# Module 10:  Apache Oozie and Hadoop Project

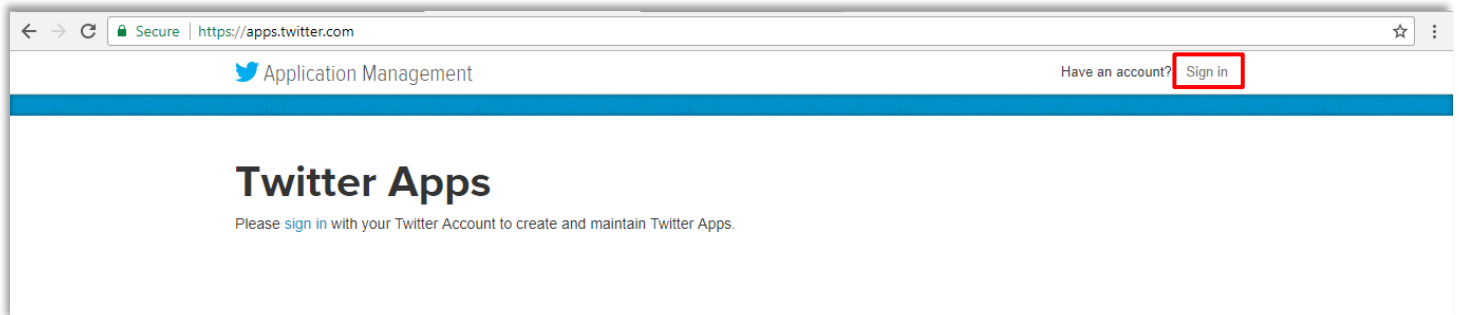## Twitter Data Streaming using Apache Flume

