*University of Pittsburgh*
*School of Computing and Information*

# INFSCI 2750: Cloud Computing

# Spring 2020

**Mini Project 1**



**Submitted by:**

| SL. No. | Name | Email address |
|---|---|---|
| 1 | Debdas Ghosh | deg107@pitt.edu |
| 2 | Piu Mallick | pim16@pitt.edu |
| 3 | Timothy M. Mizerak | tmizerak@pc.pitt.edu |

## Objective

The objective of this mini project is to get familiar with setting up the Hadoop system and to start programming in Hadoop.

## Part 1: Setting up Hadoop

The first part of the project requires setting up Hadoop on the VMs. Each member of our team was assigned a VM, on which we installed and ran Hadoop. The goal of this part of the project is to build a Hadoop cluster. As the first step of our project, we made sure that Hadoop is set up properly on each of the VMs assigned to all the team members by running single node set up. Then, as the second step, we proceeded towards setting a 3-node cluster, with the assistance of the tutorial provided.
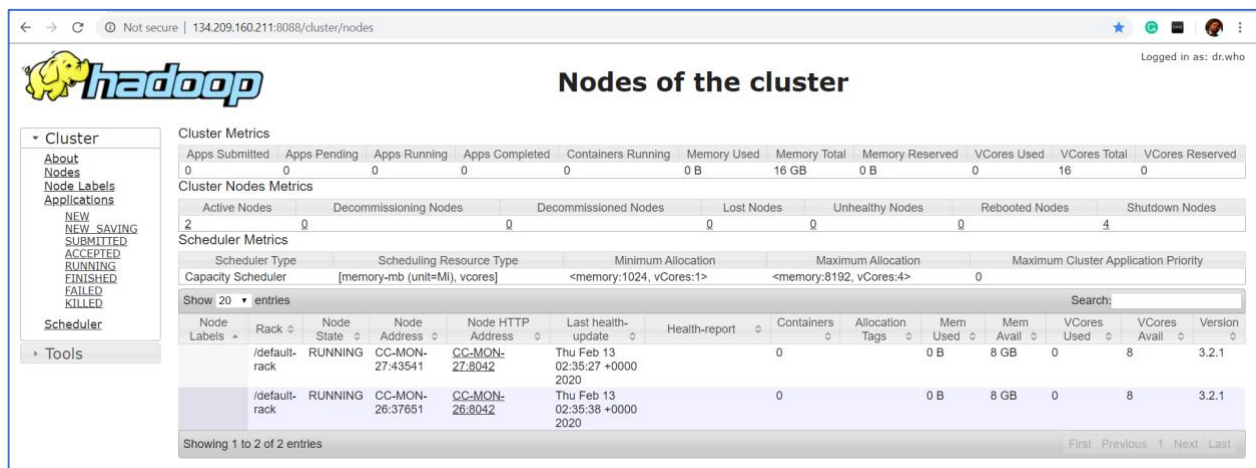
**Segregating the 3 VMs as Master & Slaves:**

**Master →** 134.209.160.211 (CC-MON-25)

**Slave 1 →** 64.225.17.118 (CC-MON-26)

**Slave 2 →** 138.197.96.66 (CC-MON-27)

**Screenshots showing cluster set-up and active status of the nodes**

**Hadoop**  Overview  Datanodes  Datanode Volume Failures  Snapshot  Startup Progress  Utilities ▾

# Overview 'CC-MON-25:9000' (active)

| Started: | Wed Feb 12 21:07:10 -0500 2020 |
|---|---|
| Version: | 3.2.1, rb3cbbb467e22ea829b3808f4b7b01d07e0bf3842 |
| Compiled: | Tue Sep 10 11:56:00 -0400 2019 by rohithsharmaks from branch-3.2.1 |
| Cluster ID: | CID-4998941e-333b-45d5-bf1b-f6266988aede |
| Block Pool ID: | BP-188256919-134.209.160.211-1581559557236 |

# Summary

Security is off.

Safemode is off.

1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem object(s).

