# BUAN 6346.501

# Bonus Project - Flume Twitter

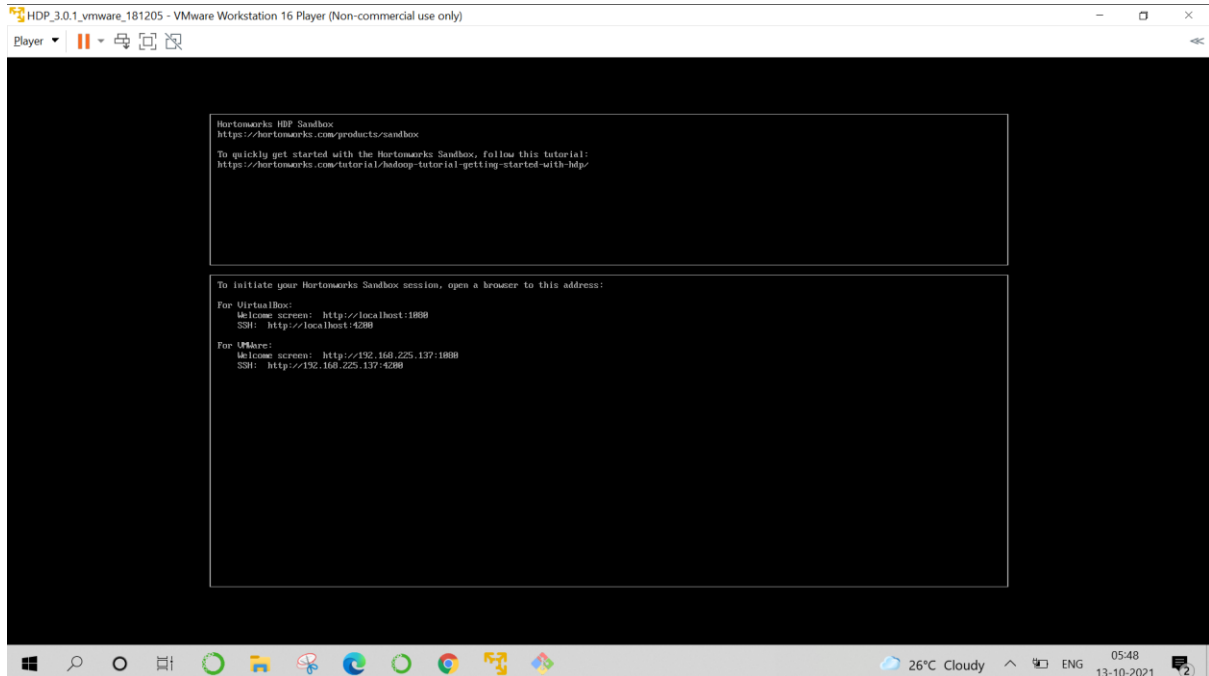**Arul Chakravarthy**

**Twitter-Flume Integration**

**Step 1:** New Twitter Developer Account was created using https://developer.twitter.com/

**Step 2:** For Hadoop ecosystem, Sandbox HDP Horton VMWare V3.0.1 was downloaded and installed from Cloudera



From above, my Sandbox IP is 192.168.225.137

**Step 3:** For Windows, Hosts.cfg has been configured with above HDP Sandbox IP and its host name mapping as follows

```
192.168.225.137 localhost sandbox-hdp.hortonworks.com sandbox-hdf.hortonworks.com
```
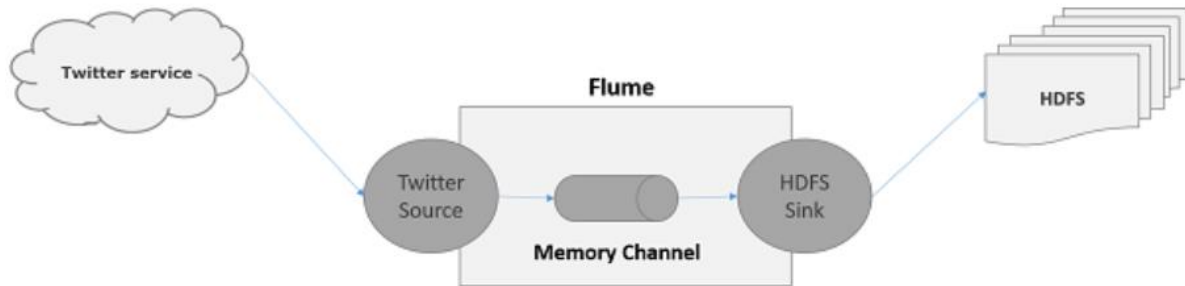
**Hosts.cfg configuration:**

```
# Copyright (c) 1993-2009 Microsoft Corp.
#
# This is a sample HOSTS file used by Microsoft TCP/IP for Windows.
#
# This file contains the mappings of IP addresses to host names. Each
# entry should be kept on an individual line. The IP address should
# be placed in the first column followed by the corresponding host name.
# The IP address and the host name should be separated by at least one
# space.
#
# Additionally, comments (such as these) may be inserted on individual
# lines or following the machine name denoted by a '#' symbol.
#
# For example:
#
#      102.54.94.97     rhino.acme.com          # source server
#       38.25.63.10     x.acme.com              # x client host

# localhost name resolution is handled within DNS itself.
#       127.0.0.1       localhost
#       ::1             localhost
192.168.225.137 localhost sandbox-hdp.hortonworks.com sandbox-hdf.hortonworks.com
```

Arul Chakravarthy

# Twitter-Flume Integration

Let's see in this document on how we can set-up Flume to read the tweets into HDFS

**Destination HDFS Directory:** /loudacre/tweets



**Step 4:** For Windows, Git Bash was used throughout to login into HDP-Sandbox by SSH as root user via port 2222

```