**edureka!**
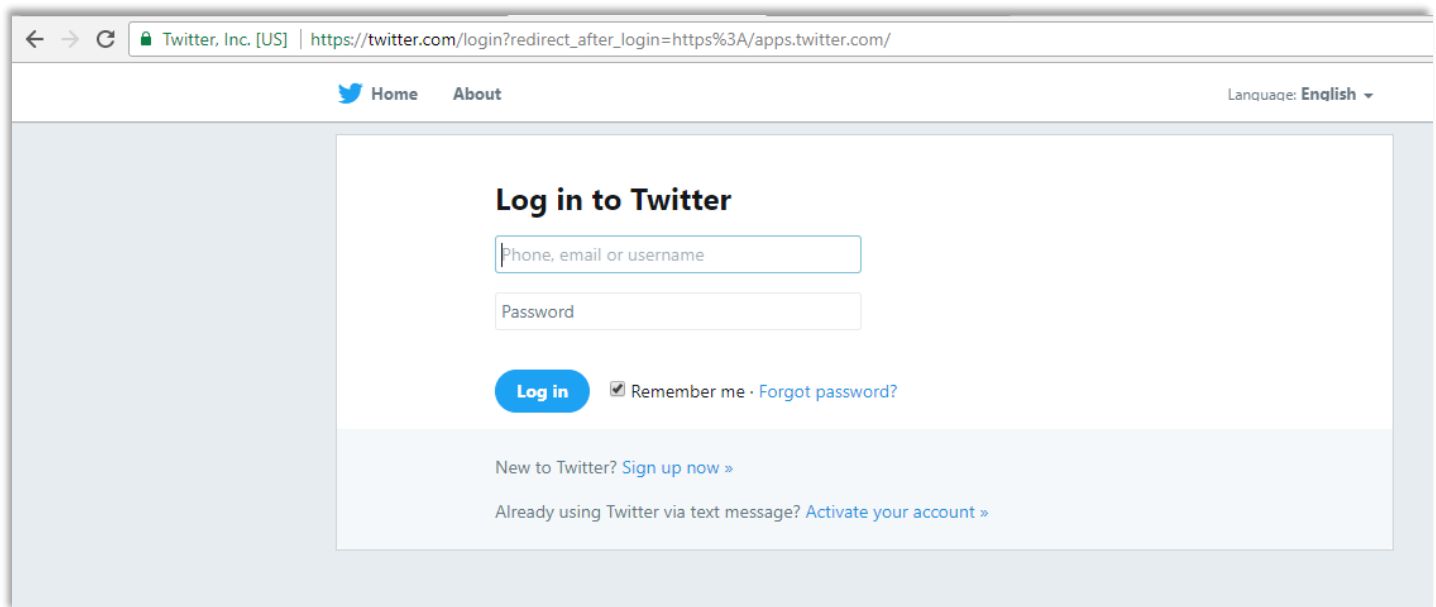
**edureka!**

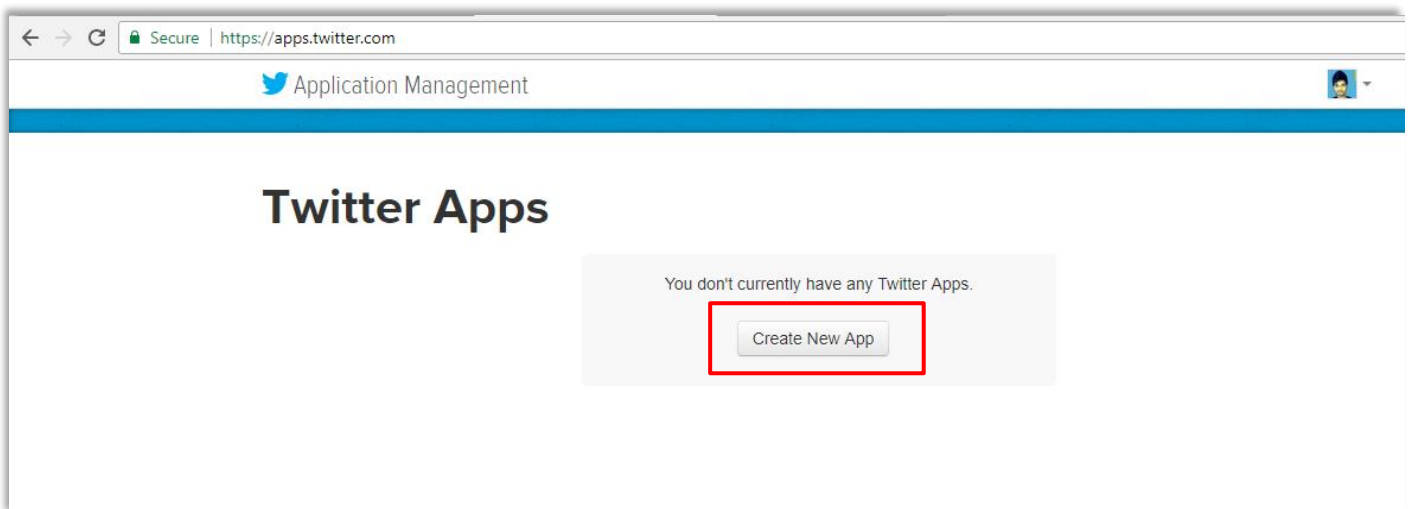**Step 1:** Go to "apps.twitter.com" and click on "Sign in"



**Step 2:** Log in to your twitter account

**Step 3:** Click on "Create New App"



**Step 4:** Enter the application details

**Step 5:** Accept the Developer Agreement and click on "Create your Twitter application"



**Step 6:** Click on "Keys and Access Tokens"

**Step 7:** Note the "Consumer Key", "Consumer Secret" that is generated and click on "create my access token"

**Step 8:** Here is the generated "Access Token" and "Access Token Secret" -



**Step 9:** Write the following configurations in a file "flume_twitter.conf"

**Note:** Here we will use the "Consumer Key", "Consumer Secret", "Access Token" and "Access Token Secret" generated in the above steps.

