



BIG DATA HADOOP & SPARK TRAINING

Assignment on Oozie and Flume



APRIL 8, 2018
RASHMI KRISHNA

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

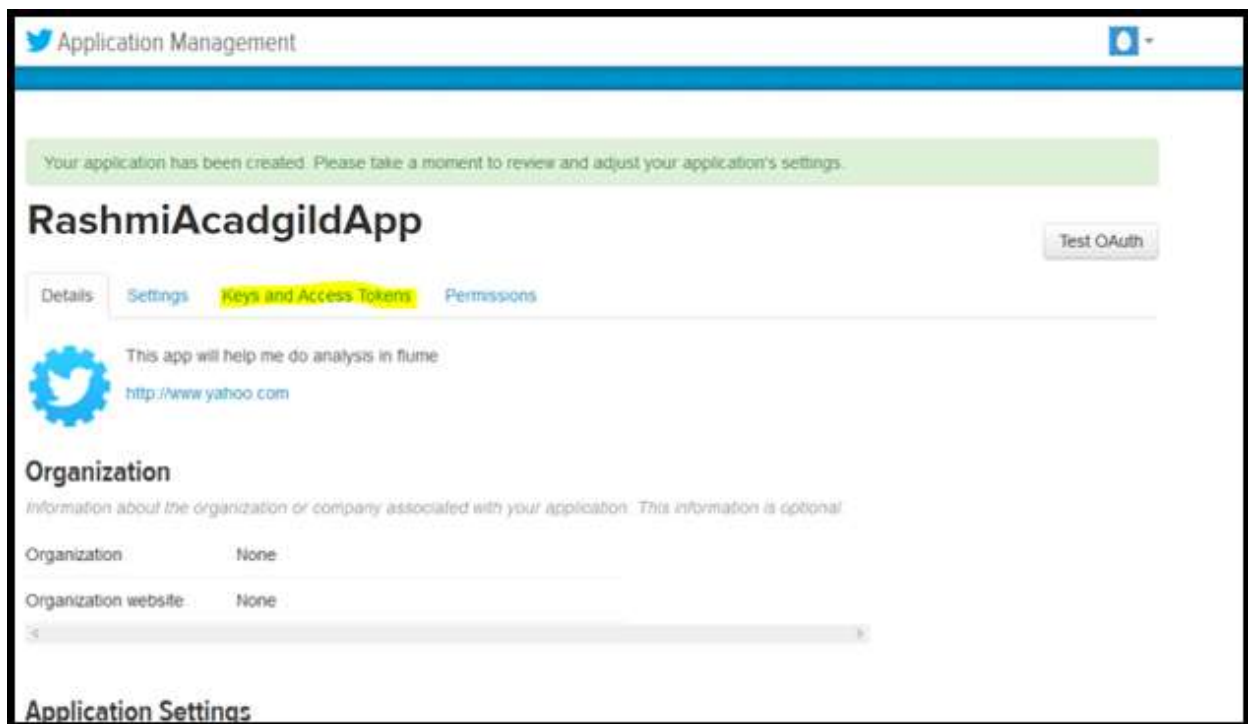
To create a flume agent, to stream the data from Twitter, first go to your twitter account and create an application from the below link:

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



Provide all the required credentials to create a app.

On successful credential authorization, you will get your application screen, from which go to "Keys and Access Tokens" tab.



In this screen copy the required information, i.e. copy consumer key, consumer secret, token and token secret

Secure <https://apps.twitter.com/app/15057038/keys>

Acadgild Retail Bat: gethue google A Simple Explanation GitHub - anjjava16 Free-Hadoop-Books Aviation Data Analy: Hive Date Functions facebook-hive-udfs

RashmiAcadgildApp

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.

Consumer Key (API Key)	y5hO1Ya4LnmWUmtA3Pr1R2bY0
Consumer Secret (API Secret)	mFTQy1cKQmdeSGxHOSc9vupOWSd1mfGp9xDOrM92GCD9WHdph

Access Level	Read and write (modify app permissions)
Owner	k_rashmi11
Owner ID	322711239

Application Actions

Regenerate Consumer Key and Secret Change App Permissions

Secure <https://apps.twitter.com/app/15057038/keys>

Acadgild Retail Bat: gethue google A Simple Explanation GitHub - anjjava16 Free-Hadoop-Books Aviation Data Analy: Hive Date Functions

Regenerate Consumer Key and Secret Change App Permissions

Your Access Token

This access token can be used to make API requests on your own account's behalf. Do not share your access token secret with anyone.

