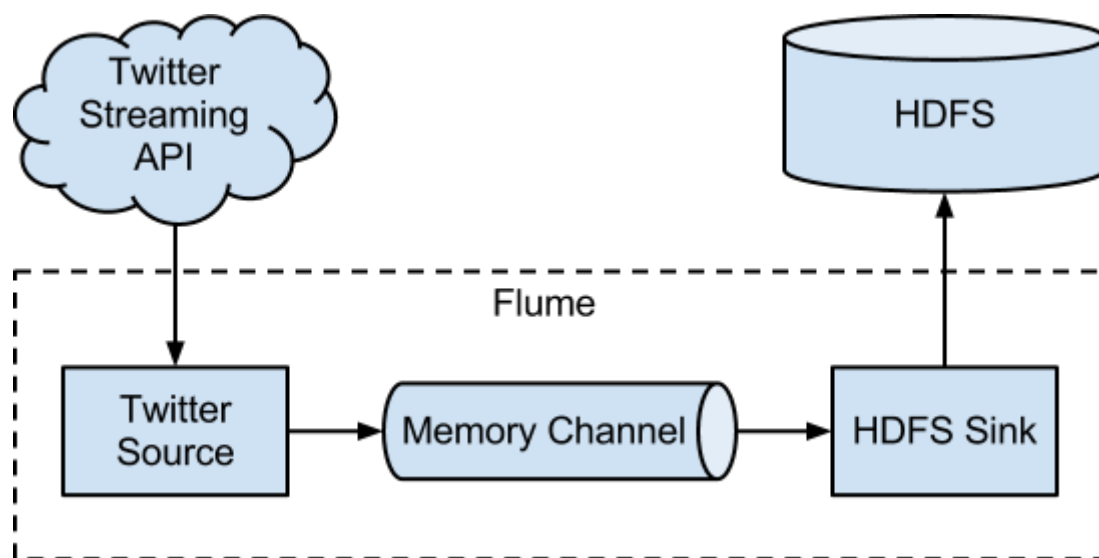# Documentation of Streaming of Twitter data using Apache Flume

-Prepared by Vignesh.R (15CSE107)

**Abstract:**

**Apache Flume** is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of log data. It has a simple and flexible architecture based on streaming data flows. It is robust and fault tolerant with tunable reliability mechanisms and many failover and recovery mechanisms. It uses a simple extensible data model that allows for online analytic application. This document depicts the streaming of Twitter data using Flume engine.



There are 3 major components, namely: Source, Channel, and Sink, which are involved in ingesting data, moving data and storing data, respectively.

Below is the breakdown of the parts applicable in this scenario:

- **Event** – A singular unit of data that is transported by Flume (typically a single log entry).
- **Source** – The entity through which data enters into the Flume. Sources either actively samples the data or passively waits for data to be delivered to them. A variety of sources such as log4j logs and syslogs, allows data to be collected.
- **Sink** – The unit that delivers the data to the destination. A variety of sinks allow data to be streamed to a range of destinations. Example: HDFS sink writes events to the HDFS.
- **Channel** – It is the connection between the Source and the Sink. The Source ingests Event into the Channel and the Sink drains the Channel.
- **Agent** – Any physical Java virtual machine running Flume. It is a collection of Sources, Sinks and Channels.
- **Client** – It produces and transmits the Event to the Source operating within the Agent

