# Big Data Hadoop and Spark Developer
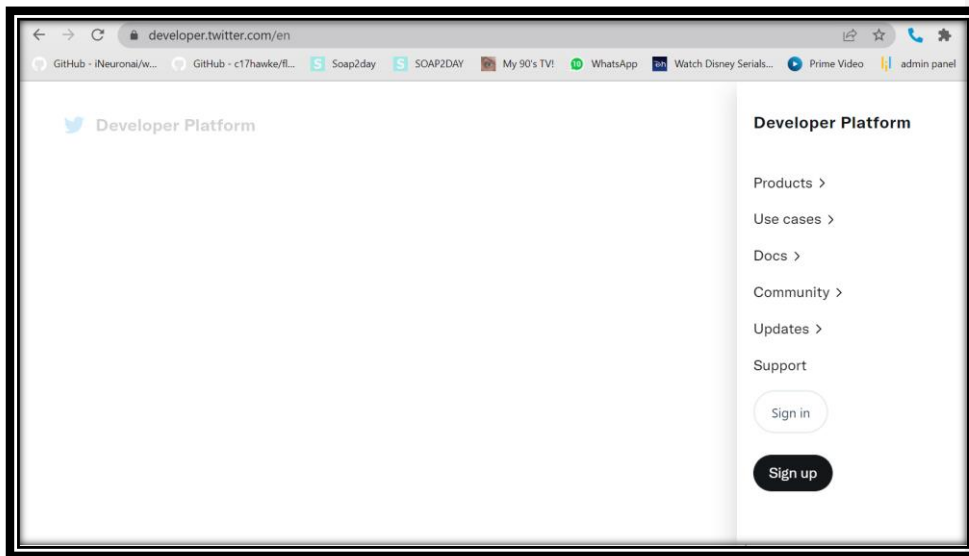
**simplilearn**

Get Certified. Get Ahead.

# Twitter Data Ingestion with Flume

**Steps to Perform:**

**Step 1:** Set up the Twitter account

1.2.1   Sign up for a developer account from the link mentioned:

https://developer.twitter.com/en



1.2   Create a project and connect an App

1.2.1 In the developer portal, click on **Create a "New Project"**

1.2.2   Provide a suitable name for the project

1.2.3   Select an appropriate use-case

1.2.4   Provide a project description

1.2.5   Create a new App or connect to an existing App

**Note:** An App is a container for API Keys that needs to make an HTTP request to the Twitter API.

1.2.6 Click on '**Create a new App instead'** and give your App a unique name

**Note:** Once you click complete, you will get your API Keys and the Bearer Token that can be used to connect to the new endpoints in the Twitter API v2.



1.2.7 Click the **(+)** next to API Key, API Secret Key, and Bearer Token and copy it in a safe place on your local machine

**Note 1:** You will need these to make the API calls in the next step.

**Note 2:** The keys in the screenshot above are hidden. However, in your own developer portal, you will be able to see the actual values for the API Key, API Secret Key, and Bearer Token.