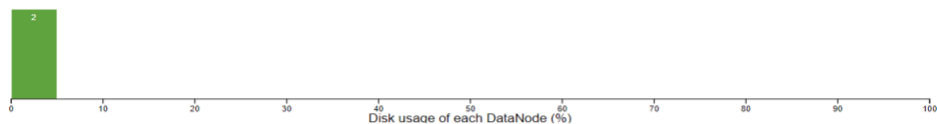Heap Memory used 104.1 MB of 233.5 MB Heap Memory. Max Heap Memory is 878.5 MB.

Non Heap Memory used 50.2 MB of 51.59 MB Commited Non Heap Memory. Max Non Heap Memory is <unbounded>.

| Configured Capacity: | 154.73 GB |
|---|---|
| Configured Remote Capacity: | 0 B |
| DFS Used: | 52 KB (0%) |

**Hadoop**  Overview  Datanodes  Datanode Volume Failures  Snapshot  Startup Progress  Utilities ▾

# Datanode Information

✔ In service  ❶ Down  ⊘ Decommissioning  ⊘ Decommissioned  ⏻ Decommissioned & dead
🔧 Entering Maintenance  🔧 In Maintenance  🔧 In Maintenance & dead

### Datanode usage histogram



Disk usage of each DataNode (%)

### In operation

Show 25 ▾ entries                                                                Search: _____

| Node | Http Address | Last contact | Last Block Report | Capacity | Blocks | Block pool used | Version |
|---|---|---|---|---|---|---|---|
| ✔CC-MON-26:9866 (64.225.17.118:9866) | http://CC-MON-26:9864 | 0s | 0m | 77.36 GB | 0 | 24 KB (0%) | 3.2.1 |
| ✔CC-MON-27:9866 (138.197.96.66:9866) | http://CC-MON-27:9864 | 0s | 28m | 77.36 GB | 0 | 28 KB (0%) | 3.2.1 |

**Part 2 → Requires us to build one Docker image to learn how to quick deploy a large-scale Hadoop Cluster in several minutes.**

As a first step, we learnt what a '**Docker image**' is and the steps or the procedure to build it.

**Running the default wordcount program that comes as part of the Hadoop package. Given an input file, the wordcount program prints the number of occurrences of each word in the file.**