## Streaming Twitter data from Flume:-

**Step 1:**

Remove protobuf-java-2.4.1.jar and guava-10.1.1.jar from library directory of flume-ng. Super user permission is given for this modification. If no such directories are available, then ignore it.
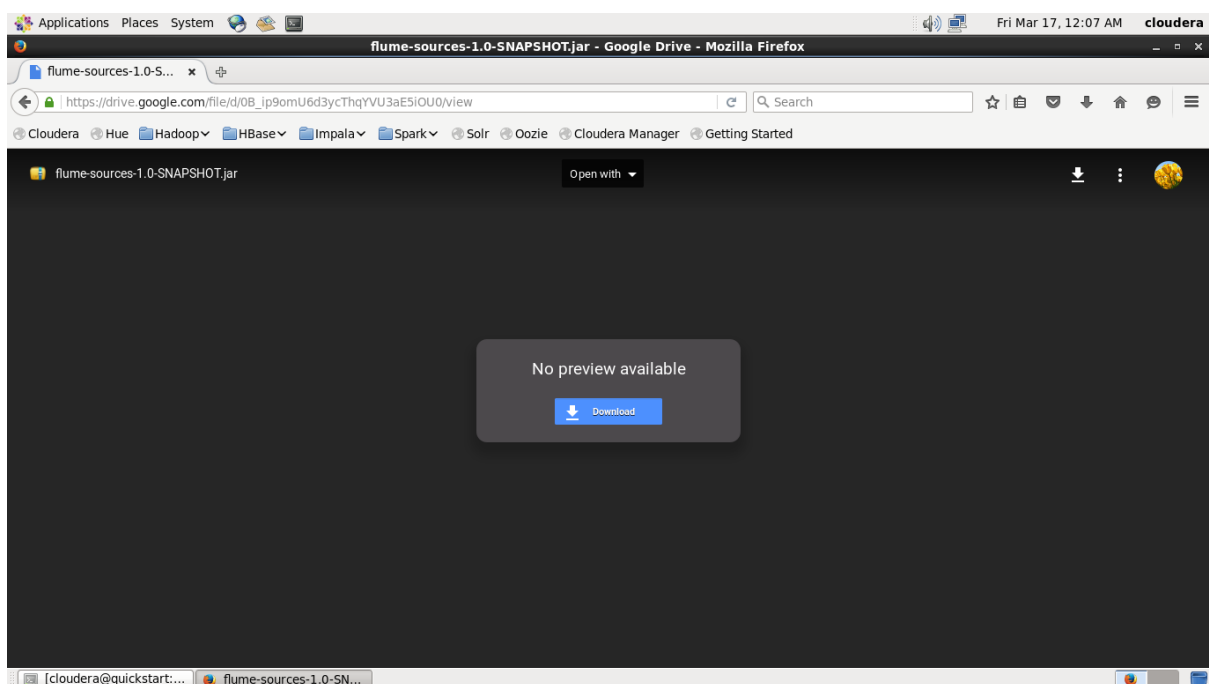
**Command:** sudo rm /usr/lib/flume-ng/lib/protobuf-java-2.4.1.jar /usr/lib/flume-ng/lib/guava-10.0.1.jar



**Step 2 :**

The flume-sources-1.0-SNAPSHOT.jar file is downloaded from https://drive.google.com/file/d/0B_ip9omU6d3ycThqYVU3aE5iOU0/view?usp=sharing link.

**Step 3:**

Click on Download option and then save the jar file.



**Step 4:**

Move the flume SNAPSHOT.jar file to the lib folder of flume from downloads.

**Command:** sudo mv /home/cloudera/Downloads/flume-sources-1.0-SNAPSHOT.jar/usr/lib/flume-ng/lib

**Step 5:**

Check whether SNAPSHOT.jar file is moved to the library folder of flume-ng.

**Command:** ls /usr/lib/flume-ng/lib/flume-sources-1.0-SNAPSHOT.jar

**Step 6:**

Then the present working directory is moved to the flume-ng directory.

**Command:** cd /usr/lib/flume-ng/



**Step 7:**

The flume-env.sh.template content is copied to flume-env.sh. Super user permission is given for this modification.