TwitterAgent.sources = Twitter

TwitterAgent.channels = MemChannel

TwitterAgent.sinks = HDFS

TwitterAgent.sources.Twitter.type = org.apache.flume.source.twitter.TwitterSource

TwitterAgent.sources.Twitter.channels = MemChannel

TwitterAgent.sources.Twitter.consumerKey = zSeSMnxGBLvzhHp1

TwitterAgent.sources.Twitter.consumerSecret = SybN56o9OlPfl23OpJL6C8ahMSSdGuOmLOlUL8Lh

TwitterAgent.sources.Twitter.accessToken = 296124655-s0YY73prdUmZatDvJeMsTKv4ZJ4OfT

TwitterAgent.sources.Twitter.accessTokenSecret = FmBJAadWnOqY2i4Ck9ecssaQWYsCPLhmw4jjb


TwitterAgent.sources.Twitter.keywords = spark, scientist, hadoop, big data, analytics, bigdata, cloudera, data science, data scientist, business intelligence, mapreduce, data warehouse, data warehousing, mahout, hbase, nosql, newsql, businessintelligence, cloudcomputing


TwitterAgent.sinks.HDFS.channel = MemChannel

TwitterAgent.sinks.HDFS.type = hdfs

TwitterAgent.sinks.HDFS.hdfs.path = hdfs://nameservice1/user/edureka_249489/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

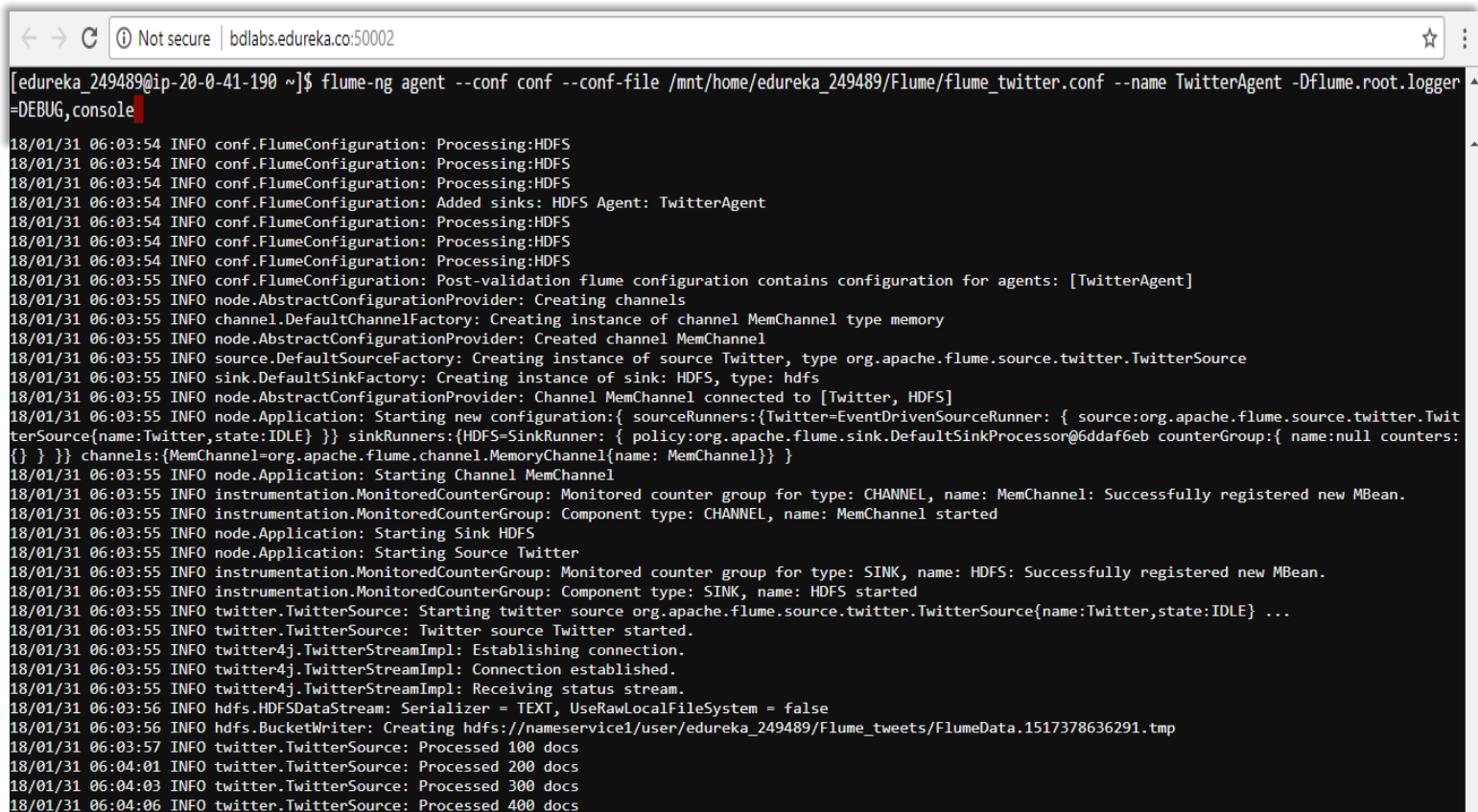TwitterAgent.sinks.HDFS.hdfs.rollInterval = 600


