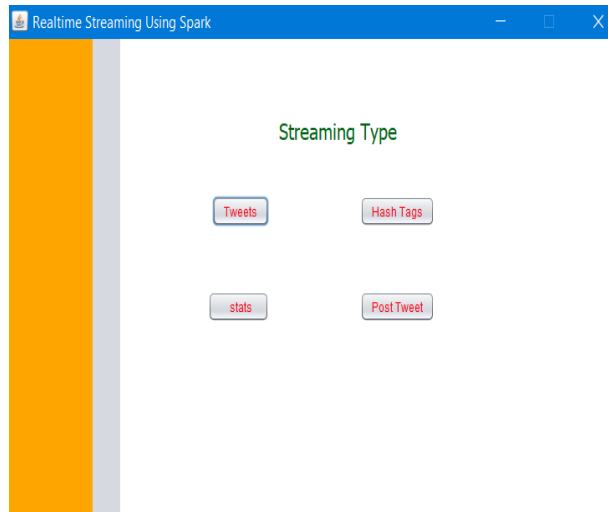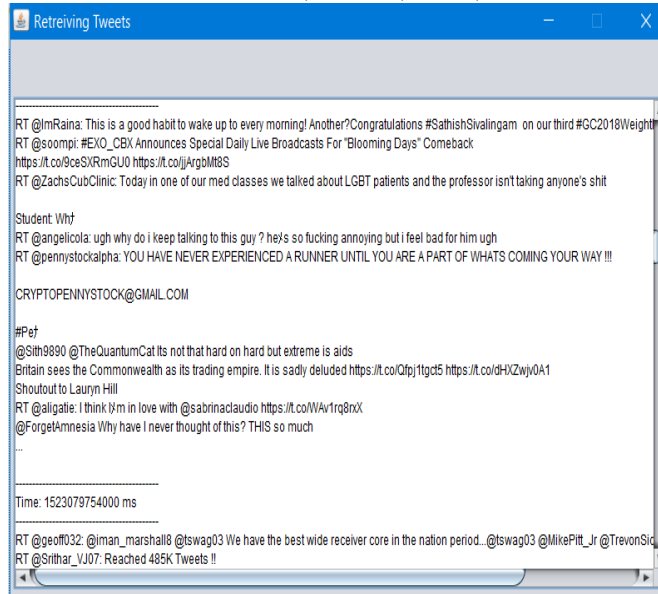
Figure 12

Here in this scenario we are redirecting the console output to a JTextArea as

```java
public class JTextAreaOutputStream extends OutputStream
{
    JTextArea ta;
    //caret.setUpdatePolicy(DefaultCaret.ALWAYS_UPDATE);

    public JTextAreaOutputStream(JTextArea t) {
        super();
        ta = t;
        ta.setCaretPosition(ta.getDocument().getLength());
        DefaultCaret caret = (DefaultCaret)ta.getCaret();
        caret.setUpdatePolicy(DefaultCaret.ALWAYS_UPDATE);
    }

    public void write(int i) {
        ta.append(Character.toString((char)i));
    }

    public void write(char[] buf, int off, int len) {
        String s = new String(buf, off, len);
        ta.append(s);
    }

}
```