**Command:** sudo cp conf/flume-env.sh.template conf/flume-env.sh

**Step 8:**

Then the JAVA_HOME and the FLUME_CLASSPATH of the copied flume-env.sh file is edited using gedit command. Super user permission is given for this modification.
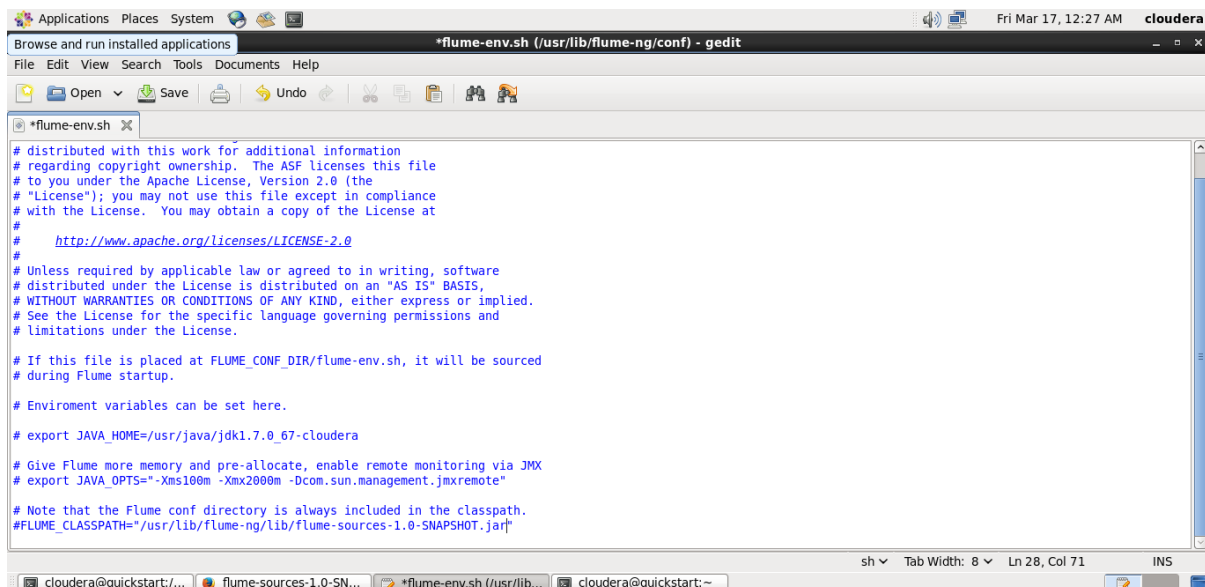
**Command:** sudo gedit conf/flume-env.sh

**Step 9:**

Here, the java path has to be identified from the local machine and java path is set. The class path of flume is set to the downloaded SNAPSHOT.jar file location. Then save the file and close it.

**Command:** JAVA_HOME= /usr/java/jdk1.7.0_67-cloudera

**Command:** FLUME_CLASSPATH="/usr/lib/flume-ng/lib/flume-sources-1.0-SNAPSHOT.jar"



**Step 10:**

Open the web browser and sign in to your account in twitter. Then enter your login credentials and open your account.

**Website link:** https://twitter.com/login

**Step 11:**

Then switch to the apps site from your twitter account.

**Website link:** https://apps.twitter.com

**Step 12:**

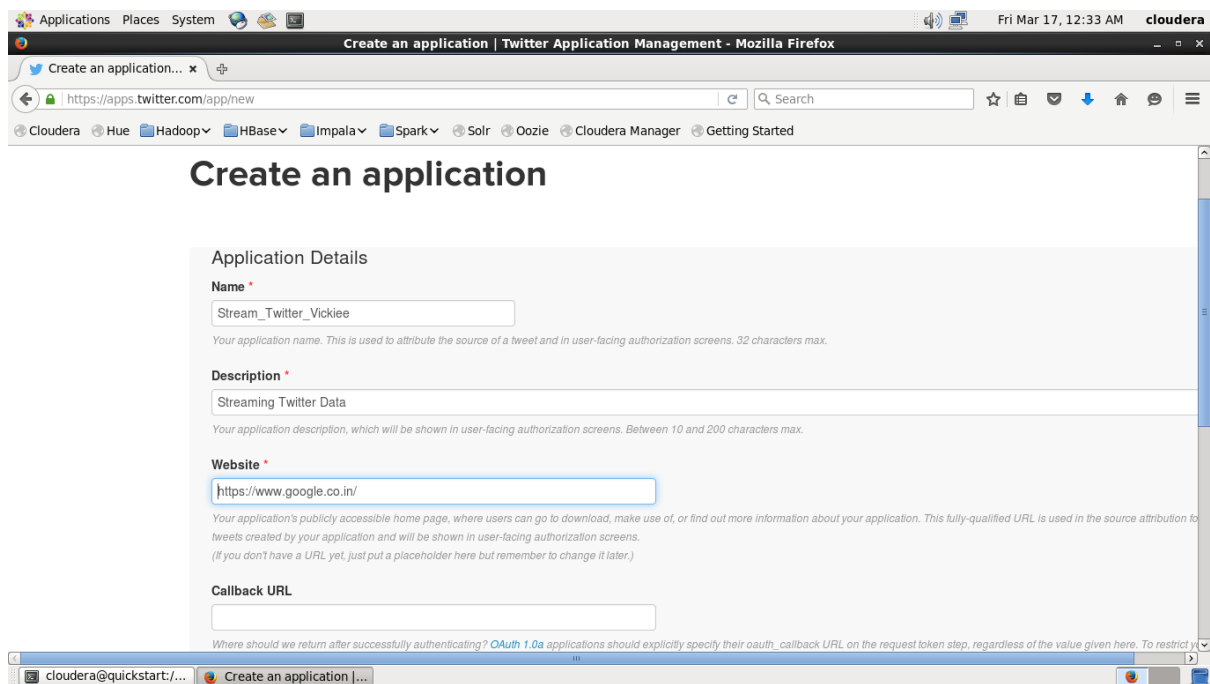Application window appears. Click on the Create Application option.