# Building Docker Image

```
MINGW64:/d/docker                                                                                                    —  □  ×
Jebda@DESKTOP-TU68JG7 MINGW64 /d/docker
$ docker build  -t sequenceiq/hadoop-docker:2.7.1 .
Sending build context to Docker daemon  90.11kB
Step 1/58 : FROM sequenceiq/pam:centos-6.5
centos-6.5: Pulling from sequenceiq/pam
Image docker.io/sequenceiq/pam:centos-6.5 uses outdated schema1 manifest format. Please upgrade to a schema2 image for better future compatibility. More information at https://docs.docker.com/registry/spec/depre
cated-schema-v1/
b253335dcf03: Pull complete                                                                          a3ed95caeb02: Pull complete                        8d2023764774: Pull complete
                    69623ef05416: Pull complete                                                                                                                          ff0696749bf6: Pull complete
                                    0c3c0ff61963: Pull complete                            72accdc282f3: Pull complete           f252bbba6bda: Pull complete                                                       5298
ddb3b339: Pull complete
                    3984257f0553: Pull complete                                                                                              26343a20fa29: Pull complete
                                                            Digest: sha256:8f39a90603f88d2d372e55514162e844054459b91d9f82313af3ff16cd5d6a23
Status: Downloaded newer image for sequenceiq/pam:centos-6.5
 ---> f280c19a2e91
Step 2/58 : MAINTAINER SequenceIQ
 ---> Running in dd707d0620c1
Removing intermediate container dd707d0620c1
 ---> 930e503a0699
Step 3/58 : USER root
 ---> Running in 411cb2e174e1
Removing intermediate container 411cb2e174e1
 ---> 4ea16d4aed0f
Step 4/58 : RUN yum clean all;     rpm --rebuilddb;     yum install -y curl which tar sudo openssh-server openssh-clients rsync
 ---> Running in 6df29609118e
Loaded plugins: fastestmirror, keys, protect-packages, protectbase
Cleaning repos: base extras updates
Cleaning up Everything
Cleaning up list of fastest mirrors
Loaded plugins: fastestmirror, keys, protect-packages, protectbase
Setting up Install Process
Determining fastest mirrors
 * base: mirror.trouble-free.net
 * extras: distro.ibiblio.org
 * updates: ftp.osuosl.org
0 packages excluded due to repository protections
Resolving Dependencies
--> Running transaction check
---> Package curl.x86_64 0:7.19.7-40.el6_6.4 will be updated
---> Package curl.x86_64 0:7.19.7-54.el6_10 will be an update
--> Processing Dependency: libcurl = 7.19.7-54.el6_10 for package: curl-7.19.7-54.el6_10.x86_64
---> Package openssh-clients.x86_64 0:5.3p1-124.el6_10 will be installed
--> Processing Dependency: openssh = 5.3p1-124.el6_10 for package: openssh-clients-5.3p1-124.el6_10.x86_64
--> Processing Dependency: libfipscheck.so.1()(64bit) for package: openssh-clients-5.3p1-124.el6_10.x86_64
--> Processing Dependency: libedit.so.0()(64bit) for package: openssh-clients-5.3p1-124.el6_10.x86_64
---> Package openssh-server.x86_64 0:5.3p1-124.el6_10 will be installed
--> Processing Dependency: libwrap.so.0()(64bit) for package: openssh-server-5.3p1-124.el6_10.x86_64
---> Package rsync.x86_64 0:3.0.6-12.el6 will be installed
---> Package sudo.x86_64 0:1.8.6p3-29.el6_10.2 will be installed
```

```
MINGW64:/d/docker                                                                                                    —  □  ×
Jebda@DESKTOP-TU68JG7 MINGW64 /d/docker
$ docker pull sequenceiq/hadoop-docker:2.7.1
2.7.1: Pulling from sequenceiq/hadoop-docker
Image docker.io/sequenceiq/hadoop-docker:2.7.1 uses outdated schema1 manifest format. Please upgrade to a schema2 image for better future compatibility. More information at https://docs.docker.com/registry/spec/
deprecated-schema-v1/
b253335dcf03: Already exists                                                                        a3ed95caeb02: Pulling fs layer                     8d2023764774: Already exists
                    69623ef05416: Already exists                                                                                                                         ff0696749bf6: Already exists
                                    0c3c0ff61963: Already exists                           72accdc282f3: Already exists          f252bbba6bda: Already exists                                                      5298
ddb3b339: Already exists
                    3984257f0553: Already exists                                                                                             26343a20fa29: Already exists
                                    f3e272e0e801: Pulling fs layer                                                                                                       ad78a593ca62: Pulling fs layer
                                                            673712aa7667: Pulling fs layer                                                                                                           aaf06cd0
aa6e: Pulling fs layer                                                                   fed9c9377250: Pulling fs layer
                    d4385c519f63: Pulling fs layer                                                                   49ca93868354: Pulling fs layer
                                    98e62c38a660: Pull complete                                                                            3679d1cf91a0: Pull complete                   13605254d8c3
 : Pull complete                                               31ae294be02b: Pull complete                  a54805751dfa: Pull complete
                    38537e9c387f: Pull complete                                                                      dc639853e053: Pull complete
                                    e267620cd7fd: Pull complete                                                                    93990a6b26ca: Pull complete
                                                            11ffe2baf32d: Pull complete              adede6edfea0: Pull complete                                                                c91b10bf3a44: Pu
ll complete
                    4afb2f219fa7: Pull complete                                                                      0335bc4000de: Pull complete                   3bb2b06400be: Pull complete
                                    e6c5265506dc: Pull complete                                                                                                                                 2a1a28b12647: Pull c
omplete                                                        d9665143ac9a: Pull complete                  5c175609cbf3: Pull complete
                    e2a7d6798159: Pull complete                                                                      88d87e462c71: Pull complete
                                    3a404fc6437e: Pull complete                                                                    5517052ef612: Pull complete
                                                            fa61c616ddd1: Pull complete                                                                                                             d4ab0c19cb91: Pull compl
ete                                                            9aa826a9ca93: Pull complete
                    b2ecd44f6d78: Pull complete                                                                      824658b0b14c: Pull complete
                                    e5c31d8cbbce: Pull complete                                                                    Digest: sha256:2da37e4eeea57bc99dd64987391ce9e1384c63b4fa
56b7525a60849a758fb950
Status: Downloaded newer image for sequenceiq/hadoop-docker:2.7.1
docker.io/sequenceiq/hadoop-docker:2.7.1

Jebda@DESKTOP-TU68JG7 MINGW64 /d/docker
$ docker run -it sequenceiq/hadoop-docker:2.7.1 /etc/bootstrap.sh -bash
/
Starting sshd:                                              [  OK  ]
20/02/12 14:33:15 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Starting namenodes on [0e8195ea77e1]
0e8195ea77e1: starting namenode, logging to /usr/local/hadoop/logs/hadoop-root-namenode-0e8195ea77e1.out
localhost: starting datanode, logging to /usr/local/hadoop/logs/hadoop-root-datanode-0e8195ea77e1.out
Starting secondary namenodes [0.0.0.0]
0.0.0.0: starting secondarynamenode, logging to /usr/local/hadoop/logs/hadoop-root-secondarynamenode-0e8195ea77e1.out
20/02/12 14:33:52 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
starting yarn daemons
starting resourcemanager, logging to /usr/local/hadoop/logs/yarn--resourcemanager-0e8195ea77e1.out
localhost: starting nodemanager, logging to /usr/local/hadoop/logs/yarn-root-nodemanager-0e8195ea77e1.out
```

```
MINGW64:/d/docker                                                                                                    —  □  ×
bash-4.1# bin/hdfs dfs -cat output/*
20/02/12 14:50:01 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
6       dfs.audit.logger
4       dfs.class
3       dfs.server.namenode.
2       dfs.period
2       dfs.audit.log.maxfilesize
2       dfs.audit.log.maxbackupindex
1       dfsmetrics.log
1       dfsadmin
1       dfs.servers
1       dfs.replication
1       dfs.file
bash-4.1# ls
LICENSE.txt  NOTICE.txt  README.txt  bin  etc  include  input  lib  libexec  logs  sbin  share
bash-4.1#
```

