

CS 5433: Bigdata Management
Programming Assignment 1
Report for PART 1

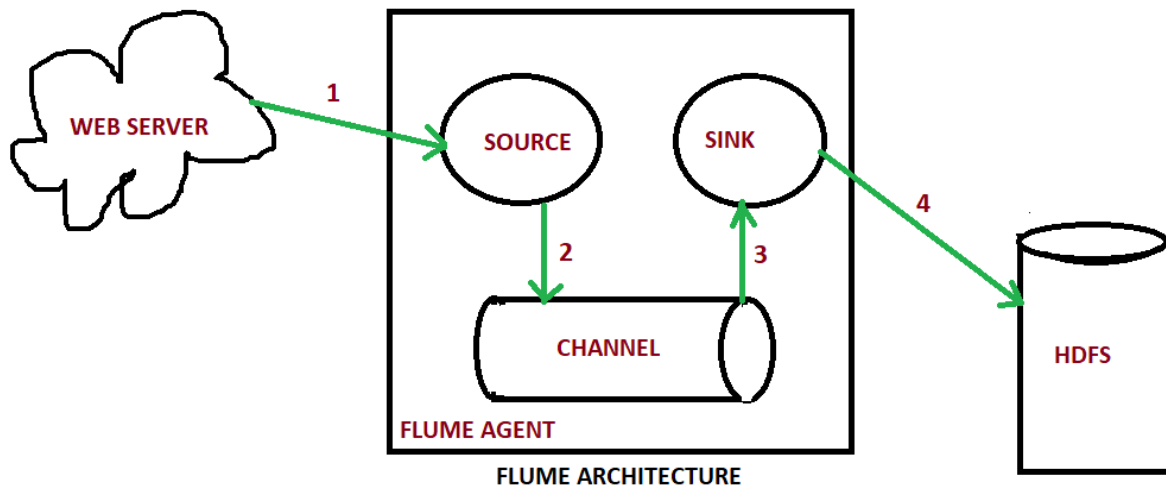
CWID: A20343337

Hadoop: It is also called as Apache Hadoop. It is an open-source framework which is used to store and process the large size of datasets (like size ranges from gigabytes to petabytes) efficiently. It mainly contains four modules,

- a) Hadoop File System (HDFS)
- b) Yet Another Resource Negotiator (YARN)
- c) MapReduce
- d) Hadoop Common

In PART1, we mainly use the HDFS which is a distributed file system that runs on low-end hardware. HDFS outperforms traditional file systems in terms of data performance, fault tolerance, and native support for huge datasets.

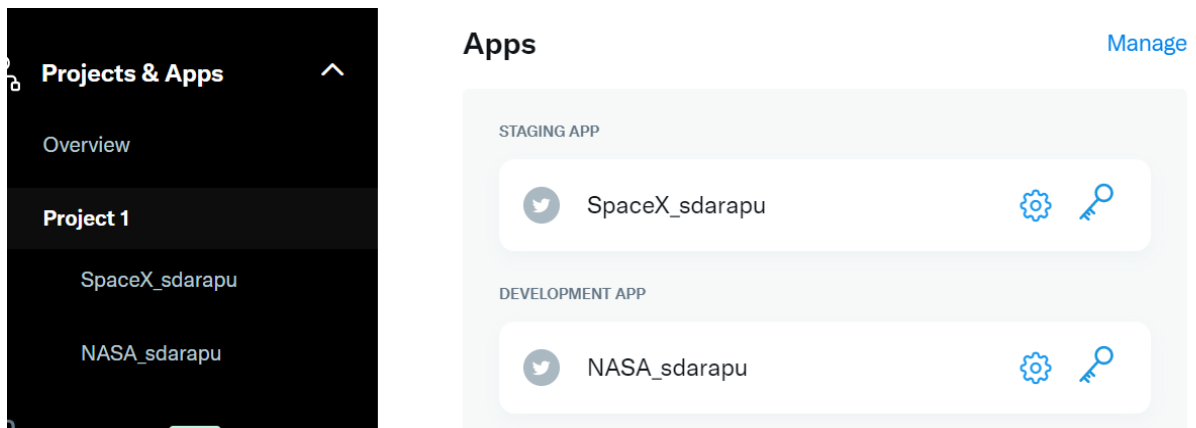
Flume: Flume is a distributed and reliable service for rapidly collecting, aggregating, and transporting massive amounts of log data. It is also called as Apache flume. It is an open-source software which is used to store the streamed data into HDFS. It has a simple and adaptable architecture based on the streamed data flows. Below is the architecture



Approach:

Below are the steps, I have used to complete the PART 1 of the assignment

1. To collect the data from twitter using flume, I have selected two keywords "NASA" and "SpaceX".
2. After selecting the keywords, I have filed a twitter application for developer account which is an elevated access. Within couple of days, my request for elevated access is approved.
3. Then I have created two Apps named "NASA_sdarapu" and "Spacex_sdarapu" to get my keys and tokens.