**Note:** To create an application, your phone number has to be added in your twitter account.



**Step 13:**

Enter the Application name, its Description and the website.

**Note:** Application name should be unique. Search engine site is given in the website field. Callback URL field is left blank.



## Step 14:

Check the 'Yes, I agree' and click on 'Create your Twitter application'.



## Step 15:

Now your application will be created.

**Step 16:**

Then click on Keys and Access tokens tab. The Consumer key and Consumer Secret keys will be generated.
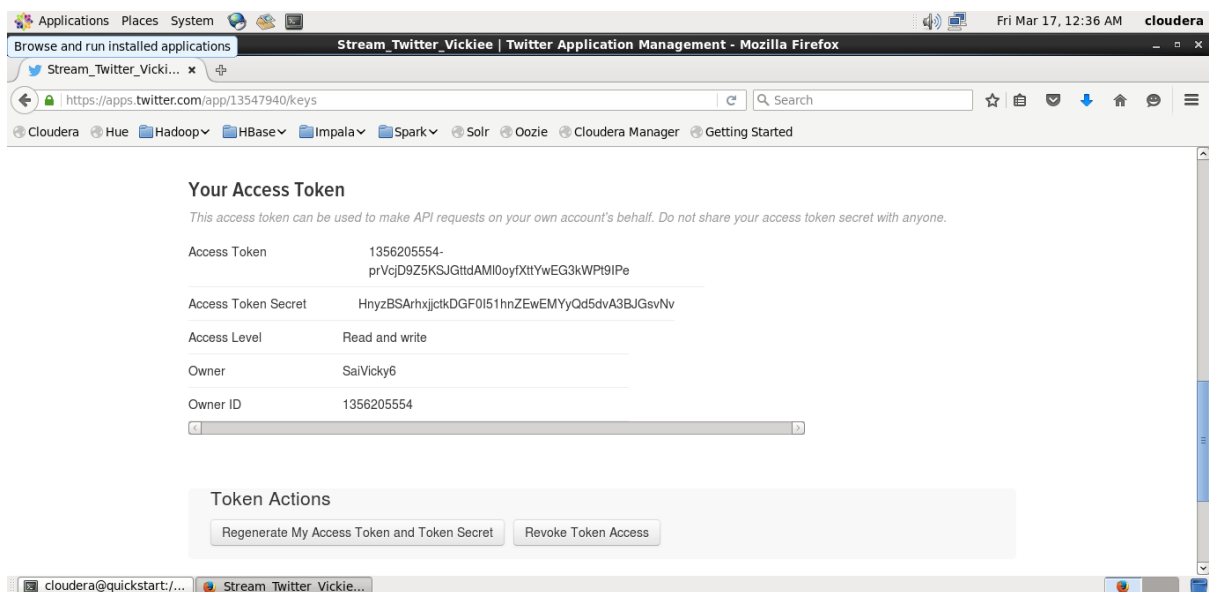