Access Token	322711239- h9wo8AnX9AKMzgcEWSMOaEJ91vXH6H0Ub6qyHhb
Access Token Secret	DyWtIf3fJTJ1UifLBP0QVpk3dZIMmTBjsRJdz2sAynhA9u

Access Level	Read and write
Owner	k_rashmi11
Owner ID	322711239

Use the information in the flume.conf_twitter file

flume.conf_twitter - Notepad

File Edit Format View Help

```
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=y5h01Ya4LnmlWUmtA3Pr1R2bY0
TwitterAgent.sources.Twitter.consumerSecret=mFT0fy1cK0mdeSGxH0So9vupOWSd1mfGp9xD0rM92GCD9WHdph
TwitterAgent.sources.Twitter.accessToken=322711239-h9wo8AnX9AKMzgcEWSMOaEJ91vXH6H01Jb6qyHhb
TwitterAgent.sources.Twitter.accessTokenSecret=DyWt1f3fJTJ1UtFLBP0QVpK3dZMmTBjsRJdz2sAynhA9u
TwitterAgent.sources.Twitter.keywords=hadoop, bigdata, mapreduce, mahout, 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:8020/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
```

Consumerkey,consumersecret, Token and Tokensecret access keys

HDFS output file path

Stream the Twitter data using flume using the below command:

flume-ng agent --conf c -n TwitterAgent -f /home/acadgild/flume/flume.conf_twitter


```

[acagild@localhost ~]$ flume-ng agent --conf c -n TwitterAgent -f /home/acagild/flume/flume.conf twitter
Info: Including Hadoop libraries found via (/home/acagild/install/hadoop/hadoop-2.6.5/bin/hadoop) for HDFS access
Info: Including HBASE libraries found via (/home/acagild/install/hbase/hbase-1.2.6/bin/hbase) for HBASE access
Info: Including Hive libraries found via (/home/acagild/install/hive/apache-hive-2.3.2-bin) for Hive access
+ exec /usr/java/jdk1.8.0_151/bin/java -Xmx20m -cp 'c:/home/acagild/install/flume/apache-flume-1.8.0-bin/lib/*:/home/acagild/install/hadoop/hadoop-2.6.5/contrib/capacity-scheduler/*:/home/acagild/install/hadoop/hadoop-2.6.5/etc/hadoop/*:/home/acagild/install/hadoop/hadoop-2.6.5/share/hadoop/common/lib/*:/home/acagild/install/hadoop/hadoop-2.6.5/share/hadoop/common/*:/home/acagild/install/hadoop/hadoop-2.6.5/share/hadoop/hdfs/lib/*:/home/acagild/install/hadoop/hadoop-2.6.5/share/hadoop/hdfs/*:/home/acagild/install/hadoop/hadoop-2.6.5/share/hadoop/yarn/lib/*:/home/acagild/install/hadoop/hadoop-2.6.5/share/hadoop/yarn/*:/home/acagild/install/hadoop/hadoop-2.6.5/share/hadoop/mapreduce/lib/*:/home/acagild/install/hadoop/hadoop-2.6.5/share/hadoop/mapreduce/*:/home/acagild/install/hbase/hbase-1.2.6/conf:/usr/java/jdk1.8.0_151/lib/tools.jar:/home/acagild/install/hbase/hbase-1.2.6:/home/acagild/install/hbase/hbase-1.2.6/lib/activation-1.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/aopalliance-1.0.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/apacheds-i18n-2.0.0-M15.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/apacheds-kerberos-codec-2.0.0-M15.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/api-asn1-api-1.0.0-M20.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/api-util-1.0.0-M20.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/asm-3.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/avro-1.7.4.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-beanutils-1.7.0.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-beanutils-core-1.8.0.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-cli-1.2.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-codec-1.9.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-collections-3.2.2.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-compress-1.4.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-configuration-1.6.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-daemon-1.0.13.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-digester-1.8.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-el-1.0.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-httpclient-3.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-io-2.4.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-lang-2.6.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-logging-1.2.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-math-2.2.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-math3-3.1.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/commons-net-3.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/disruptor-3.3.0.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/findbugs-annotations-1.3.9-1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/guava-12.0.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/guice-3.0.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/guice-servlet-3.0.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-annotations-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-auth-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-client-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-common-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-hdfs-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-mapreduce-client-app-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-mapreduce-client-common-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-mapreduce-client-core-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-mapreduce-client-jobclient-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-mapreduce-client-shuffle-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-yarn-api-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-yarn-client-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-yarn-common-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hadoop-yarn-server-common-2.5.1.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hbase-annotations-1.2.6.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hbase-annotations-1.2.6-tests.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hbase-client-1.2.6.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hbase-common-1.2.6.jar:/home/acagild/install/hbase/hbase-1.2.6/lib/hbase-common-1.2.6-tests.jar:/h