4. Now, I logged into the department's Hadoop cluster ("hadoop-nn001.cs.okstate.edu") by using my credentials.

OpenSSH SSH client

```
Microsoft Windows [Version 10.0.19042.1526]
(c) Microsoft Corporation. All rights reserved.

C:\Users\darap>ssh sdarapu@hadoop-nn001.cs.okstate.edu
sdarapu@hadoop-nn001.cs.okstate.edu's password:
Welcome to Ubuntu 20.04.3 LTS (GNU/Linux 5.4.0-91-generic x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:       https://ubuntu.com/advantage

System information as of Sat 26 Feb 2022 10:24:56 AM CST

System load:  0.06               Processes:            335
Usage of /:   3.6% of 1006.92GB   Users logged in:     3
Memory usage: 48%               IPv4 address for ens3: 192.168.122.2
Swap usage:   9%

Last login: Wed Feb 23 13:47:34 2022 from 10.200.192.236
sdarapu@hadoop-nn001:~$
```

5. After that, I have created two new directories “NASA_PA1data” and “SpaceX_PA1data” using the below command where the twitter data gets download for the keywords “NASA” and “SpaceX”.

For NASA_PA1data:

```
sdarapu@hadoop-nn001:~$ hdfs dfs -mkdir /user/sdarapu/NASA_PA1data
2022-02-26 10:26:11,640 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform...
sdarapu@hadoop-nn001:~$
```

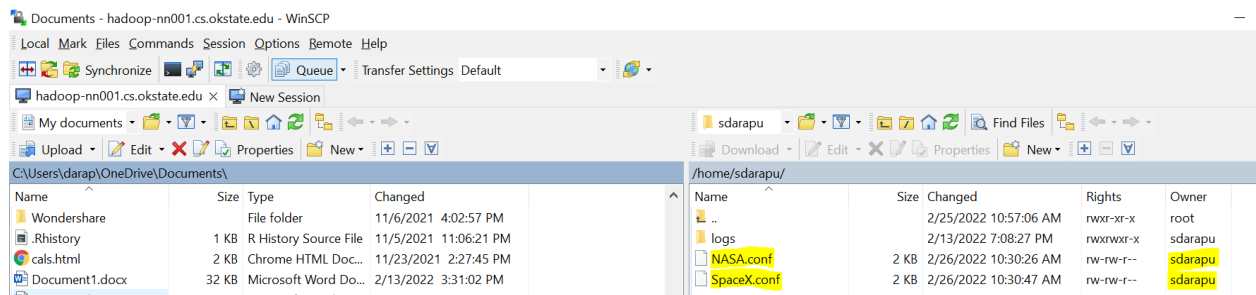
For SpaceX_PA1data:

```
sdarapu@hadoop-nn001:~$ hdfs dfs -mkdir /user/sdarapu/SpaceX_PA1data
2022-02-26 10:27:15,744 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform...
sdarapu@hadoop-nn001:~$
```

6. By using the “hdfs dfs -ls” command, I have checked whether the directories are created or not.

```
sdarapu@hadoop-nn001:~$ hdfs dfs -ls
2022-02-26 10:28:12,616 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform...
Found 2 items
drwxr-xr-x - sdarapu sdarapu      0 2022-02-26 10:26 NASA_PA1data
drwxr-xr-x - sdarapu sdarapu      0 2022-02-26 10:27 SpaceX_PA1data
sdarapu@hadoop-nn001:~$
```

7. Once directories are created, I have created two configuration files named “NASA.conf” and “SpaceX.conf” for connecting the twitter and flume agent directly in “WinSCP” instead of using “nano” command as shown below.