**Step 17:**

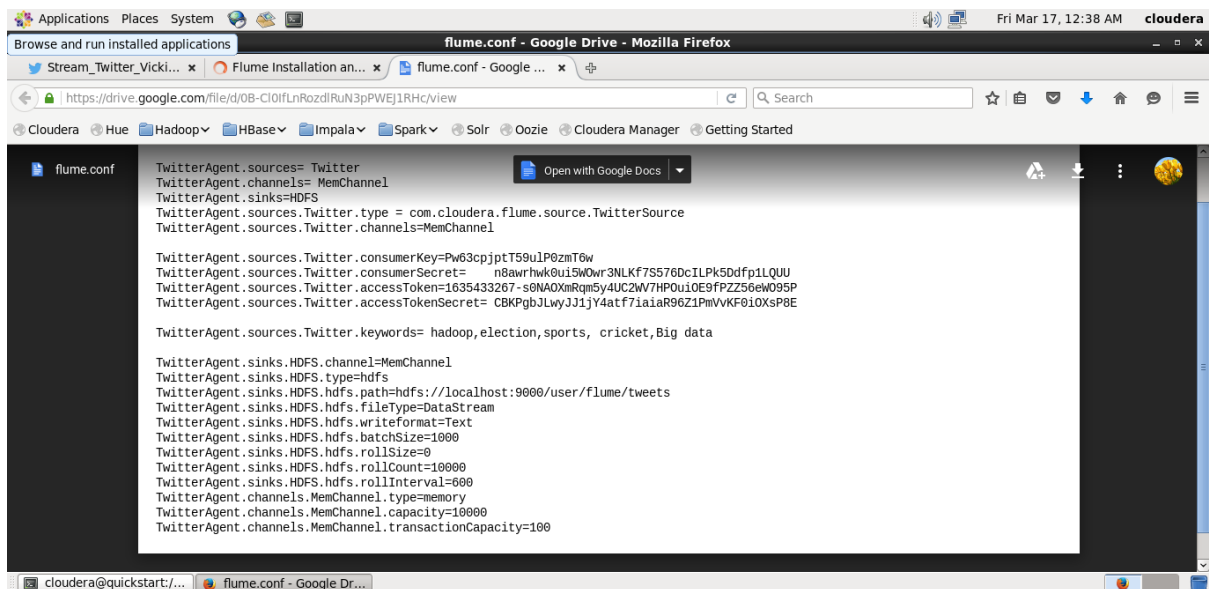Scroll down. Click on 'Create my Access token'.

**Step 18:**

The Access token and Access token Secret will be generated.



**Step 19:**

Then download the flume.conf file from https://drive.google.com/file/d/0B-Cl0IfLnRozdlRuN3pPWEJ1RHc/view?usp=sharing

**Step 20:**

Move the downloaded flume.conf file to configuration folder of flume-ng. Super user permission is given for this modification.

**Command:** sudo cp /home/cloudera/Downloads/flume.conf /usr/lib/flume-ng/conf/
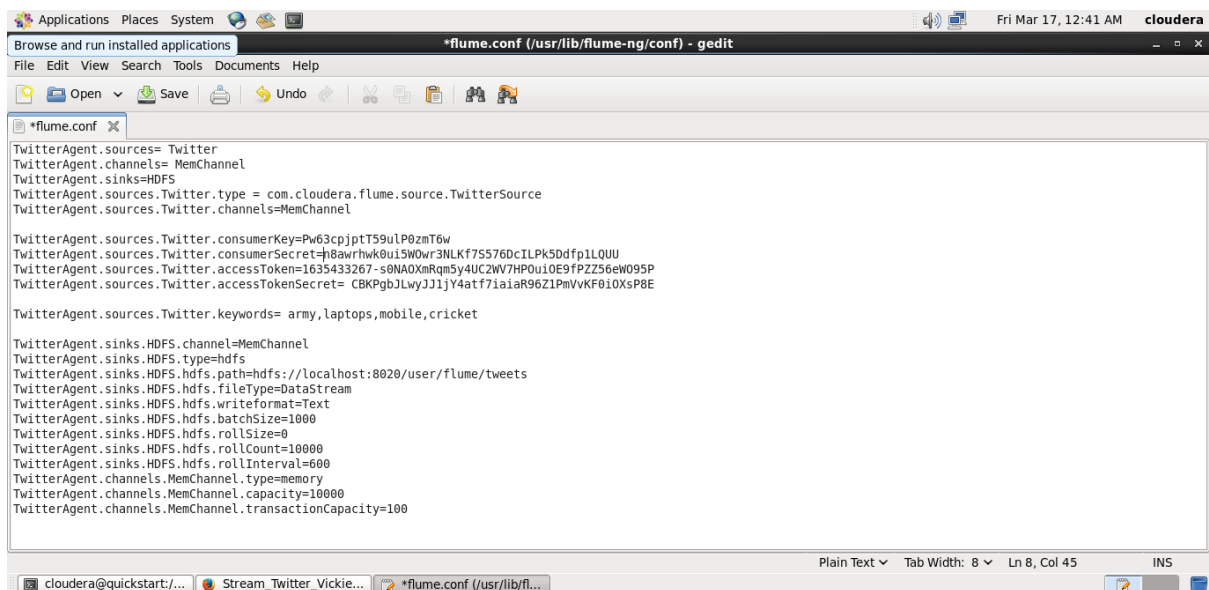


**Step 21:**

Open flume.conf file using gedit to make changes in the access credentials.