## Part 3: Developing a Hadoop program

We are required to develop a Hadoop program that produces the **ngram-frequencies** of the text in a given input file. **n-gram is a contiguous sequence of n characters from a given sequence of text**.

An **n-gram of size 1** is referred to as a "**unigram**"; size **2** is a "**bigram**" (or, less commonly, a "**digram**"); size **3** is a "**trigram**". For example, the **2-gram frequency** in the text, "**Helloworld**" is as follows:

He-1, el-1, ll-1, lo-1, ow-1, wo-1, or-1, rl-1, ld-1

The program should accept **n** as an **input parameter** and produce the **n-gram frequencies** in the **text** as an **output file**.

**Answer:**

**Source code** and **readme file** (explaining the procedure of the program) is provided separately in the zip file.

**Output of the program:**

The steps to execute the jar file are mentioned in the readme file (inside the folder provided).

## Part 4: Developing a Hadoop program to analyze real logs

In this part, we developed several MapReduce programs to analyze a real anonymous log to answer several questions based on the log. The log is in '**access_log.zip**', provided in the **Course Web**.

The log file is in Common Log Format:
10.223.157.186 - - [15/Jul/2009:15:50:35 -0700] "GET /assets/js/lowpro.js HTTP/1.1"
200 10469
%h %l %u %t \"%r\" %>s %b

```
Where:
• %h is the IP address of the client
• %l is identity of the client, or "-" if it's unavailable
• %u is username of the client, or "-" if it's unavailable
• %t is the time that the server finished processing the request. The format is
[day/month/year:hour:minute:second zone]
• %r is the request line from the client is given (in double quotes). It contains
the method, path,
query-string, and protocol or the request.
• %>s is the status code that the server sends back to the client. You will see
see mostly status
codes 200 (OK - The request has succeeded), 304 (Not Modified) and 404 (Not
Found). See
more information on status codes in W3C.org
• %b is the size of the object returned to the client, in bytes. It will be "-
" in case of status code
304.
```