For NASA.conf:

```
/home/sdarapu/NASA.conf - hadoop-nn001.cs.okstate.edu - Editor - WinSCP
# Naming the components on the current agent
TwitterAgent.sources = Twitter
TwitterAgent.channels = MemChannel
TwitterAgent.sinks = HDFS

# Describing/Configuring the source
TwitterAgent.sources.Twitter.type = com.cloudera.flume.source.TwitterSource
TwitterAgent.sources.Twitter.consumerKey = 1ohmbRhAbTzjtZUmZlTerIhLK
TwitterAgent.sources.Twitter.consumerSecret = PgZ96EnGvf0aHK7M0uzQ15zCtwIIIn9sOxIAduHtnhKQDTKKnq
TwitterAgent.sources.Twitter.accessToken = 1491145396961304586-qZvpaB3VVkmZs7qvt1WF8NC4dd3diC
TwitterAgent.sources.Twitter.accessTokenSecret = zcsNfEqdcaPGEkI6IvyRpVcX0t5cNwRdjEjuWPKOR9vg4
TwitterAgent.sources.Twitter.keywords = NASA

# Describing/Configuring the sink
TwitterAgent.sinks.HDFS.type = hdfs
TwitterAgent.sinks.HDFS.hdfs.path = hdfs://hadoop-nn001.cs.okstate.edu:9000/user/sdarapu/NASA_PA1data/%Y/%m/%d/%H
TwitterAgent.sinks.HDFS.hdfs.useLocalTimeStamp = true
TwitterAgent.sinks.HDFS.hdfs.fileType = DataStream
TwitterAgent.sinks.HDFS.hdfs.writeFormat = Text
TwitterAgent.sinks.HDFS.hdfs.batchSize = 100
TwitterAgent.sinks.HDFS.hdfs.rollSize = 0
TwitterAgent.sinks.HDFS.hdfs.rollCount = 0

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

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

For SpaceX.conf:

```
/home/sdarapu/SpaceX.conf - hadoop-nn001.cs.okstate.edu - Editor - WinSCP
# Naming the components on the current agent
TwitterAgent.sources = Twitter
TwitterAgent.channels = MemChannel
TwitterAgent.sinks = HDFS

# Describing/Configuring the source
TwitterAgent.sources.Twitter.type = com.cloudera.flume.source.TwitterSource
TwitterAgent.sources.Twitter.consumerKey = 3RJWtUQCrasVyMKZboejqB3dC
TwitterAgent.sources.Twitter.consumerSecret = 6U4hyfBrf2gH26TXv0ims8GnQBh1kPsqabNlmsVj01Dr44a5Kf
TwitterAgent.sources.Twitter.accessToken = 1491145396961304586-dGcDqkJ3lTR5x33DFLywzXgCGQOmXJ
TwitterAgent.sources.Twitter.accessTokenSecret = WZMrcmZjU7g7QJ5cNh2D8dU2S9qb7FhZH7Q505g7IC6jI
TwitterAgent.sources.Twitter.keywords = SpaceX

# Describing/Configuring the sink
TwitterAgent.sinks.HDFS.type = hdfs
TwitterAgent.sinks.HDFS.hdfs.path = hdfs://hadoop-nn001.cs.okstate.edu:9000/user/sdarapu/SpaceX_PA1data/%Y/%m/%d/%H
TwitterAgent.sinks.HDFS.hdfs.useLocalTimeStamp = true
TwitterAgent.sinks.HDFS.hdfs.fileType = DataStream
TwitterAgent.sinks.HDFS.hdfs.writeFormat = Text
TwitterAgent.sinks.HDFS.hdfs.batchSize = 100
TwitterAgent.sinks.HDFS.hdfs.rollSize = 0
TwitterAgent.sinks.HDFS.hdfs.rollCount = 0

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

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

- Once conf files are created, I have run the “nohup” command to download the twitter data into the respective directories specified in the conf files. Also, note down the process id’s to kill once the enough data is downloaded later.

Like, data related to “NASA” keyword downloads into the directory “NASA_PA1data” and data related to “SpaceX” keyword downloads into the directory “SpaceX_PA1data”.

