

ASSIGNMENT 11.3

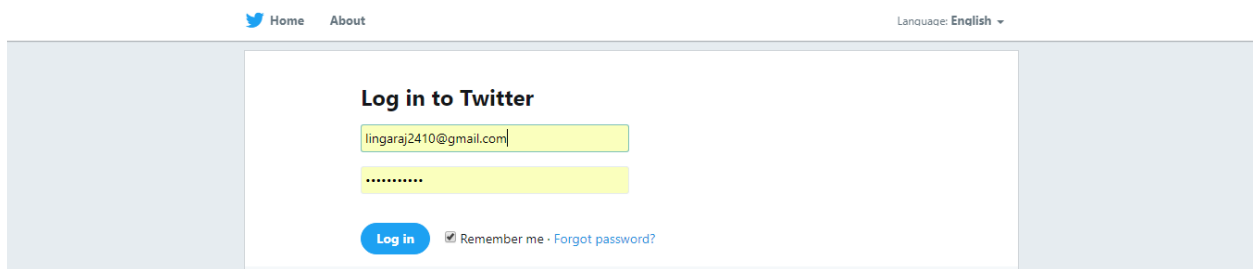
Problem Statement:

Create a flume agent that streams data from Twitter and stores in the HDFS.

Solution:

Twitter data streaming into HDFS using Flume:

Step 1: Login to the twitter account

A screenshot of the Twitter login page. At the top, there are links for 'Home' and 'About', and a language selector set to 'English'. The main heading is 'Log in to Twitter'. Below it, there are two input fields: the first contains the email 'lingaraj2410@gmail.com' and the second is masked with dots. A blue 'Log in' button is positioned below the password field. To the right of the button is a checked checkbox labeled 'Remember me' and a link for 'Forgot password?'.

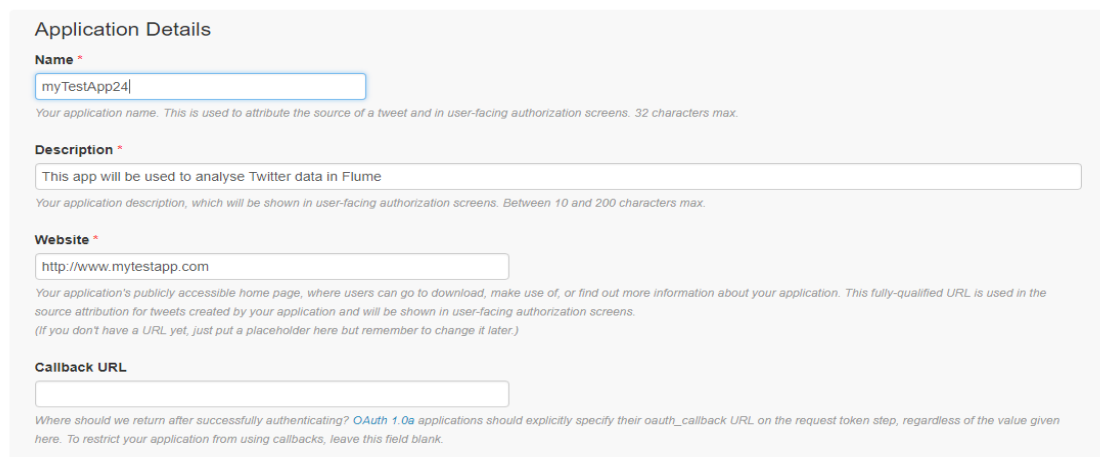
Step 2: Go to the following link and click the 'create new app' button.

<https://apps.twitter.com/app>

A screenshot of the 'Twitter Apps' management page. The header includes the Twitter logo, 'Application Management', and a user profile icon. The main heading is 'Twitter Apps'. On the right side, there is a button labeled 'Create New App'.

Step 3: Enter the necessary details.

Create an application



A screenshot of the 'Create an application' form. The form is titled 'Application Details' and contains four sections: 'Name' with the value 'myTestApp24', 'Description' with the text 'This app will be used to analyse Twitter data in Flume', 'Website' with the URL 'http://www.mytestapp.com', and 'Callback URL' which is currently empty. Each section has a small text block below it providing instructions or constraints for the input.

Step 4: Accept the developer agreement and select the ‘create your Twitter application’ button.

Developer Agreement
☒ Yes, I have read and agree to the [Twitter Developer Agreement](#).


Create your Twitter application

Step 5: Select the ‘Keys and Access Token’ tab.

 Application Management 

myTestApp24 Test OAuth

[Details](#) [Settings](#) [Keys and Access Tokens](#) [Permissions](#)

 This app will be used to analyse Twitter data in Flume
<http://www.mytestapp.com>

Organization
Information about the organization or company associated with your application. This information is optional.
Organization
Organization website

Application Settings
Your application's Consumer Key and Secret are used to [authenticate](#) requests to the Twitter Platform.
Access level

Step 6: Copy the consumer key and the consumer secret code.