arulc@LAPTOP-70QSL9Q7 MINGW64 ~
$ ssh root@sandbox-hdp.hortonworks.com -p 2222
root@sandbox-hdp.hortonworks.com's password:
You are required to change your password immediately (root enforced)
Last login: Wed Oct 13 22:33:18 2021 from 172.18.0.3
Changing password for root.
(current) UNIX password:
New password:
Retype new password:
```

**Step 4:** A simple check whether Hadoop is installed

```
[root@sandbox-hdp ~]# hadoop version
Hadoop 3.1.1.3.0.1.0-187
Source code repository git@github.com:hortonworks/hadoop.git -r 2820e4d6fc7ec31ac42187083ed5933c823e9784
Compiled by jenkins on 2018-09-19T10:19Z
Compiled with protoc 2.5.0
From source with checksum 889327faf5a6ca5fc06fcf97c13af29
This command was run using /usr/hdp/3.0.1.0-187/hadoop/hadoop-common-3.1.1.3.0.1.0-187.jar
```

**Step 5:** Latest version of Apache-Flume was downloaded from the below link in my local laptop https://downloads.apache.org/flume/1.7.0/ . From local laptop, apache-flume zip file has been transferred into HDP Sandbox using SCP

```
arulc@LAPTOP-70QSL9Q7 MINGW64 ~/Downloads
$ scp -P 2222 apache-flume-1.7.0-bin.tar.gz root@sandbox-hdp.hortonworks.com:/root
root@sandbox-hdp.hortonworks.com's password:
apache-flume-1.7.0-bin.tar.gz                                    100%   53MB  77.4MB/s   00:00
```

```
[root@sandbox-hdp ~]# ls
anaconda-ks.cfg  apache-flume-1.7.0-bin.tar.gz
```

**Step 6:** Apache-flume zip file has been unzipped as follows using sudo access

```
[root@sandbox-hdp ~]# sudo tar xzf apache-flume-1.7.0-bin.tar.gz
[root@sandbox-hdp ~]#
```

Arul Chakravarthy

**Step 7:** Unzipped apache-flume was moved into /usr/local/flume/ of my sandbox

```
[root@sandbox-hdp ~]# sudo mv apache-flume-1.7.0-bin  /usr/local/flume/
[root@sandbox-hdp ~]# |
```

