**ORIEN IT**

*Mr.Kalyan, Apache Contributor, Cloudera CCA175 Certified Consultant,*
*6+ years of Big Data exp, IIT Kharagpur, Gold Medalist*

# Kalyan Big Data Projects – Project 6
# How To Stream CSV Data Into Hadoop Using
# Apache Flume  - Kafka Channel

**Pre-Requisites of Flume Project:**

hadoop-2.6.0
flume-1.6.0
kafka-0.9.0
java-1.7

**NOTE:** Make sure that install all the above components

**Flume Project Download Links:**

`hadoop-2.6.0.tar.gz` ==> link
(https://archive.apache.org/dist/hadoop/core/hadoop-2.6.0/hadoop-2.6.0.tar.gz)

`apache-flume-1.6.0-bin.tar.gz` ==> link
(https://archive.apache.org/dist/flume/1.6.0/apache-flume-1.6.0-bin.tar.gz)

`kafka_2.11-0.9.0.0.tgz` ==> link
(https://archive.apache.org/dist/kafka/0.9.0.0/kafka_2.11-0.9.0.0.tgz)

`kalyan-bigdata-examples.jar` ==> link
(https://github.com/kalyanhadooptraining/kalyan-bigdata-realtime-
projects/blob/master/kalyan/kalyan-bigdata-examples.jar)

`kalyan-kafka-channel-agent.conf` ==> link
(https://github.com/kalyanhadooptraining/kalyan-bigdata-realtime-
projects/blob/master/kafka/project6-flume-kafka-channel/kalyan-kafka-channel-agent.conf)

--------------------------------------------------------------------------------------------------------------------
**Learnings of this Project:**
--------------------------------------------------------------------------------------------------------------------
➢ We will learn Flume Configurations and Commands
➢ Flume Agent
    1.  Source (Exec Source)
    2.  Channel (Kafka Channel)
    3.  Sink (Hdfs Sink)
➢ We will learn Kafka Configurations and Commands
➢ Kafka Information
    1.  Kalyan Util (CSV data generator)

**ORIEN IT**

*Mr.Kalyan, Apache Contributor, Cloudera CCA175 Certified Consultant,*
*6+ years of Big Data exp, IIT Kharagpur, Gold Medalist*

2. Kafka Producer (Listen on CSV data)
3. Kafka Consumer (Recieves the data from Kafka Producer)
4. Flume Kafka Channel (Will Recieves the Kafka Channel data from Flume Source)

➢ Major project in Real Time `Product Log Analysis`
1. We are extracting the data from server logs
2. This data will be useful to do analysis on product views
3. CSV is the output format

➢ We can use hive / pig / mapreduce to analyze this data

1. explore hive query to analysis
2. explore pig scripts to analysis
3. explore mapreduce to analysis

--------------------------------------------------------------------------------------------------------------------

1. create "**kalyan-kafka-channel-agent.conf**" file with below content

```
agent.sources = EXEC
agent.channels = KAFKA
agent.sinks = HDFS

agent.sources.EXEC.type = exec
agent.sources.EXEC.command = tail -F /tmp/users.csv
agent.sources.EXEC.channels = KAFKA

agent.sinks.HDFS.type = hdfs
agent.sinks.HDFS.channel = KAFKA
agent.sinks.HDFS.hdfs.path = hdfs://localhost:8020/user/kafka/messages
agent.sinks.HDFS.hdfs.fileType = DataStream
agent.sinks.HDFS.hdfs.writeFormat = Text
agent.sinks.HDFS.hdfs.batchSize = 10
agent.sinks.HDFS.hdfs.rollSize = 0
agent.sinks.HDFS.hdfs.rollCount = 10
agent.sinks.HDFS.hdfs.useLocalTimeStamp = true

agent.channels.KAFKA.type = org.apache.flume.channel.kafka.KafkaChannel
agent.channels.KAFKA.brokerList = localhost:9092
agent.channels.KAFKA.zookeeperConnect = localhost:2181
agent.channels.KAFKA.kafka.consumer.timeout.ms = 100
```