Step 7: Scroll down further and select the ‘create my access token’ button.

Your Access Token

You haven't authorized this application for your own account yet.

By creating your access token here, you will have everything you need to make API calls right away. The access token generated will be assigned your application's current permission level.

Token Actions

Create my access token

Now, you will receive a message stating “that you have successfully generated your application access token”.

Status
Your application access token has been successfully generated. It may take a moment for changes you've made to reflect.
[Refresh](#) if your changes are not yet indicated.

myTestApp24 Test OAuth

[Details](#) [Settings](#) [Keys and Access Tokens](#) [Permissions](#)

Application Settings
Keep the "Consumer Secret" a secret. This key should never be human-readable in your application.

Step 8: Copy the Access Token and Access token Secret code.

Step 9: Create a configuration file that is used to specify Flume properties in compliance with Twitter streaming. Save this file in the ‘conf’ deirectory of Flume home path.

Here is the properties file I have created that includes Twitter app credentials as well as the HDFS path into which Tweets will be stored:

```
TwitterAgent.sources = Twitter
TwitterAgent.channels = MemChannel
TwitterAgent.sinks = HDFS

# Describing/Configuring the source
TwitterAgent.sources.Twitter.type = org.apache.flume.source.twitter.TwitterSource
TwitterAgent.sources.Twitter.consumerKey=0r1e13F3XNV3XtyZ2XKKK490ac
TwitterAgent.sources.Twitter.consumerSecret=L2aNWtNxxeKUGHME71RDIPg6UjCB44il0fGfmmQqHV68my7KFd
TwitterAgent.sources.Twitter.accessToken=2221479913-T6KiphqccmFYjzIL6k6AZDG0RS6R8iINMD6SDW4
TwitterAgent.sources.Twitter.accessTokenSecret=IW6Rm3eDnaKeu4leiJT4tzmr7bx0heJmJAvpwIQsNQHQQL
TwitterAgent.sources.Twitter.keywords=hadoop, bigdata, mapreduce, spark, hbase, nosql

# Describing/Configuring the sink

TwitterAgent.sources.Twitter.keywords= hadoop,election,sports, cricket,Big data

TwitterAgent.sinks.HDFS.channel=MemChannel
TwitterAgent.sinks.HDFS.type=hdfs
TwitterAgent.sinks.HDFS.hdfs.path=hdfs://localhost:9000/user/flume/tweets
TwitterAgent.sinks.HDFS.hdfs.fileType=DataStream
TwitterAgent.sinks.HDFS.hdfs.writeFormat=Text
TwitterAgent.sinks.HDFS.hdfs.batchSize=1000
TwitterAgent.sinks.HDFS.hdfs.rollSize=0
TwitterAgent.sinks.HDFS.hdfs.rollCount=10000
TwitterAgent.sinks.HDFS.hdfs.rollInterval=600

TwitterAgent.channels.MemChannel.type=memory
TwitterAgent.channels.MemChannel.capacity=10000
TwitterAgent.channels.MemChannel.transactionCapacity=1000

TwitterAgent.sources.Twitter.channels = MemChannel
TwitterAgent.sinks.HDFS.channel = MemChannel
```

We can mention which keywords tweet data to be collected from the Twitter application. We can change the keywords in the TwitterAgent.sources.Twitter.keywords property. In our example, we are fetching tweet data related to Hadoop, election, sports, cricket and Big data.

Step 10: Open a new terminal and start all the Hadoop daemons, before running the flume command to fetch the twitter data. Use the ‘jps’ command to see the running Hadoop daemons.

```
[acadgild@localhost sbin]$ ./start-all.sh
This script is Deprecated. Instead use start-dfs.sh and start-yarn.sh
18/01/23 16:48:06 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java cl
asses where applicable
Starting namenodes on [localhost]
localhost: starting namenode, logging to /usr/local/hadoop-2.6.0/logs/hadoop-acadgild-namenode-localhost.localdomain.out
localhost: datanode running as process 2810. Stop it first.
Starting secondary namenodes [0.0.0.0]
0.0.0.0: starting secondarynamenode, logging to /usr/local/hadoop-2.6.0/logs/hadoop-acadgild-secondarynamenode-localhost.localdomain.out
18/01/23 16:48:32 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java cl
asses where applicable
starting yarn daemons
resourcemanager running as process 3101. Stop it first.
localhost: nodemanager running as process 3204. Stop it first.
[acadgild@localhost sbin]$ jps
4082 Jps
3204 NodeManager
3881 SecondaryNameNode
2810 DataNode
3101 ResourceManager
3663 NameNode
[acadgild@localhost sbin]$ █
```