For NASA data:

```
sdarapu@hadoop-nn001:~$ nohup $FLUME_HOME/bin/flume-ng agent -n TwitterAgent -f /home/sdarapu/NASA.conf --conf /usr/local/flume/conf &
[1] 367095
sdarapu@hadoop-nn001:~$ nohup: ignoring input and appending output to 'nohup.out'
```

Process Id: 367095 → For nohup process of NASA data

For SpaceX data:

```
sdarapu@hadoop-nn001:~$ nohup $FLUME_HOME/bin/flume-ng agent -n TwitterAgent -f /home/sdarapu/SpaceX.conf --conf /usr/local/flume/conf &
[1] 369858
sdarapu@hadoop-nn001:~$ nohup: ignoring input and appending output to 'nohup.out'
```

Process Id: 369858 → For nohup process of SpaceX data

- After running the nohup command, I have stopped it after waiting some time to download the enough data and then I have killed the nohup process once the enough data is downloaded from twitter. To get the process id, we can use the “ps -x” command as shown below.

For NASA data:

```
sdarapu@hadoop-nn001:~$ ps -x
  PID TTY          STAT TIME  COMMAND
 365215 ?            Ss   0:00 /lib/systemd/systemd --user
 365218 ?            S    0:00 (sd-pam)
 365225 ?          Ss1   0:00 /usr/bin/pulseaudio --daemonize=no --log-target=journal
 365244 ?            Ss   0:00 /usr/bin/dbus-daemon --session --address=systemd: --nofork --nopidfile --systemd-activation --syslog-only
 365302 ?            S    0:00 sshd: sdarapu@notty
 365303 ?            Ss   0:00 /usr/lib/openssh/sftp-server
 365928 ?            R    0:00 sshd: sdarapu@pts/5
 365930 pts/5      Ss   0:00 -bash
 366625 ?            S    0:00 sshd: sdarapu@notty
 366626 ?            Ss   0:00 /usr/lib/openssh/sftp-server
 367095 pts/5      S1   0:38 /usr/bin/java -Xmx60m -Dtwitter4j.streamBase.URL=https://stream.twitter.com/1.1/ -cp /usr/local/flume/conf:
 368773 pts/5      R+   0:00 ps -x
sdarapu@hadoop-nn001:~$
```

```
sdarapu@hadoop-nn001:~$ kill -9 367095
sdarapu@hadoop-nn001:~$
```

For SpaceX data:

```
sdarapu@hadoop-nn001:~$ ps -x
  PID TTY          STAT TIME  COMMAND
 365215 ?            Ss   0:00 /lib/systemd/systemd --user
 365218 ?            S    0:00 (sd-pam)
 365225 ?          Ss1   0:00 /usr/bin/pulseaudio --daemonize=no --log-target=journal
 365244 ?            Ss   0:00 /usr/bin/dbus-daemon --session --address=systemd: --nofork --nopidfile --systemd-activation --syslog-only
 365302 ?            S    0:00 sshd: sdarapu@notty
 365303 ?            Ss   0:00 /usr/lib/openssh/sftp-server
 365928 ?            S    0:00 sshd: sdarapu@pts/5
 365930 pts/5      Ss   0:00 -bash
 366625 ?            S    0:00 sshd: sdarapu@notty
 366626 ?            Ss   0:00 /usr/lib/openssh/sftp-server
 369858 pts/5      S1   0:41 /usr/bin/java -Xmx60m -Dtwitter4j.streamBase.URL=https://stream.twitter.com/1.1/ -cp /usr/local/flume/conf:
 372112 pts/5      R+   0:00 ps -x
sdarapu@hadoop-nn001:~$ kill -9 369858
sdarapu@hadoop-nn001:~$
```