**Command:** sudo gedit conf/flume.conf

**Step 22:**

Replace the TwitterAgent.sources.Twitter keys and tokens with the one provided by the twitter in the application Keys and Access tokens steps. Replace the keywords of TwitterAgent.sources.Twitter.keywords to the one which we are analyzing. Also change the

localhost number of TwitterAgent.sinks.HDFS.hdfs.path to the localhost number of your host. Save the file and close it.



## Step 23:

Change the permissions for flume directory. Super user permission is given for this modification.

**Command:** sudo chmod -R 755 /usr/lib/flume-ng/



## Step 24:

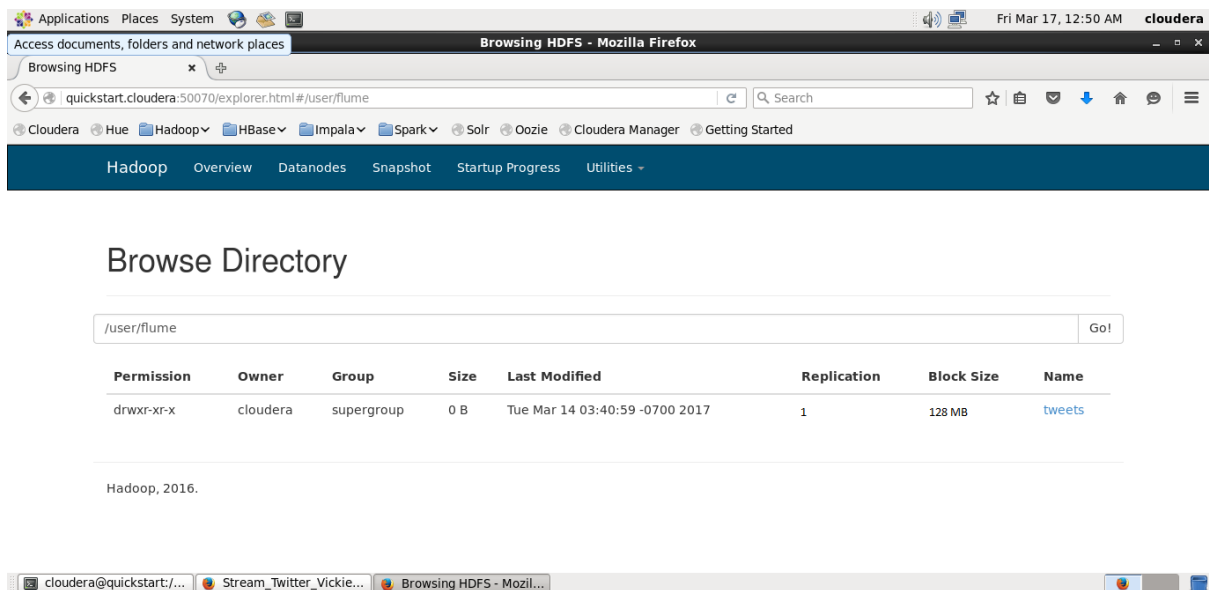The final step is to start fetching the data from twitter.

**Command:** ./bin/flume-ng agent -n TwitterAgent -c conf -f /usr/lib/ flume-ng/conf/flume.conf

Now the live data is streamed from Twitter. Wait for about 50 to 60 seconds to fetch the streaming data from twitter. Ctrl+C is pressed to terminate the streaming. The system may throw few exceptions as the process is terminated in the middle but ignore it.

**Step 25:**

Now the browser is opened in the the VM machine. Locate the tweets data path in HDFS at /user/flume/tweets.



**Step 26:**

The data is yielded as shown in the snapshot and modified according to demands of the user.

**Conclusion:**

Thus the Twitter data data is successfully streamed from the Twitter using Apache Flume engine.