**Answer:**

**Source code** and **readme file** (explaining the procedure of the program) is provided separately in the zip file.


**Problems:**

1. **How many hits were made to the website item "/assets/img/home-logo.png"?**

   **Answer**: 98744

## 2. How many hits were made from the IP: 10.153.239.5?

**Answer**: 547

```
Launched reduce tasks=1
Data-local map tasks=5
Total time spent by all maps in occupied slots (ms)=84546
Total time spent by all reduces in occupied slots (ms)=16078
Total time spent by all map tasks (ms)=84546
Total time spent by all reduce tasks (ms)=16078
Total vcore-milliseconds taken by all map tasks=84546
Total vcore-milliseconds taken by all reduce tasks=16078
Total megabyte-milliseconds taken by all map tasks=86575104
Total megabyte-milliseconds taken by all reduce tasks=16463872
Map-Reduce Framework
Map input records=4477843
Map output records=547
Map output bytes=9299
Map output materialized bytes=43
Input split bytes=460
Combine input records=547
Combine output records=1
Reduce input groups=1
Reduce shuffle bytes=43
Reduce input records=1
Reduce output records=1
Spilled Records=2
Shuffled Maps =4
Failed Shuffles=0
Merged Map outputs=4
GC time elapsed (ms)=1155
CPU time spent (ms)=31650
Physical memory (bytes) snapshot=2092666880
Virtual memory (bytes) snapshot=13151399936
Total committed heap usage (bytes)=1786773504
Peak Map Physical memory (bytes)=487890944
Peak Map Virtual memory (bytes)=2630201344
Peak Reduce Physical memory (bytes)=148803584
Peak Reduce Virtual memory (bytes)=2637877248
Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
Bytes Read=504953820
File Output Format Counters
Bytes Written=17
student@CC-MON-25:~/hadoop$ bin/hdfs dfs -cat output/*
2020-02-15 00:41:36,473 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
10.153.239.5    547
student@CC-MON-25:~/hadoop$
```

3.  **Which path in the website has been hit most? How many hits were made to the path?**