10. To check the data files are downloaded in the respective folders, I have executed the below command.

For NASA data:

```
sdarapu@hadoop-nn001:~$ hdfs dfs -ls /user/sdarapu/NASA_PA1data
2022-02-26 10:51:46,112 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where
Found 1 items
drwxr-xr-x - sdarapu sdarapu 0 2022-02-26 10:34 /user/sdarapu/NASA_PA1data/2022
[1]+ Killed
nohup $FLUME_HOME/bin/Flume-ng agent -n TwitterAgent -f /home/sdarapu/NASA.conf --conf /usr/local/Flume/conf
sdarapu@hadoop-nn001:~$
```

```
sdarapu@hadoop-nn001:~$ hdfs dfs -ls /user/sdarapu/NASA_PA1data/2022/02/26/10
2022-02-26 10:52:34,713 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes
Found 73 items
-rw-r--r-- 3 sdarapu sdarapu 49031 2022-02-26 10:34 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266371
-rw-r--r-- 3 sdarapu sdarapu 69846 2022-02-26 10:34 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266372
-rw-r--r-- 3 sdarapu sdarapu 65033 2022-02-26 10:35 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266373
-rw-r--r-- 3 sdarapu sdarapu 48972 2022-02-26 10:35 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266374
-rw-r--r-- 3 sdarapu sdarapu 74761 2022-02-26 10:35 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266375
-rw-r--r-- 3 sdarapu sdarapu 69502 2022-02-26 10:35 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266376
-rw-r--r-- 3 sdarapu sdarapu 45463 2022-02-26 10:35 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266377
-rw-r--r-- 3 sdarapu sdarapu 59067 2022-02-26 10:36 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266378
-rw-r--r-- 3 sdarapu sdarapu 51947 2022-02-26 10:36 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266379
-rw-r--r-- 3 sdarapu sdarapu 49956 2022-02-26 10:36 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266380
-rw-r--r-- 3 sdarapu sdarapu 57291 2022-02-26 10:36 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266381
-rw-r--r-- 3 sdarapu sdarapu 75806 2022-02-26 10:36 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266382
-rw-r--r-- 3 sdarapu sdarapu 62993 2022-02-26 10:37 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266383
-rw-r--r-- 3 sdarapu sdarapu 57919 2022-02-26 10:37 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266384
-rw-r--r-- 3 sdarapu sdarapu 51321 2022-02-26 10:37 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266385
-rw-r--r-- 3 sdarapu sdarapu 58128 2022-02-26 10:37 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266386
-rw-r--r-- 3 sdarapu sdarapu 53974 2022-02-26 10:38 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266387
-rw-r--r-- 3 sdarapu sdarapu 53476 2022-02-26 10:38 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266388
-rw-r--r-- 3 sdarapu sdarapu 54231 2022-02-26 10:38 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266389
-rw-r--r-- 3 sdarapu sdarapu 48141 2022-02-26 10:38 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266390
-rw-r--r-- 3 sdarapu sdarapu 61988 2022-02-26 10:38 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266391
-rw-r--r-- 3 sdarapu sdarapu 58997 2022-02-26 10:38 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266392
-rw-r--r-- 3 sdarapu sdarapu 45280 2022-02-26 10:39 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266393
-rw-r--r-- 3 sdarapu sdarapu 43000 2022-02-26 10:39 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266394
-rw-r--r-- 3 sdarapu sdarapu 54230 2022-02-26 10:39 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266395
-rw-r--r-- 3 sdarapu sdarapu 68000 2022-02-26 10:39 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266396
-rw-r--r-- 3 sdarapu sdarapu 67001 2022-02-26 10:39 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266397
-rw-r--r-- 3 sdarapu sdarapu 52100 2022-02-26 10:40 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266398
-rw-r--r-- 3 sdarapu sdarapu 49569 2022-02-26 10:40 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266399
-rw-r--r-- 3 sdarapu sdarapu 47433 2022-02-26 10:40 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266400
-rw-r--r-- 3 sdarapu sdarapu 46874 2022-02-26 10:40 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266401
-rw-r--r-- 3 sdarapu sdarapu 43956 2022-02-26 10:40 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266402
-rw-r--r-- 3 sdarapu sdarapu 53729 2022-02-26 10:41 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266403
-rw-r--r-- 3 sdarapu sdarapu 50057 2022-02-26 10:41 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266404
-rw-r--r-- 3 sdarapu sdarapu 54626 2022-02-26 10:41 /user/sdarapu/NASA_PA1data/2022/02/26/10/FlumeData.1645893266405
```


For SpaceX data:

[illegible]

- 12.** Now I used below command to check the size of the data that has been download after killing the respective processes.

For NASA data:

```
sdarapu@hadoop-nn001:~$ hdfs dfs -du -h /user/sdarapu/NASA_PA1data
2022-02-26 10:54:04,582 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform...
4.0 M 395.8 M /user/sdarapu/NASA_PA1data/2022
sdarapu@hadoop-nn001:~$
```

Size of the keyword NASA data downloaded is 4.0 MB

For SpaceX data:

```
sdarapu@hadoop-nn001:~$ hdfs dfs -du -h /user/sdarapu/SpaceX_PA1data
2022-02-26 11:29:14,490 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform...
3.8 M    11.4 M    /user/sdarapu/SpaceX_PA1data/2022
sdarapu@hadoop-nn001:~$
```

Size of the keyword SpaceX data downloaded is 3.8 MB