Step 11: Create a new directory inside HDFS path, where the Twitter tweet data should be stored.

```
[acadgild@localhost sbin]$ hadoop fs -mkdir /user/acadgild/twitter_data
18/01/23 16:53:27 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java cl
asses where applicable
[acadgild@localhost sbin]$ hadoop fs -ls /user/acadgild/
18/01/23 16:53:48 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java cl
asses where applicable
Found 6 items
drwxr-xr-x - acadgild supergroup          0 2015-11-20 11:46 /user/acadgild/Pictures
drwxr-xr-x - acadgild supergroup          0 2018-01-23 12:32 /user/acadgild/_sqoop
drwxr-xr-x - acadgild supergroup          0 2018-01-23 10:13 /user/acadgild/emp_info
drwxr-xr-x - acadgild supergroup          0 2015-11-17 02:03 /user/acadgild/oozie-acad
drwxr-xr-x - acadgild supergroup          0 2015-11-17 02:00 /user/acadgild/share
drwxr-xr-x - acadgild supergroup          0 2018-01-23 16:53 /user/acadgild/twitter_data
[acadgild@localhost sbin]$
```

```
$ flume-ng agent -n TwitterAgent -f /usr/local/flume/conf/acadgild.conf
```

We can derive from the logs above that the Flume agent is able to fetch Tweets from Twitter and storing the related data into HDFS directory which we have mentioned in the configuration file.

```
[acadgild@localhost ~]$ hadoop ls -ls /user/acadgild/twitter_data
18/01/25 18:28:58 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Found 1 items
-rw-r--r-- 1 acadgild supergroup 3917637 2018-01-25 18:19 /user/acadgild/twitter_data/FlumeData.1516884393277
[acadgild@localhost ~]$
```

트 | 판영 넬영 웅영 책영 Wanna One || ♥발한소녀단 || 액소♥ || 세븐틴 (이중석·박보검) (서라,린) ♥[chanbaek #kookmin #mino
(필크보이) [jokik 67(2018-01-25T18:19:12Z@RT @BESTFRIEND_92: 180125 BAEKHYIN IG LIVE
ထေ့ထေ့ထေ့ : လာဘ်လော့မာကတုတ်? ကိုယ့်ပုံပေါ်ပေါ်ပေါ်ပေါ် '။' [baekhyin href="http://twitter.com/download/iphone" rel="nofollow">Twitter for
iPhone [956509257209229319 [국 국: 발한소녀단: 02:サササ] デス: @BTS_twt [956509257209229319
8:19:12Z@jiminnnn_1013 めちゃくそ美人よね。がんばってみる。 [baekhyin href="http://twitter.com/download/iphone" rel="no
follow">Twitter for iPhone [9565092572260720640 [一ピット [紳士です。フォローRT等お気軽に。
2017年までゲーム業界人でしたがマゾで楽しみたいので脱社畜が目標。
豆腐メンタル/Destiny大好き /代表作(雑用) ネットハイP /乱気力グラほか高木P作品 /幕末Rock
/ブログhttp://bit.ly/Gentleblog
/動画http://bit.ly/GentleYT [9565092572260720640 [AmNo1Gentleman(2018-01-25T18:19:12Z@DirtyMeganeMAN 配信しながらや
っていいですか! [baekhyin href="http://twitter.com/" rel="nofollow">Twitter Web Client [9565092572260720640 [9565092572260720640
응원하다 정용화 CNBLUE [baekhyin 한사랑 [타소리+드림소리+피아노소리 [BAKEUP Realize
♥연재나하이팅 [baekhyin 그대 [baekhyin 날래나래 [baekhyin 나RaeJYH(2018-01-25T18:19:12Z@RT @CNlove_binbin: #경희대 진실을 밝혀
We will stay young forever [baekhyin href="http://twitter.com/download/android" rel="nofollow">Twitter for Android [9565092572260720640 [9565092572260720640
0000 0T
/
{ "type": "record", "name": "Doc", "doc": "adoc", "fields": { "name": "id", "type": "string", "name": "user_friends_count", "type": ["int",
"null"], "name": "user_location", "type": ["string", "null"], "name": "user_description", "type": ["string", "null"], "name": "use
r_statuses_count", "type": ["int", "null"], "name": "user_followers_count", "type": ["int", "null"], "name": "user_name", "type": ["s
tring", "null"], "name": "user_screen_name", "type": ["string", "null"], "name": "created_at", "type": ["string", "null"], "name": "
text", "type": ["string", "null"], "name": "retweet_count", "type": ["long", "null"], "name": "retweeted", "type": ["boolean", "null"],
"name": "in_reply_to_user_id", "type": ["long", "null"], "name": "source", "type": ["string", "null"], "name": "in_reply_to statu
s_id", "type": ["long", "null"], "name": "media_url_https", "type": ["string", "null"], "name": "expanded_url", "type": ["string", "nu
ll"] } } [956509257226055681 [みあか \$全裸で着衣泳 [956509257226055681
nwdxf(2018-01-25T18:19:12Z@RT @vuzuki lestokvo: 歌ってのゆづきさん

Since we haven't mentioned any specified language for Tweets, we are seeing tweets from various languages. This is how we can fetch Twitter data in the real time with the use of **Apache Flume**.