**Answer**: /assets/css/combined.css      117348

```
student@CC-MON-25:~/hadoop$ bin/hadoop jar AccessLog3.jar  input1/ output/
2020-02-15 00:25:15,172 INFO client.RMProxy: Connecting to ResourceManager at CC-MON-25/134.209.160.211:8032
2020-02-15 00:25:15,578 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to re
medy this.
2020-02-15 00:25:15,594 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/student/.staging/job_1581557217626_0012
2020-02-15 00:25:15,702 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2020-02-15 00:25:15,905 INFO input.FileInputFormat: Total input files to process : 1
2020-02-15 00:25:16,011 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2020-02-15 00:25:16,056 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2020-02-15 00:25:16,067 INFO mapreduce.JobSubmitter: number of splits:4
2020-02-15 00:25:16,205 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2020-02-15 00:25:16,234 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1581557217626_0012
2020-02-15 00:25:16,235 INFO mapreduce.JobSubmitter: Executing with tokens: []
2020-02-15 00:25:16,447 INFO conf.Configuration: resource-types.xml not found
2020-02-15 00:25:16,447 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2020-02-15 00:25:16,716 INFO impl.YarnClientImpl: Submitted application application_1581557217626_0012
2020-02-15 00:25:16,764 INFO mapreduce.Job: The url to track the job: http://CC-MON-25:8088/proxy/application_1581557217626_0012/
2020-02-15 00:25:16,765 INFO mapreduce.Job: Running job: job_1581557217626_0012
2020-02-15 00:25:25,967 INFO mapreduce.Job: Job job_1581557217626_0012 running in uber mode : false
2020-02-15 00:25:25,982 INFO mapreduce.Job:  map 0% reduce 0%
2020-02-15 00:25:39,146 INFO mapreduce.Job:  map 25% reduce 0%
2020-02-15 00:25:47,216 INFO mapreduce.Job:  map 38% reduce 0%
2020-02-15 00:25:48,225 INFO mapreduce.Job:  map 57% reduce 0%
2020-02-15 00:25:53,276 INFO mapreduce.Job:  map 61% reduce 0%
2020-02-15 00:25:54,285 INFO mapreduce.Job:  map 73% reduce 0%
2020-02-15 00:25:56,303 INFO mapreduce.Job:  map 82% reduce 0%
2020-02-15 00:25:58,322 INFO mapreduce.Job:  map 90% reduce 0%
2020-02-15 00:25:59,328 INFO mapreduce.Job:  map 100% reduce 67%
2020-02-15 00:26:00,347 INFO mapreduce.Job:  map 100% reduce 100%
2020-02-15 00:26:01,369 INFO mapreduce.Job: Job job_1581557217626_0012 completed successfully
2020-02-15 00:26:01,514 INFO mapreduce.Job: Counters: 55
        File System Counters
                FILE: Number of bytes read=4836356
                FILE: Number of bytes written=10800260
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=504954288
                HDFS: Number of bytes written=32
                HDFS: Number of read operations=17
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
                HDFS: Number of bytes read erasure-coded=0
        Job Counters
                Killed map tasks=1
                Launched map tasks=5
                Launched reduce tasks=1
                Data-local map tasks=5
                Total time spent by all maps in occupied slots (ms)=105108
                Total time spent by all reduces in occupied slots (ms)=18948
```



```
                Launched reduce tasks=1
                Data-local map tasks=5
                Total time spent by all maps in occupied slots (ms)=105108
                Total time spent by all reduces in occupied slots (ms)=18948
                Total time spent by all map tasks (ms)=105108
                Total time spent by all reduce tasks (ms)=18948
                Total vcore-milliseconds taken by all map tasks=105108
                Total vcore-milliseconds taken by all reduce tasks=18948
                Total megabyte-milliseconds taken by all map tasks=107630592
                Total megabyte-milliseconds taken by all reduce tasks=19402752
        Map-Reduce Framework
                Map input records=4477843
                Map output records=4477843
                Map output bytes=202784119
                Map output materialized bytes=4836374
                Input split bytes=468
                Combine input records=4477843
                Combine output records=78913
                Reduce input groups=42386
                Reduce shuffle bytes=4836374
                Reduce input records=78913
                Reduce output records=1
                Spilled Records=157826
                Shuffled Maps =4
                Failed Shuffles=0
                Merged Map outputs=4
                GC time elapsed (ms)=1293
                CPU time spent (ms)=44300
                Physical memory (bytes) snapshot=2125549568
                Virtual memory (bytes) snapshot=13163786240
                Total committed heap usage (bytes)=1784152064
                Peak Map Physical memory (bytes)=491540480
                Peak Map Virtual memory (bytes)=2631966720
                Peak Reduce Physical memory (bytes)=168316928
                Peak Reduce Virtual memory (bytes)=2638016512
        Shuffle Errors
                BAD_ID=0
                CONNECTION=0
                IO_ERROR=0
                WRONG_LENGTH=0
                WRONG_MAP=0
                WRONG_REDUCE=0
        File Input Format Counters
                Bytes Read=504953820
        File Output Format Counters
                Bytes Written=32
student@CC-MON-25:~/hadoop$ bin/hdfs dfs -cat output/*
2020-02-15 00:26:19,575 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
/assets/css/combined.css        117348
student@CC-MON-25:~/hadoop$
```