**Step 2:** Set up the Flume configuration file on the "**Webconsole"**

2.1   Open the course "**Big Data Hadoop and Spark Developer"**

2.2   Click on the "**PRACTICE LABS"** tab on the left side and select "**LAUNCH LAB"**



**2.3** Click on the **"LAUNCH LAB"** button

**2.4** Log in to the **"Webconsole"** lab



**2.5** Run the command below to create a configuration file:

**Command:**

vi flume.conf

```
[testdemomay1301mailinator@bdh-cluster2-edgenode10 ~]$ vi flume.conf
[testdemomay1301mailinator@bdh-cluster2-edgenode10 ~]$ ▯
```
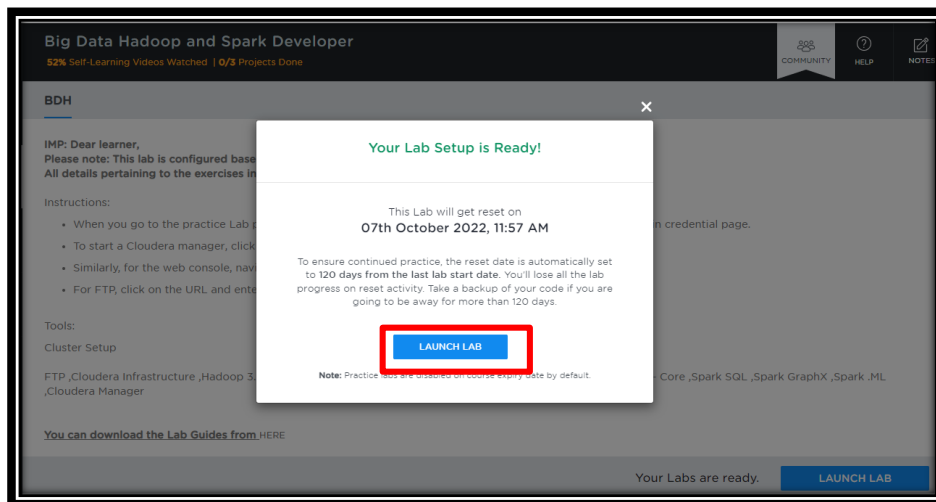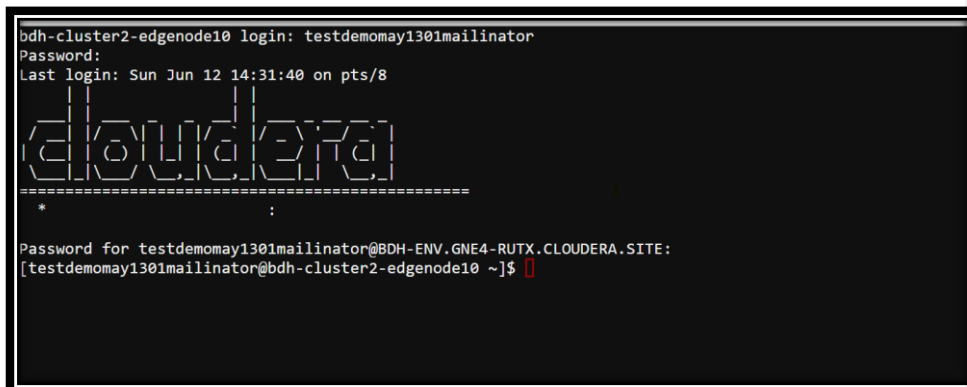
**2.6** Add the below configurations with your secret access tokens received from the twitter API:

# 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 = <your info>

TwitterAgent.sources.Twitter.consumerSecret = <your info>

TwitterAgent.sources.Twitter.accessToken = <your info>

TwitterAgent.sources.Twitter.accessTokenSecret = <your info>

**TwitterAgent.sources.Twitter.keywords = ISIS,Terrorist,global_threat**

# Describing/Configuring the sink

```
TwitterAgent.sinks.HDFS.type = hdfs

TwitterAgent.sinks.HDFS.hdfs.path = /user/testdemomay1301mailinator

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.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
```

```
# 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 = <your info>
TwitterAgent.sources.Twitter.consumerSecret = <your info>
TwitterAgent.sources.Twitter.accessToken = <your info>
TwitterAgent.sources.Twitter.accessTokenSecret = <your info>
TwitterAgent.sources.Twitter.keywords = ISIS,Terrorist,global_threat

# Describing/Configuring the sink
TwitterAgent.sinks.HDFS.type = hdfs
TwitterAgent.sinks.HDFS.hdfs.path = /user/testdemomay1301mailinator
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.capacity = 10000
TwitterAgent.channels.MemChannel.transactionCapacity = 100
```

**2.7** Run the file on Console using the command below and verify the output on HDFS:

flume-ng agent  -f /mnt/home/testdemomay1301mailinator/flume-demo/flume.conf -Dflume.root.logger=DEBUG,console -n TwitterAgent

| | Name | Size | User | Group | Permissions | Date |
|---|---|---|---|---|---|---|
| ▪ | ↥ | | singh25novgmail | hadoop | drwxrwx--- | April 04, 2022 12:16 PM |
| ▪ | . | | singh25novgmail | hadoop | drwxrwx--- | April 04, 2022 12:17 PM |
| 🗋 | FlumeData.1649099814826 | 1.2 KB | singh25novgmail | hadoop | -rw-r--r-- | April 04, 2022 12:16 PM |
| 🗋 | FlumeData.1649099814827 | 19.5 KB | singh25novgmail | hadoop | -rw-r--r-- | April 04, 2022 12:16 PM |
| 🗋 | FlumeData.1649099814828 | 19.5 KB | singh25novgmail | hadoop | -rw-r--r-- | April 04, 2022 12:16 PM |
| 🗋 | FlumeData.1649099814829 | 17.8 KB | singh25novgmail | hadoop | -rw-r--r-- | April 04, 2022 12:16 PM |
| 🗋 | FlumeData.1649099814830 | 17.2 KB | singh25novgmail | hadoop | -rw-r--r-- | April 04, 2022 12:16 PM |
| 🗋 | FlumeData.1649099814831 | 22.0 KB | singh25novgmail | hadoop | -rw-r--r-- | April 04, 2022 12:17 PM |
| 🗋 | FlumeData.1649099814832 | 16.4 KB | singh25novgmail | hadoop | -rw-r--r-- | April 04, 2022 12:17 PM |
| 🗋 | FlumeData.1649099814833 | 17.5 KB | singh25novgmail | hadoop | -rw-r--r-- | April 04, 2022 12:17 PM |