**Step 8:** Necessary environment variables has been set as FLUME_PATH and the bin directory of flume has been added to the PATH variable which will be used in the further steps

```
[root@sandbox-hdp bin]# export FLUME_HOME=/usr/local/flume/apache-flume-1.7.0-bin/
[root@sandbox-hdp bin]# export PATH=$PATH:$FLUME_HOME/bin
```

**Step 9:** Check for flume-ng version

```
[root@sandbox-hdp bin]# flume-ng version
Flume 1.7.0
Source code repository: https://git-wip-us.apache.org/repos/asf/flume.git
Revision: 511d868555dd4d16e6ce4fedc72c2d1454546707
Compiled by bessbd on Wed Oct 12 20:51:10 CEST 2016
From source with checksum 0d21b3ffdc55a07e1d08875872c00523
```

**Step 10:** Once Flume installation was done, Twitter application was created using the following settings

**Edit App details**

**App name**

BUAN6346.501_FlumeTwitter_Proj

Maximum length 32 characters                                      2

**App icon**

Upload

Maximum size of 700k, JPG, GIF, PNG

**Description**

Briefly describe your App.

This app was created to use the Twitter API to injest feeds into HDFS using Flume for BUAN6346.501 project

Between 10 and 200 characters                                      94

Arul Chakravarthy

↩ **App permissions**                    ✏ **Edit**

Read and Write

Read + Post Tweets and profile information

✓ **Authentication settings**                    ✏ **Edit**

3-legged OAuth is disabled                    👤ˣ

Use 3-legged OAuth for Sign in with Twitter, posting Tweets on behalf of other accounts and more. Get more information in the docs.

**Step11:** Destination directory has been created as /loudacre/tweets in Sandbox HDFS to store the tweets as follows

```
[root@sandbox-hdp ~]# sudo -u hdfs hadoop fs -mkdir -p /loudacre/tweets
[root@sandbox-hdp ~]# |
```

Changed the HDFS destination directory owner, as root for root user to write tweets into /loudacre/tweets

```
[root@sandbox-hdp conf]# sudo -u hdfs hadoop fs -chown root /loudacre/tweets
```

**Step 12:** flume-env.sh configuration file of Apache Flume was editied using flume-env.sh.template pre-created during installation as reference and added JAVA_PATH and CLASSPATH as environment variables

```
[root@sandbox-hdp conf]# cd $FLUME_HOME/conf/
[root@sandbox-hdp conf]# cp flume-env.sh.template flume-env.sh
[root@sandbox-hdp conf]# vi flume-env.sh
```

Arul Chakravarthy

**Flume-env.sh**

```
root@sandbox-hdp:/usr/local/flume/apache-flume-1.7.0-bin/conf
# Licensed to the Apache Software Foundation (ASF) under one
# or more contributor license agreements.  See the NOTICE file
# distributed with this work for additional information
# regarding copyright ownership.  The ASF licenses this file
# to you under the Apache License, Version 2.0 (the
# "License"); you may not use this file except in compliance
# with the License.  You may obtain a copy of the License at
#
#      http://www.apache.org/licenses/LICENSE-2.0
#
# Unless required by applicable law or agreed to in writing, software
# distributed under the License is distributed on an "AS IS" BASIS,
# WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
# See the License for the specific language governing permissions and
# limitations under the License.

# If this file is placed at FLUME_CONF_DIR/flume-env.sh, it will be sourced
# during Flume startup.

# Enviroment variables can be set here.

# export JAVA_HOME=/usr/lib/jvm/java-6-sun

# Give Flume more memory and pre-allocate, enable remote monitoring via JMX
# export JAVA_OPTS="-Xms100m -Xmx2000m -Dcom.sun.management.jmxremote"

# Let Flume write raw event data and configuration information to its log files for debugging
# purposes. Enabling these flags is not recommended in production,
# as it may result in logging sensitive user information or encryption secrets.
# export JAVA_OPTS="$JAVA_OPTS -Dorg.apache.flume.log.rawdata=true -Dorg.apache.flume.log.printconfig=true "

# Note that the Flume conf directory is always included in the classpath.
#FLUME_CLASSPATH=""
export JAVA_HOME=/usr/lib/jvm/java-1.8.0-openjdk
export CLASSPATH=$CLASSPATH:/usr/local/flume/apache-flume-1.7.0-bin/lib/*
~
~
~
~
~
"flume-env.sh" 35L, 1687C
```