```

```

18/04/08 16:35:28 INFO node.PollingPropertiesFileConfigurationProvider: Configuration provider starting
18/04/08 16:35:28 INFO node.PollingPropertiesFileConfigurationProvider: Reloading configuration file:/home/acagild/flume/flume.conf_twitter
18/04/08 16:35:28 INFO conf.FlumeConfiguration: Processing:HDFS
18/04/08 16:35:28 INFO conf.FlumeConfiguration: Processing:HDFS
18/04/08 16:35:28 INFO conf.FlumeConfiguration: Processing:HDFS
18/04/08 16:35:28 INFO conf.FlumeConfiguration: Processing:HDFS
18/04/08 16:35:28 INFO conf.FlumeConfiguration: Processing:HDFS
18/04/08 16:35:28 INFO conf.FlumeConfiguration: Added sinks: HDFS Agent: TwitterAgent
18/04/08 16:35:28 INFO conf.FlumeConfiguration: Processing:HDFS
18/04/08 16:35:28 INFO conf.FlumeConfiguration: Processing:HDFS
18/04/08 16:35:28 INFO conf.FlumeConfiguration: Post-validation flume configuration contains configuration for agents: [TwitterAgent]
18/04/08 16:35:28 INFO node.AbstractConfigurationProvider: Creating channels
18/04/08 16:35:28 INFO channel.DefaultChannelFactory: Creating instance of channel MemChannel type memory
18/04/08 16:35:28 INFO node.AbstractConfigurationProvider: Created channel MemChannel
18/04/08 16:35:28 INFO source.DefaultSourceFactory: Creating instance of source Twitter, type org.apache.flume.source.twitter.TwitterSource
18/04/08 16:35:29 INFO sink.DefaultSinkFactory: Creating instance of sink: HDFS, type: hdfs
18/04/08 16:35:29 INFO node.AbstractConfigurationProvider: Channel MemChannel connected to [Twitter, HDFS]
18/04/08 16:35:29 INFO node.Application: Starting new configuration:{ sourceRunners:{Twitter=EventDrivenSourceRunner: { source:org.apache.flume.source.twitter.TwitterSource{name:Twitter,state:IDLE} }} sinkRunners:{HDFS=SinkRunner: { policy:org.apache.flume.sink.DefaultSinkProcessor@1e62b4a counterGroup:{ name:null counters:{ } }} channels:{MemChannel=org.apache.flume.channel.MemoryChannel{name:MemChannel}}} }
18/04/08 16:35:29 INFO node.Application: Starting Channel MemChannel
18/04/08 16:35:29 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: CHANNEL, name: MemChannel: Successfully registered new MBean.
18/04/08 16:35:29 INFO instrumentation.MonitoredCounterGroup: Component type: CHANNEL, name: MemChannel started
18/04/08 16:35:29 INFO node.Application: Starting Sink HDFS
18/04/08 16:35:29 INFO node.Application: Starting Source Twitter
18/04/08 16:35:29 INFO twitter.TwitterSource: Starting twitter source org.apache.flume.source.twitter.TwitterSource{name:Twitter,state:IDLE} ...
18/04/08 16:35:29 INFO twitter.TwitterSource: Twitter source Twitter started.
18/04/08 16:35:29 INFO twitter4j.TwitterStreamImpl: Establishing connection.
18/04/08 16:35:29 INFO instrumentation.MonitoredCounterGroup: Monitored counter group for type: SINK, name: HDFS: Successfully registered new MBean.
18/04/08 16:35:29 INFO instrumentation.MonitoredCounterGroup: Component type: SINK, name: HDFS started