Figure 13

System.setOut(new PrintStream(new JTextAreaOutputStream(console)));
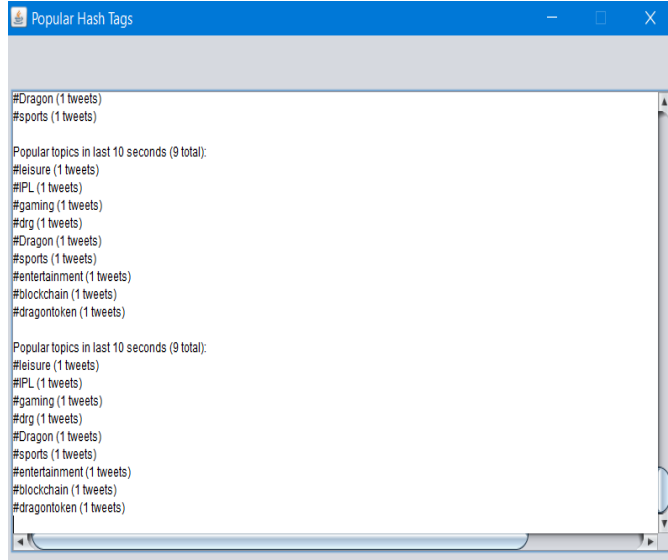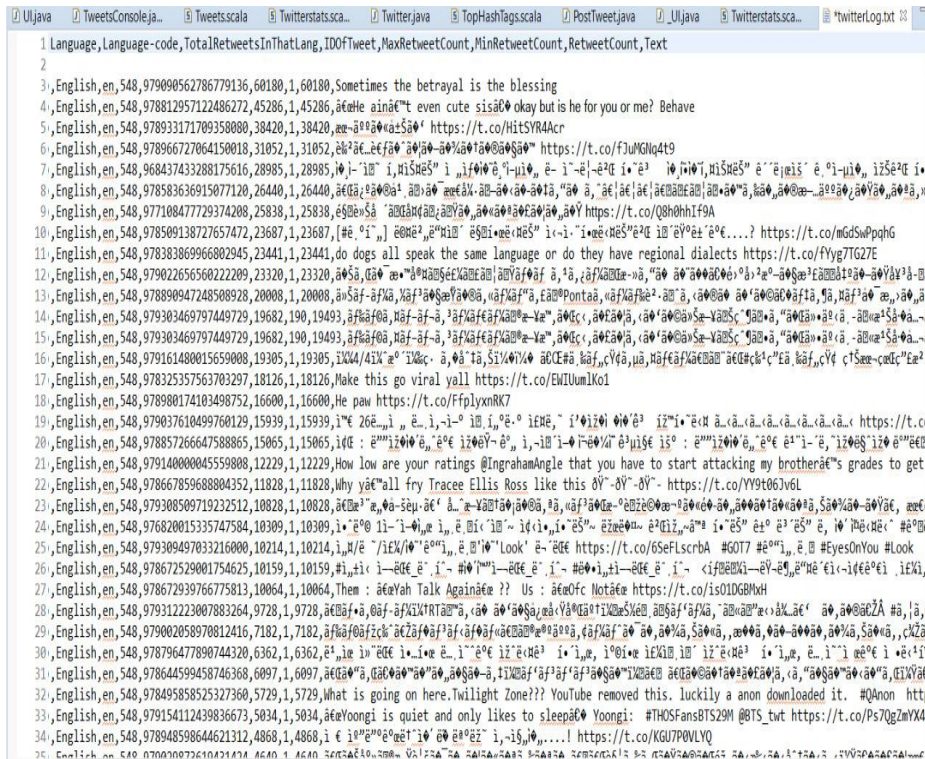
Popular HashTags