2. Copy "**kalyan-kafka-channel-agent.conf**" file into "**$FUME_HOME/conf**" folder
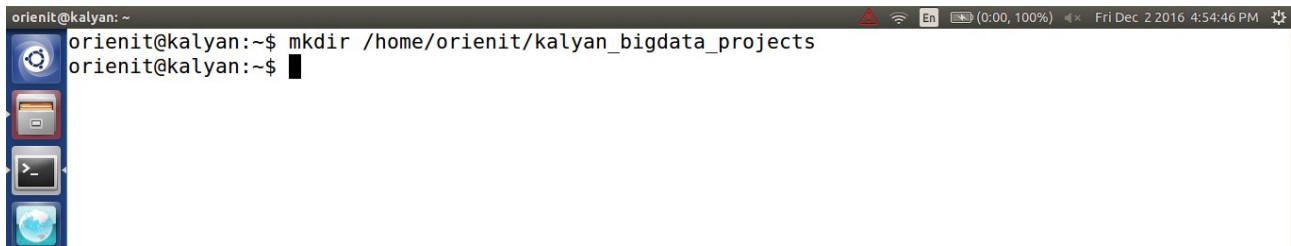
3. Generate Large Amount of Sample CSV data follow this article.

(http://kalyanbigdatatraining.blogspot.com/2016/12/how-to-generate-large-amount-of-sample.html)

4. Follow below steps...

# ORIEN IT

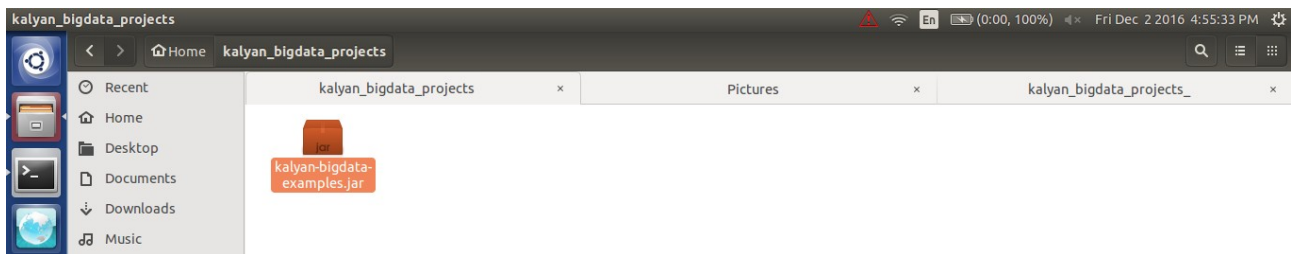*Mr.Kalyan, Apache Contributor, Cloudera CCA175 Certified Consultant,*
*6+ years of Big Data exp, IIT Kharagpur, Gold Medalist*

i) Create '**kalyan_bigdata_projects**' folder in **user home** (i.e **/home/orienit**)
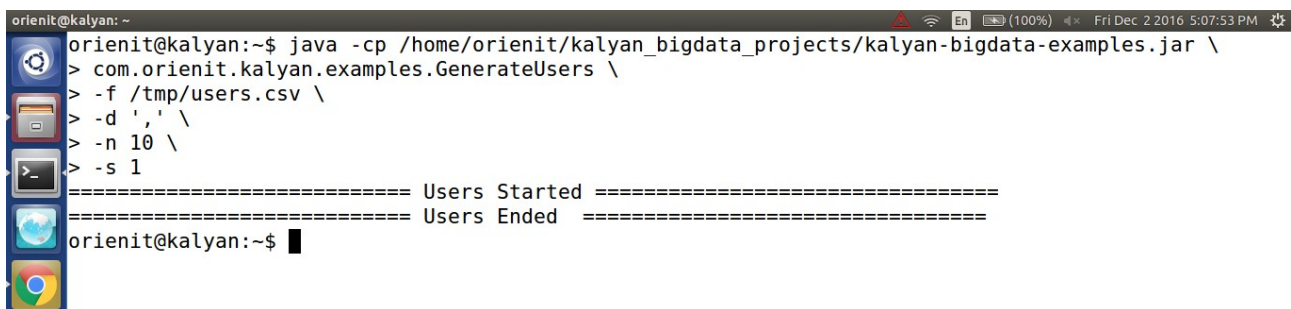
**Command:**  *mkdir /home/orienit/kalyan_bigdata_projects*



ii)  Copy '**kalyan-bigdata-examples.jar**' jar file into '**/home/orienit/kalyan_bigdata_projects**' folder



iii)  Execute Below Command to Generate Sample CSV data with 100 lines. Increase this number to get more data ...