**4. Which IP accesses the website most? How many accesses were made by it?**

**Answer**: 10.216.113.172               158614

```
student@CC-MON-25: ~/hadoop
student@CC-MON-25:~/hadoop$ bin/hadoop jar AccessLog4.jar  input1/ output/
2020-02-15 00:12:29,512 INFO client.RMProxy: Connecting to ResourceManager at CC-MON-25/134.209.160.211:8032
2020-02-15 00:12:29,921 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to re
medy this.
2020-02-15 00:12:29,934 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/student/.staging/job_1581557217626_0011
2020-02-15 00:12:30,049 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2020-02-15 00:12:30,320 INFO input.FileInputFormat: Total input files to process : 1
2020-02-15 00:12:30,361 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2020-02-15 00:12:30,409 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2020-02-15 00:12:30,417 INFO mapreduce.JobSubmitter: number of splits:4
2020-02-15 00:12:30,541 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2020-02-15 00:12:30,583 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1581557217626_0011
2020-02-15 00:12:30,584 INFO mapreduce.JobSubmitter: Executing with tokens: []
2020-02-15 00:12:30,775 INFO conf.Configuration: resource-types.xml not found
2020-02-15 00:12:30,776 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2020-02-15 00:12:30,842 INFO impl.YarnClientImpl: Submitted application application_1581557217626_0011
2020-02-15 00:12:30,883 INFO mapreduce.Job: The url to track the job: http://CC-MON-25:8088/proxy/application_1581557217626_0011/
2020-02-15 00:12:30,885 INFO mapreduce.Job: Running job: job_1581557217626_0011
2020-02-15 00:12:40,030 INFO mapreduce.Job: Job job_1581557217626_0011 running in uber mode : false
2020-02-15 00:12:40,031 INFO mapreduce.Job:  map 0% reduce 0%
2020-02-15 00:12:52,183 INFO mapreduce.Job:  map 25% reduce 0%
2020-02-15 00:13:01,262 INFO mapreduce.Job:  map 67% reduce 0%
2020-02-15 00:13:05,302 INFO mapreduce.Job:  map 75% reduce 0%
2020-02-15 00:13:07,323 INFO mapreduce.Job:  map 83% reduce 0%
2020-02-15 00:13:09,348 INFO mapreduce.Job:  map 83% reduce 17%
2020-02-15 00:13:10,360 INFO mapreduce.Job:  map 100% reduce 17%
2020-02-15 00:13:12,376 INFO mapreduce.Job:  map 100% reduce 100%
2020-02-15 00:13:13,396 INFO mapreduce.Job: Job job_1581557217626_0011 completed successfully
2020-02-15 00:13:13,512 INFO mapreduce.Job: Counters: 55
        File System Counters
                FILE: Number of bytes read=6840959
                FILE: Number of bytes written=14809466
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=504954288
                HDFS: Number of bytes written=22
                HDFS: Number of read operations=17
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
                HDFS: Number of bytes read erasure-coded=0
        Job Counters
                Killed map tasks=1
                Launched map tasks=5
                Launched reduce tasks=1
                Data-local map tasks=5
                Total time spent by all maps in occupied slots (ms)=97480
                Total time spent by all reduces in occupied slots (ms)=17258
                Total time spent by all map tasks (ms)=97480
                Total time spent by all reduce tasks (ms)=17258
```



```
student@CC-MON-25: ~/hadoop
                Launched reduce tasks=1
                Data-local map tasks=5
                Total time spent by all maps in occupied slots (ms)=97480
                Total time spent by all reduces in occupied slots (ms)=17258
                Total time spent by all map tasks (ms)=97480
                Total time spent by all reduce tasks (ms)=17258
                Total vcore-milliseconds taken by all map tasks=97480
                Total vcore-milliseconds taken by all reduce tasks=17258
                Total megabyte-milliseconds taken by all map tasks=99819520
                Total megabyte-milliseconds taken by all reduce tasks=17672192
        Map-Reduce Framework
                Map input records=4477843
                Map output records=4477843
                Map output bytes=79641285
                Map output materialized bytes=6840977
                Input split bytes=468
                Combine input records=4477843
                Combine output records=347168
                Reduce input groups=333923
                Reduce shuffle bytes=6840977
                Reduce input records=347168
                Reduce output records=1
                Spilled Records=694336
                Shuffled Maps =4
                Failed Shuffles=0
                Merged Map outputs=4
                GC time elapsed (ms)=1079
                CPU time spent (ms)=44580
                Physical memory (bytes) snapshot=2155749376
                Virtual memory (bytes) snapshot=13166747648
                Total committed heap usage (bytes)=1834483712
                Peak Map Physical memory (bytes)=493744128
                Peak Map Virtual memory (bytes)=2635468800
                Peak Reduce Physical memory (bytes)=199041024
                Peak Reduce Virtual memory (bytes)=2637692928
        Shuffle Errors
                BAD_ID=0
                CONNECTION=0
                IO_ERROR=0
                WRONG_LENGTH=0
                WRONG_MAP=0
                WRONG_REDUCE=0
        File Input Format Counters
                Bytes Read=504953820
        File Output Format Counters
                Bytes Written=22
student@CC-MON-25:~/hadoop$ bin/hdfs dfs -cat output/*
2020-02-15 00:13:20,921 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
10.216.113.172  158614
student@CC-MON-25:~/hadoop$
```

The steps to execute the jar files are mentioned in the readme file (inside the folder provided).