**Step13:** Under conf folder of apache-flume a new configuration file has been created as twitter.conf with necessary Flume Source, sink and channel configurations

**Twitter.conf**

```
arulc@LAPTOP-70QSL9Q7 MINGW64 ~
$ ssh root@sandbox-hdp.hortonworks.com -p 2222
root@sandbox-hdp.hortonworks.com's password:
Last login: Thu Oct 14 00:22:30 2021 from 172.18.0.3
[root@sandbox-hdp ~]# cat /usr/local/flume/apache-flume-1.7.0-bin/conf/twitter.conf
# Naming the components on the current agent.
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 = sfGUZ
TwitterAgent.sources.Twitter.consumerSecret = M3sty
TwitterAgent.sources.Twitter.accessToken = 1320
TwitterAgent.sources.Twitter.accessTokenSecret = FR0Nt9C
TwitterAgent.sources.Twitter.keywords = hadoop,hive, bigdata, mapreduce, sqoop, hbase, pig

# Describing/Configuring the sink
TwitterAgent.sinks.HDFS.type = hdfs
TwitterAgent.sinks.HDFS.hdfs.path = /loudacre/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

# Describing/Configuring the channel TwitterAgent.channels.MemChannel.type = memory
TwitterAgent.channels.MemChannel.type = memory
TwitterAgent.channels.MemChannel.capacity = 10000
TwitterAgent.channels.MemChannel.transactionCapacity = 100

# Binding the source and sink to the channel
TwitterAgent.sources.Twitter.channels = MemChannel
TwitterAgent.sinks.HDFS.channel = MemChannel
```

*ConsumerKey and AccessTokens are masked in the above snapshot owing to confidentiality*

Arul Chakravarthy

**Step14:** Flume Agent was run using the below command with parameters as *-n TwitterAgent* and *–conf-file <path-to-twitter.conf-file>* as follows to read the tweets from twitter

```
flume-ng agent -name TwitterAgent --conf-file $FLUME_PATH/conf/twitter.conf
```



