| Name: Subrato Tapaswi | Class/Roll No.: D16AD/60 | Grade: |
|---|---|---|

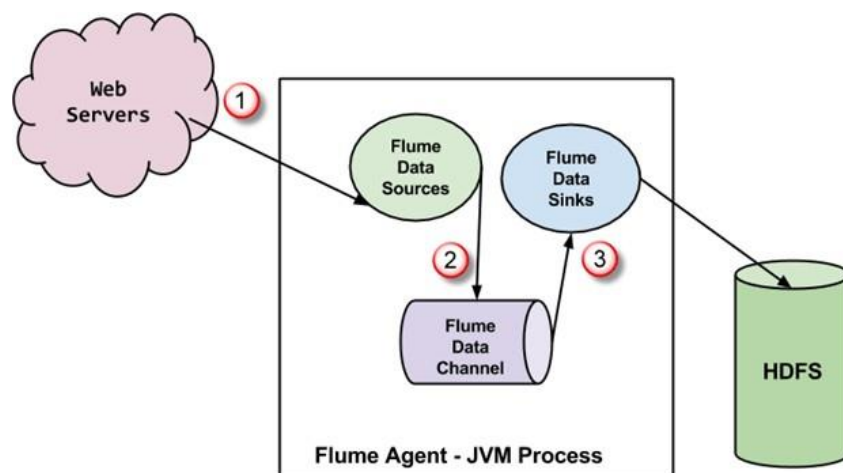**Title of Experiment:** To study Twitter data analysis using Flume.

**Objective of Experiment:** To teach the fundamental techniques and principles in achieving big data analytics with scalability and streaming capability. To enable students to have skills that will help them to solve complex real-world problems in decision support.

**Outcome of Experiment:** Solve Complex real-world problems in various applications like recommender systems, social media applications, health and medical systems, etc.

**Problem Statement:** To study Twitter data analysis using Flume.

**Theory:**

- **Apache Flume:** Apache Flume is a reliable and distributed system for collecting, aggregating and moving massive quantities of log data. It has a simple yet flexible architecture based on streaming data flows. Apache Flume is used to collect log data present in log files from web servers and aggregate it into HDFS for analysis.

- **Flume Architecture:** A Flume agent is a JVM process which has 3 components –Flume Source, Flume Channel and Flume Sink– through which events propagate after initiated at an external source.
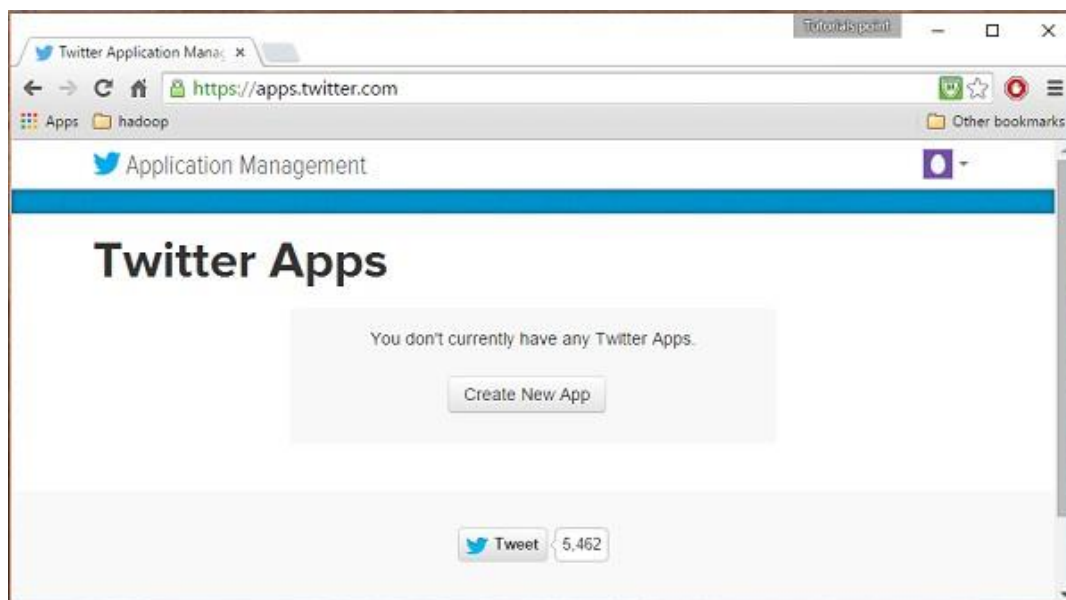
➢ **Flume Architecture:**