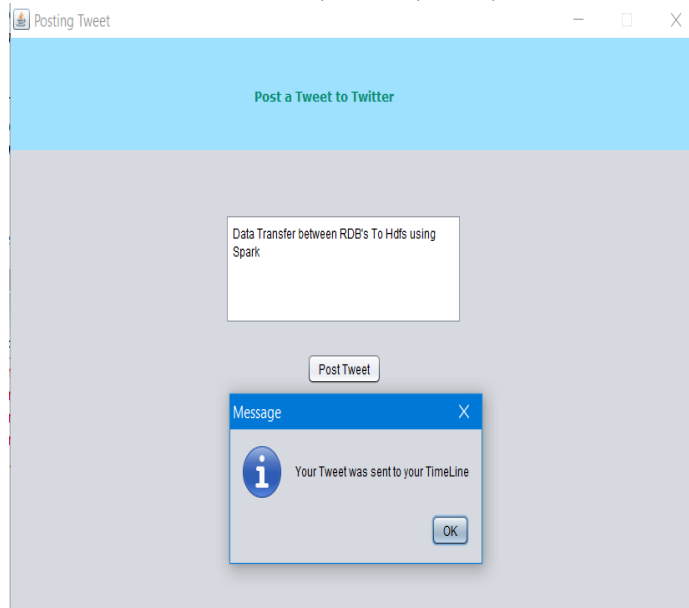Figure 14

Twitter stats



Figure 15

Figure 16

Checkthe Twitter Timeline so that the post will be seen on the account who have logged in.

## VI.    Conclusion

 Apache Hadoop became a synonym for Bigdata due to its scalability, cost effectiveness, replication etc., to process voluminous data. Even though Hadoop is very fast in processing and generating efficient results, there is still incompleteness. Data ingestion into Hadoop cluster is a crucial task in today's fast moving world.The components of Hadoop Ecosystem like Apache Sqoop, Apache Flume, KAFKA, Cassandra can complete the task but the streaming data generated needs a much more reliable and efficient tool like Apache Spark which is 10 times faster than Hadoop. Our work concludes explaining the procedure to connect to a Twitter using OAuth, Access token and ingest data into cluster, and process to find out insights from generated data.

Apache Spark is preferred on top of Apache flume for its varied applications.A Hadoop user can easily understand the working of Apache Spark and Scala in the explained work.

## Future Work

Technology is in continuous advancement due to large amount of data generated from online networking sites. Gathering data and storing it in Hadoop cluster is explained in our work. This work can be more precisely implemented using Apache Oozie which can gather data at regular time intervals allotted. As we have feasibility to increase the nodes in the developed cluster, as per data requirements we can parallely process the data.

Below are some other interesting dimensions which we can add to this project in the future.

- Creating dashboards to monitor all incoming tweets, user interaction and present results in a visually appealing manner.
- Geo-spatial maps could be plotted, signifying the location (and sentiment) of the tweets
- Creating a pipeline to dump all processed tweets into HDFS, Databases like Cassandra, Mongo DB, PostgreSQL or MySQL.

**References**

I.      Altti Ilari Maarala, Mika Rautiainen, Miikka Salmi, Susanna Pirttikangas and Jukka Riekki", Low latency analytics for streaming traffic data with Apache Spark" IEEE International Conference on Big Data (2015).

II.     Anand Gupta, Hardeo Kumar Thakur " A Big Data Analysis Framework Using Apache Spark and Deep Learning", IEEE International Conference on Data Mining Workshops (2017).

III.    Babak Yadranjiaghdam, Seyedfaraz Yasrobi, Nasseh Tabrizi "Developing a Real-time Data Analytics Framework For Twitter Streaming Data", IEEE 6th International Congress on Big Data (2017).

IV.     Hassan Nazeer, Waheed Iqbal, Fawaz Bokhari, Shuja Ur Rehman Baig " Real-time Text Analytics Pipeline Using Open-source Big Data Tools", arXiv:1712.04344, Dec (2017).

V.      Marouane Birjalia, Abderrahim Beni-Hssane, Mohammed Erritali "Analyzing Social Media through Big Data using InfoSphere BigInsights and Apache Flume ",  The 8th International Conference on Emerging Ubiquitous Systems and Pervasive Networks Elsevier (2017).

VI.     Ramkrushna C. Maheshwar, D. Haritha "Survey on High Performance Analytics of Bigdata with Apache Spark", International Conference on Advanced Communication Control and Computing Technologies (ICACCCT) (2016).

VII.    Sangeeta "Twitter Data Analysis Using FLUME & HIVE on Hadoop Framework", Special Issue on International Journal of Recent Advances in Engineering & Technology (IJRAET) V-4 I-2February (2016).

VIII.   S. Cha and M. Wachowicz. "Developing a real-time data analytics framework using Hadoop", IEEE International Congress on Big Data June (2015).