java -cp /home/orienit/kalyan_bigdata_projects/kalyan-bigdata-examples.jar \
com.orienit.kalyan.examples.GenerateUsers \
-f /tmp/users.csv \
-d ',' \
-n 10 \
-s 1



5. Verify the Sample CSV data in Console, using below command

cat /tmp/users.csv

ORIEN IT

*Mr.Kalyan, Apache Contributor, Cloudera CCA175 Certified Consultant,*
*6+ years of Big Data exp, IIT Kharagpur, Gold Medalist*

```
orienit@kalyan:~$ cat /tmp/users.csv
1,user1,user1,user1@gmail.com,US,Washington,Seattle,2016-07-02 05:07:48
2,user2,user2,user2@gmail.com,US,Florida,Orlando,2016-07-02 05:07:48
3,user3,user3,user3@gmail.com,US,New York,Little Falls,2016-07-02 05:07:49
4,user4,user4,user4@gmail.com,India,Karnataka,Mangaluru,2016-07-02 05:07:49
5,user5,user5,user5@gmail.com,US,Hawaii,Hanapepe,2016-07-02 05:07:49
6,user6,user6,user6@gmail.com,India,Chennai,Kottur,2016-07-02 05:07:49
7,user7,user7,user7@gmail.com,India,Andhra Pradesh,Kakinada,2016-07-02 05:07:49
8,user8,user8,user8@gmail.com,US,Hawaii,East Honolulu,2016-07-02 05:07:49
9,user9,user9,user9@gmail.com,US,Florida,Hollywood,2016-07-02 05:07:49
10,user10,user10,user10@gmail.com,US,Washington,Bellevue,2016-07-02 05:07:49
orienit@kalyan:~$
```

6. Start the `zookeeper` using below command (New Terminal)

$KAFKA_HOME/bin/zookeeper-server-start.sh $KAFKA_HOME/config/zookeeper.properties

```
orienit@kalyan:~$ $KAFKA_HOME/bin/zookeeper-server-start.sh $KAFKA_HOME/config/zookeeper.properties
[2017-01-03 16:41:44,986] INFO Reading configuration from: /home/orienit/work/kafka_2.11-0.9.0.0/co
nfig/zookeeper.properties (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2017-01-03 16:41:45,024] INFO autopurge.snapRetainCount set to 3 (org.apache.zookeeper.server.Data
dirCleanupManager)
[2017-01-03 16:41:45,024] INFO autopurge.purgeInterval set to 0 (org.apache.zookeeper.server.Datadi
rCleanupManager)
[2017-01-03 16:41:45,024] INFO Purge task is not scheduled. (org.apache.zookeeper.server.DatadirCle
anupManager)
[2017-01-03 16:41:45,025] WARN Either no config or no quorum defined in config, running  in standal
one mode (org.apache.zookeeper.server.quorum.QuorumPeerMain)
[2017-01-03 16:41:45,086] INFO Reading configuration from: /home/orienit/work/kafka_2.11-0.9.0.0/co
nfig/zookeeper.properties (org.apache.zookeeper.server.quorum.QuorumPeerConfig)
[2017-01-03 16:41:45,086] INFO Starting server (org.apache.zookeeper.server.ZooKeeperServerMain)
[2017-01-03 16:41:45,117] INFO Server environment:zookeeper.version=3.4.6-1569965, built on 02/20/2
014 09:09 GMT (org.apache.zookeeper.server.ZooKeeperServer)
```

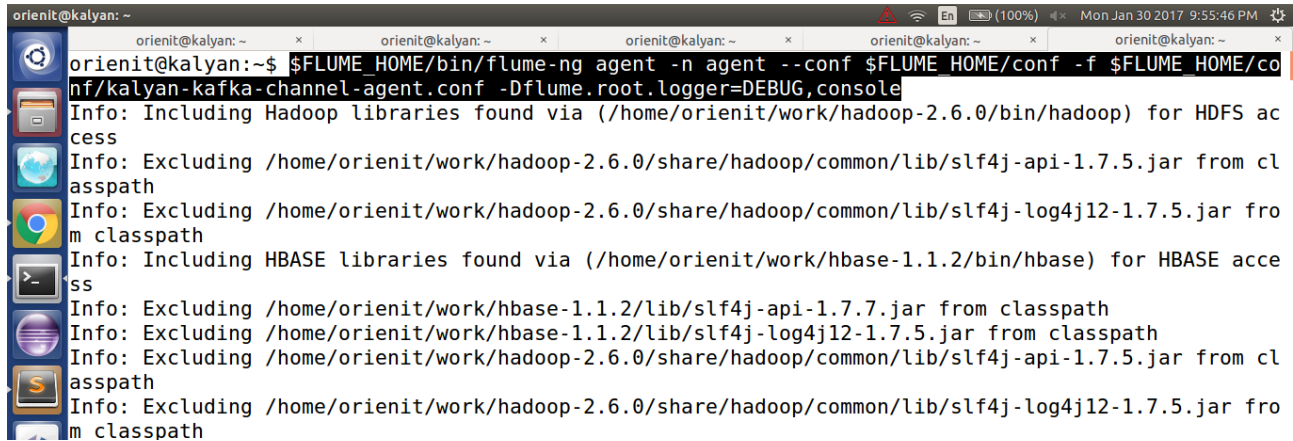7. Start the `kafka server` using below command (New Terminal)