1. In the above diagram, the events generated by an external source (WebServer) are consumed by Flume Data Source. The external source sends events to the Flume source in a format that is recognized by the target source.

2. Flume Source receives an event and stores it into one or more channels. The channel acts as a store which keeps the event until it is consumed by the flume sink. This channel may use a local file system in order to store these events.

3. Flume sink removes the event from a channel and stores it into an external repository like e.g., HDFS. There could be multiple flume agents, in which case flume sinks forwards the event to the flume source of the next flume agent in the flow.

## Output:

**Step 1:** To create a Twitter application, click on the following link
https://apps.twitter.com/



**Step 2:** Click on the Create New App button. You will be redirected to a window where you will get an application form in which you have to fill in your details in order to create the App.

**Step 3:** Under keys and Access Tokens tab at the bottom of the page,



**Step 4: Install / Verify Hadoop:** Install Hadoop. If Hadoop is already installed in your system, then verify it. It will show output in the following manner:

```
$ hadoop version
```

```
Hadoop 2.6.0
Subversion https://git-wip-us.apache.org/repos/asf/hadoop.git -r
e3496499ecb8d220fba99dc5ed4c99c8f9e33bb1
Compiled by jenkins on 2014-11-13T21:10Z
Compiled with protoc 2.5.0
From source with checksum 18e43357c8f927c0695f1e9522859d6a
This command was run using /home/Hadoop/hadoop/share/hadoop/comm
```

**Step 5: Starting Hadoop:** Browse through the sbin directory of Hadoop and start yarn and Hadoop dfs (distributed filesystem) as shown below.

```
cd /$Hadoop_Home/sbin/
$ start-dfs.sh
localhost: starting namenode, logging to
   /home/Hadoop/hadoop/logs/hadoop-Hadoop-namenode-localhost.locald
localhost: starting datanode, logging to
   /home/Hadoop/hadoop/logs/hadoop-Hadoop-datanode-localhost.localde
Starting secondary namenodes [0.0.0.0]
starting secondarynamenode, logging to
   /home/Hadoop/hadoop/logs/hadoop-Hadoop-secondarynamenode-local

$ start-yarn.sh
starting yarn daemons
starting resourcemanager, logging to
   /home/Hadoop/hadoop/logs/yarn-Hadoop-resourcemanager-localhost.lc
localhost: starting nodemanager, logging to
   /home/Hadoop/hadoop/logs/yarn-Hadoop-nodemanager-localhost.local
```

```
$cd /$Hadoop_Home/bin/
$ hdfs dfs -mkdir hdfs://localhost:9000/user/Hadoop/twitter_data
```

**Step 6: Configuring Flume:** We have to configure the source, the channel, and the sink using the configuration file in the conffolder.



```
export CLASSPATH=$CLASSPATH:/FLUME_HOME/lib/*
```

**Step 7: Execution:** Browse through the Flume home directory and execute the application as shown below.

```
$ cd $FLUME_HOME
$ bin/flume-ng agent --conf ./conf/ -f conf/twitter.conf
Dflume.root.logger=DEBUG,console -n TwitterAgent
```

**Verifying HDFS**





**CONCLUSION:** The study concluded that Flume technology can be used to extract real-time data from Twitterand store it in HDFS for further analysis.