TwitterAgent was started as below.,



**Step15: Successfully Flume has ingested Tweets from Twitter into HDFS /loudacre/tweets directory as follows**

**Step16:** Finally, a sanity check to display the contents of tweets from above HDFS destination directory
using `hdfs dfs -cat /loudacre/tweets/FlumeData.1634170937877 | head`

```
[root@sandbox-hdp ~]# hdfs dfs -cat /loudacre/tweets/FlumeData.1634170937877 | head
{"type":"record","name":"Doc","doc":"adoc","fields":[{"name":"id","type":"string"},{"name":"user_friends_count","type":["int","null"]},{"name":"user_location","type":["stri
ng","null"]},{"name":"user_description","type":["string","null"]},{"name":"user_statuses_count","type":["int","null"]},{"name":"user_followers_count","type":["int","null"]}
,{"name":"user_name","type":["string","null"]},{"name":"user_screen_name","type":["string","null"]},{"name":"created_at","type":["string","null"]},{"name":"text","type":["s
tring","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_status_id","type":["long","null"]},{"name":"media_url_https","type":["string","null"]},{"name":"expanded_url","type"
:["string","null"]}]}}▓y="Bp▓▒C▓Y◊y&1448443990005665792▓
,Santa Catarina, Brasil▓Los dias que tu juegas son todo lo que soy | https://www.instagram.com/aloize_/ RMCF ♛▓
Aloize NetoAloizeNeto(2021-10-14T00:22:12Z@Vontade de dar uma volta por ai▓<a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a>&14484439900015
28840▓(Dois Irmãos, Brasil▓Meu sangue sempre foi e sempre será azul.▒▒Lucas ₁₁lsgallert(2021-10-14T00:22:12Z▓Senhor eu só peço um atacante que finalize com força pq esses
cara parece que tem cupim nas pernas▓<a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a>&1448443989984747522▓
Piraju, Brasil▓breca não, deixa vim // @corinthians ♡▓▒▓mazinhamarissinhaaaa(2021-10-14T00:22:12Z▓eu odeio ter q ver jogo no futemax e essas muie na tela falando q tá a 5
km d mim▓<a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a>&1448443990014144513\                    <FF1122com(2021-10-14T00:22:12Z▓#

a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a>&1448443990014058499▓#▓V▓@amigo solo gente mas de 22 añosfuegopo(2021-1>▓@
0-14T00:22:12ZXRT @TikTokMenn: ♧ https://t.co/gq1sul1Bco▓<a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a>https://pbs.twimg.com/ext_tw_vi
deo_thumb/1447771340166078474/pu/img/swiB8Gi8K7IJO217.jpg▓https://twitter.com/TikTokMenn/status/1447771558135705604/video/1&1448443989988888579▓2Leonina
_bogado(2021-10-14T00:22:12Z▓RT @GenesisMora90: la ventaja de estar conmigo es que ando caliente 24/7▓<a href="http://twitter.com/download/android" rel="nofollow">Twitter f
or Android</a>&1448443989993000962▓VRChatにどハマり。 cyber系な感じ Design好き好きマン 駆け出し初心者3 Dモデラー VRCID anima1020 リツイートいっぱいします。 ▓▓anima ▓ani
ma102014(2021-10-14T00:22:12Z▓RT @RIORAO: 絶望の場面(暦廃墟)  https://t.co/ZL6L7RyaBW▓<a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a>^ht
tps://pbs.twimg.com/media/FBmbzpTUUAAu8H1.jpgzhttps://twitter.com/RIORAO/status/1448360749747748865/photo/1&1448443990005587971▓思った事をつぶやきます。
無言フォローすみません。▓{▓
              ⋅わ gzMT4wftKtrhtXo(2021-10-14T00:22:12Z▓RT @I5McNNdZnCXsI5u: 息子本人の希望により、昨日から午前中のみ別室登校始めました。母は祈るのみ、朝の会だけ
教室の端で参加。長期休んでいたのにオンライン繋がっていたからスッと溶け込めた。先生に感謝。オンライン授業は子どもたちの学びの命綱って本当なんです。コロナ…▓<a href="http://tw
itter.com/download/android" rel="nofollow">Twitter for Android</a>&1448443989997228034▓
YQ770eT▓<a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a>^https://pbs.twimg.com/media/FBmiMW7VUAEDgSk.jpgzhttps://twitter.com/ndc_rm/status
/1448367166017794050/photo/1&1448443989984690183▓Gatlin, NE▓▓▒▒▒▒▓Tim of the Corntim5bags(2021-10-14T00:22:12Z▓@FinnertyUSA @newsmax @FairfieldU looks like we're not s
ending our brightest out to set the world ablaze ▓<a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a>◊1448443990014005249&)▓ ▓>akun ini d
iawasi oleh orang tua▓
kajofbrenkseg(2021-10-14T00:22:12Z▓@jevierasz love u moreee jeje <3▓▓▓▓'▓<a href="http://twitter.com/download/android" rel="nofollow">Twitter for Android</a>▓▓(&144
8443990005612547▓柏白杵山踏の大阪(嬢つけ) ▓74式次男坊でチタン大好き。
「日々是精進・森羅万象に多情多恨」がモットーです。
自転車 釣り プラモ＆ガレージキット 等オタクネタメインにつぶやきます。▓
                                                      ★ﾅｸﾞﾀｾﾄｷiliketitanium(2021-10-14T00:22:12Z▓@dh3to5 ふふふふの 布団を 丁寧に 前脚でこねてぺったんこにして
くれる職人▓{▓<a href="http://twitter.com/download/android" rel="nofollow">Twitter for Android</a>▓▓▓(&1448443989988806656▓ ▓♥IG, tiktok,youtube ♂Devin Millan
auto follow back ko ♥♥♥
cat: Unable to write to output stream.
```

*** END***