```

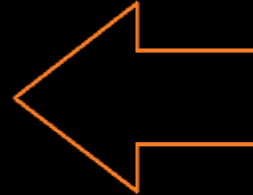


```

18/04/08 16:35:39 INFO twitter.TwitterSource: Processed 200 docs
18/04/08 16:35:41 INFO twitter.TwitterSource: Processed 300 docs
18/04/08 16:35:44 INFO twitter.TwitterSource: Processed 400 docs
18/04/08 16:35:47 INFO twitter.TwitterSource: Processed 500 docs
18/04/08 16:35:50 INFO twitter.TwitterSource: Processed 600 docs
18/04/08 16:35:52 INFO twitter.TwitterSource: Processed 700 docs
18/04/08 16:35:55 INFO twitter.TwitterSource: Processed 800 docs
18/04/08 16:35:58 INFO twitter.TwitterSource: Processed 900 docs
18/04/08 16:36:01 INFO twitter.TwitterSource: Processed 1,000 docs
18/04/08 16:36:01 INFO twitter.TwitterSource: Total docs indexed: 1,000, total skipped docs: 0
18/04/08 16:36:01 INFO twitter.TwitterSource: 32 docs/second
18/04/08 16:36:01 INFO twitter.TwitterSource: Run took 31 seconds and processed:
18/04/08 16:36:01 INFO twitter.TwitterSource: 0.009 MB/sec sent to index
18/04/08 16:36:01 INFO twitter.TwitterSource: 0.27 MB text sent to index
18/04/08 16:36:01 INFO twitter.TwitterSource: There were 0 exceptions ignored:
18/04/08 16:36:03 INFO twitter.TwitterSource: Processed 1,100 docs
18/04/08 16:36:06 INFO twitter.TwitterSource: Processed 1,200 docs
18/04/08 16:36:08 INFO twitter.TwitterSource: Processed 1,300 docs
18/04/08 16:36:11 INFO twitter.TwitterSource: Processed 1,400 docs
18/04/08 16:36:14 INFO twitter.TwitterSource: Processed 1,500 docs
18/04/08 16:36:16 INFO twitter.TwitterSource: Processed 1,600 docs
18/04/08 16:36:18 INFO twitter.TwitterSource: Processed 1,700 docs
18/04/08 16:36:21 INFO twitter.TwitterSource: Processed 1,800 docs
18/04/08 16:36:23 INFO twitter.TwitterSource: Processed 1,900 docs
18/04/08 16:36:26 INFO twitter.TwitterSource: Processed 2,000 docs
18/04/08 16:36:26 INFO twitter.TwitterSource: Total docs indexed: 2,000, total skipped docs: 0
18/04/08 16:36:26 INFO twitter.TwitterSource: 35 docs/second
18/04/08 16:36:26 INFO twitter.TwitterSource: Run took 56 seconds and processed:
18/04/08 16:36:26 INFO twitter.TwitterSource: 0.01 MB/sec sent to index
18/04/08 16:36:26 INFO twitter.TwitterSource: 0.539 MB text sent to index
18/04/08 16:36:26 INFO twitter.TwitterSource: There were 0 exceptions ignored:
18/04/08 16:36:28 INFO twitter.TwitterSource: Processed 2,100 docs
18/04/08 16:36:31 INFO twitter.TwitterSource: Processed 2,200 docs
18/04/08 16:36:33 INFO twitter.TwitterSource: Processed 2,300 docs
18/04/08 16:36:36 INFO twitter.TwitterSource: Processed 2,400 docs
18/04/08 16:36:38 INFO twitter.TwitterSource: Processed 2,500 docs
18/04/08 16:36:41 INFO twitter.TwitterSource: Processed 2,600 docs
18/04/08 16:36:44 INFO twitter.TwitterSource: Processed 2,700 docs

```

Processing documents



```

18/04/08 16:54:07 INFO instrumentation.MonitoredCounterGroup: Shutdown Metric for type: SINK, name: HDFS. sink.start.time == 1523185329
663
18/04/08 16:54:07 INFO instrumentation.MonitoredCounterGroup: Shutdown Metric for type: SINK, name: HDFS. sink.stop.time == 15231866470
27
18/04/08 16:54:07 INFO instrumentation.MonitoredCounterGroup: Shutdown Metric for type: SINK, name: HDFS. sink.batch.complete == 0
18/04/08 16:54:07 INFO instrumentation.MonitoredCounterGroup: Shutdown Metric for type: SINK, name: HDFS. sink.batch.empty == 60
18/04/08 16:54:07 INFO instrumentation.MonitoredCounterGroup: Shutdown Metric for type: SINK, name: HDFS. sink.batch.underflow == 1
18/04/08 16:54:07 INFO instrumentation.MonitoredCounterGroup: Shutdown Metric for type: SINK, name: HDFS. sink.connection.closed.count
== 2
18/04/08 16:54:07 INFO instrumentation.MonitoredCounterGroup: Shutdown Metric for type: SINK, name: HDFS. sink.connection.creation.coun
t == 2
18/04/08 16:54:07 INFO instrumentation.MonitoredCounterGroup: Shutdown Metric for type: SINK, name: HDFS. sink.connection.failed.count
== 0
18/04/08 16:54:07 INFO instrumentation.MonitoredCounterGroup: Shutdown Metric for type: SINK, name: HDFS. sink.event.drain.attempt == 6
21
18/04/08 16:54:07 INFO instrumentation.MonitoredCounterGroup: Shutdown Metric for type: SINK, name: HDFS. sink.event.drain.sucess == 62
1
You have new mail in /var/spool/mail/acadgild
[acadgild@localhost ~]$ hadoop fs -ls /user/flume/tweets
18/04/08 16:54:41 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where
e applicable
Found 2 items
-rw-r--r-- 1 acadgild supergroup 11430392 2018-04-08 16:45 /user/flume/tweets/FlumeData.152318534093
-rw-r--r-- 1 acadgild supergroup 881617 2018-04-08 16:54 /user/flume/tweets/FlumeData.1523186140569
[acadgild@localhost ~]$ hadoop fs -cat /user/flume/tweets/FlumeData.152318534093

```

Tweets retrieved from twitter, saved in HDFS is read using the cat command:

[illegible]