$KAFKA_HOME/bin/kafka-server-start.sh $KAFKA_HOME/config/server.properties

```
orienit@kalyan:~$ $KAFKA_HOME/bin/kafka-server-start.sh $KAFKA_HOME/config/server.properties
[2017-01-03 16:43:14,430] INFO KafkaConfig values:
        advertised.host.name = null
        metric.reporters = []
        quota.producer.default = 9223372036854775807
        offsets.topic.num.partitions = 50
        log.flush.interval.messages = 9223372036854775807
        auto.create.topics.enable = true
        controller.socket.timeout.ms = 30000
        log.flush.interval.ms = null
        principal.builder.class = class org.apache.kafka.common.security.auth.DefaultPrincipalBuild
er
        replica.socket.receive.buffer.bytes = 65536
        min.insync.replicas = 1
        replica.fetch.wait.max.ms = 500
        num.recovery.threads.per.data.dir = 1
        ssl.keystore.type = JKS
        default.replication.factor = 1
        ssl.truststore.password = null
        log.preallocate = false
```

# ORIEN IT

*Mr.Kalyan, Apache Contributor, Cloudera CCA175 Certified Consultant,*
*6+ years of Big Data exp, IIT Kharagpur, Gold Medalist*

8. Execute the below command to `**Extract data from CSV file using KAFKA Channel**`
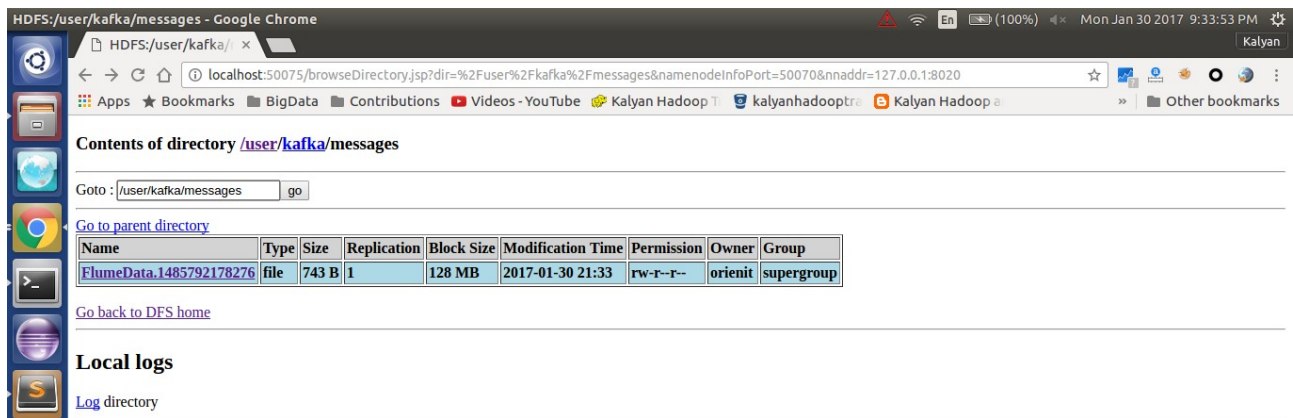
$FLUME_HOME/bin/flume-ng agent -n agent --conf $FLUME_HOME/conf -f
$FLUME_HOME/conf/kalyan-kafka-channel-agent.conf -Dflume.root.logger=DEBUG,console



9. Verify the data in hdfs location is "**hdfs://localhost:8020/user/kafka/messages**"