TwitterAgent.channels.MemChannel.type = memory

TwitterAgent.channels.MemChannel.capacity = 10000

TwitterAgent.channels.MemChannel.transactionCapacity = 100

```
flume_twitter - Notepad
File  Edit  Format  View  Help
TwitterAgent.sources = Twitter
TwitterAgent.channels = MemChannel
TwitterAgent.sinks = HDFS

TwitterAgent.sources.Twitter.type = org.apache.flume.source.twitter.TwitterSource
TwitterAgent.sources.Twitter.channels = MemChannel
TwitterAgent.sources.Twitter.consumerKey = zSeSMnxGBLvzhl
TwitterAgent.sources.Twitter.consumerSecret = SybN56o9Oll
TwitterAgent.sources.Twitter.accessToken = 296124655-s0Y
TwitterAgent.sources.Twitter.accessTokenSecret = FmBJAadl
TwitterAgent.sources.Twitter.keywords = spark, scientist, hadoop, big data, analytics, bigdata, cloudera, data science

TwitterAgent.sinks.HDFS.channel = MemChannel
TwitterAgent.sinks.HDFS.type = hdfs
TwitterAgent.sinks.HDFS.hdfs.path = hdfs://nameservice1/user/edureka_249489/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 = 100
```

Step 10: Run the following command to capture twitter data

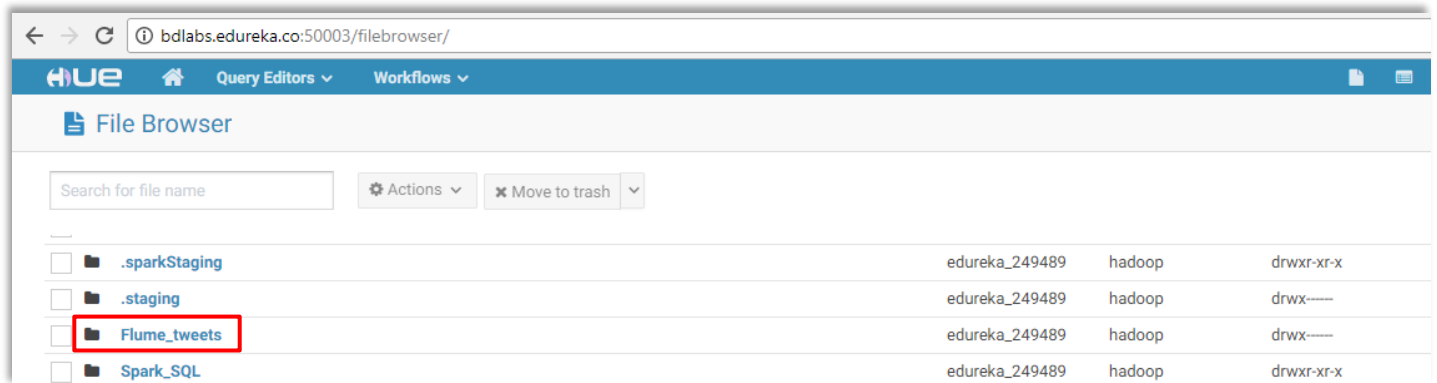**Command:** `flume-ng agent --conf conf --conf-file /mnt/home/edureka_249489/Flume/flume_twitter.conf --name TwitterAgent -Dflume.root.logger=DEBUG,console`
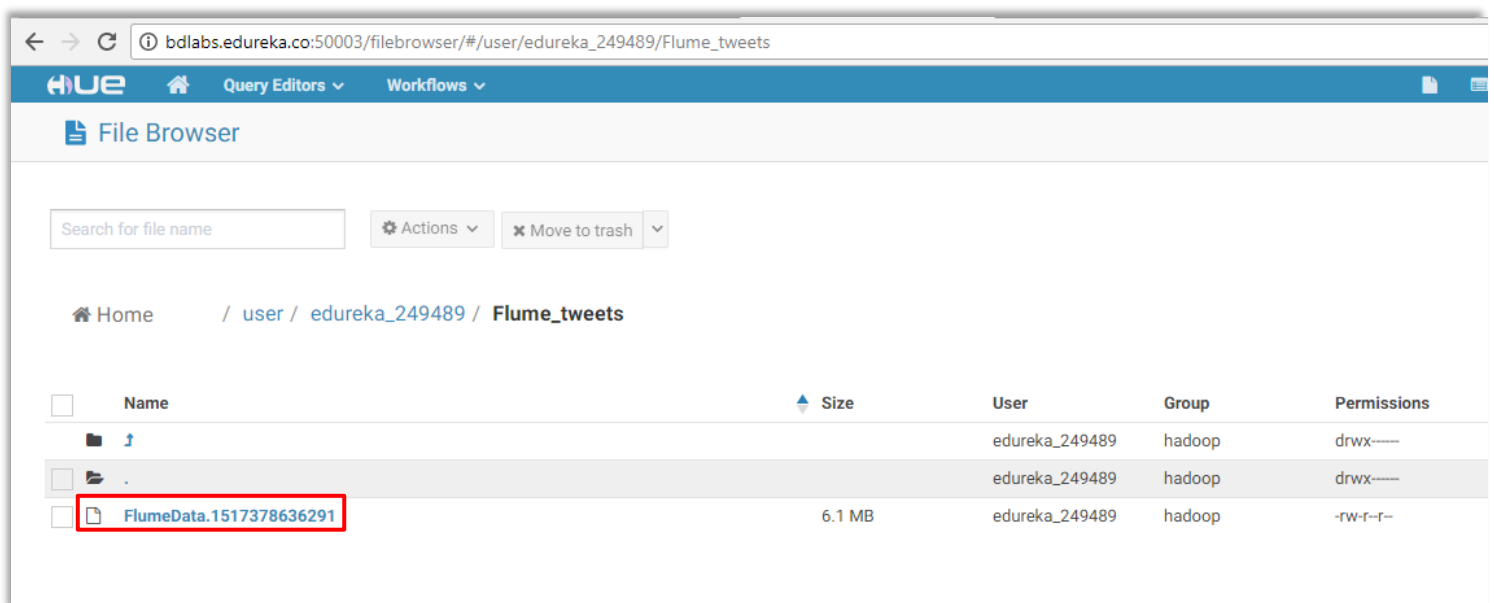
**Step 11:** Go to Hue and browse to the HDFS directory and click on "Flume_tweets"



**Step 12:** Click on "FlumeData.1517378636291" file

**Step13:** Check the output. This is the data that has been downloaded from Twitter.



**Note:** Press 'ctrl + c' on webconsole to stop the generation of tweets.