

LAB 2 INSTRUCTIONS

DS8003 – MGT OF BIG DATA AND TOOLS

Ryerson University

Instructor: Kanchana Padmanabhan

Lab & Assignments

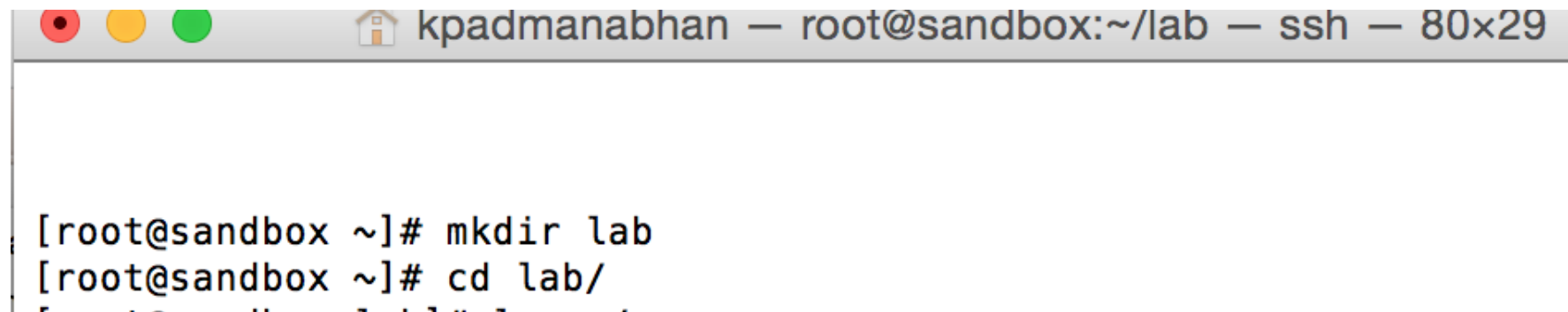
- **HDP Sandbox (VIRTUAL BOX LOGIN)**
 - ▣ SSH using Putty (host: localhost, port: 2222)
 - ▣ Username: root
 - ▣ Password: depends on what you set it
- Instructions for setup are in Lab 1

Download file

- Download the “shakespeare_100.txt” file from Datasets & Scripts

Make a new directory on Virtual Machine

- SSH using Putty and Log-in to your Virtual Machine.
- Type the following commands

A screenshot of a terminal window. The title bar shows a home icon, the username 'kpadmanabhan', the host 'root@sandbox', the current directory '~/lab', and the connection type 'ssh' with dimensions '80x29'. The terminal content shows two commands being executed: '[root@sandbox ~]# mkdir lab' and '[root@sandbox ~]# cd lab/'.

```
kpadmanabhan — root@sandbox:~/lab — ssh — 80x29

[root@sandbox ~]# mkdir lab
[root@sandbox ~]# cd lab/
```

- This will make a new directory called “lab.”
- Use “cd” command to navigate into the lab folder
- Move the file “shakespere_100.txt” into the “lab” folder on VirtualBox using FileZilla (See Lab 1 for instructions)

Today's lab – HDFS Shell Commands

- Step 1: Create a directory in HDFS,
- Step 2: Upload a file & List files in Directory
- Step 3: Try few basic Linux commands
- Step 4: Download Files From HDFS to Local File System
- Step 5: Find Out Space Utilization in a HDFS Directory
- Step 6: Explore Three Advanced Features
- Step 7: Use Help Command

Hadoop FileSystem Shell

Hadoop Filesystem Command	Description
<code>hadoop fs -mkdir</code>	create a new directory in hdfs
<code>hadoop fs -ls</code>	list files in a directory
<code>hadoop fs -put</code>	to copy file from local to hdfs
<code>hadoop fs -cat</code>	to preview the content of an hdfs file
<code>hadoop fs -get</code>	to move file from hdfs to local
<code>hadoop fs -rmdir</code>	to delete a directory
<code>hadoop fs -cp</code>	to make a copy of an hdfs file
<code>hadoop fs -du</code>	to display the size of an hdfs file
<code>hadoop fs -mv</code>	to move hdfs files from source to destination
<code>hadoop fs -tail</code>	to print the last few lines of an hdfs file/directory
<code>hadoop fs -head</code>	to print the first few lines of an hdfs file
<code>hadoop fs -getmerge</code>	to merge several hdfs files into one single file and copy to local

To learn more about Hadoop shell commands with Example Usage, check out the documentations

<http://hadoop.apache.org/docs/r2.6.0/hadoop-project-dist/hadoop-common/FileSystemShell.html>

https://hadoop.apache.org/docs/r1.0.4/file_system_shell.html

Step 1: Create a folder on HDFS that can be accessed by the user “root”

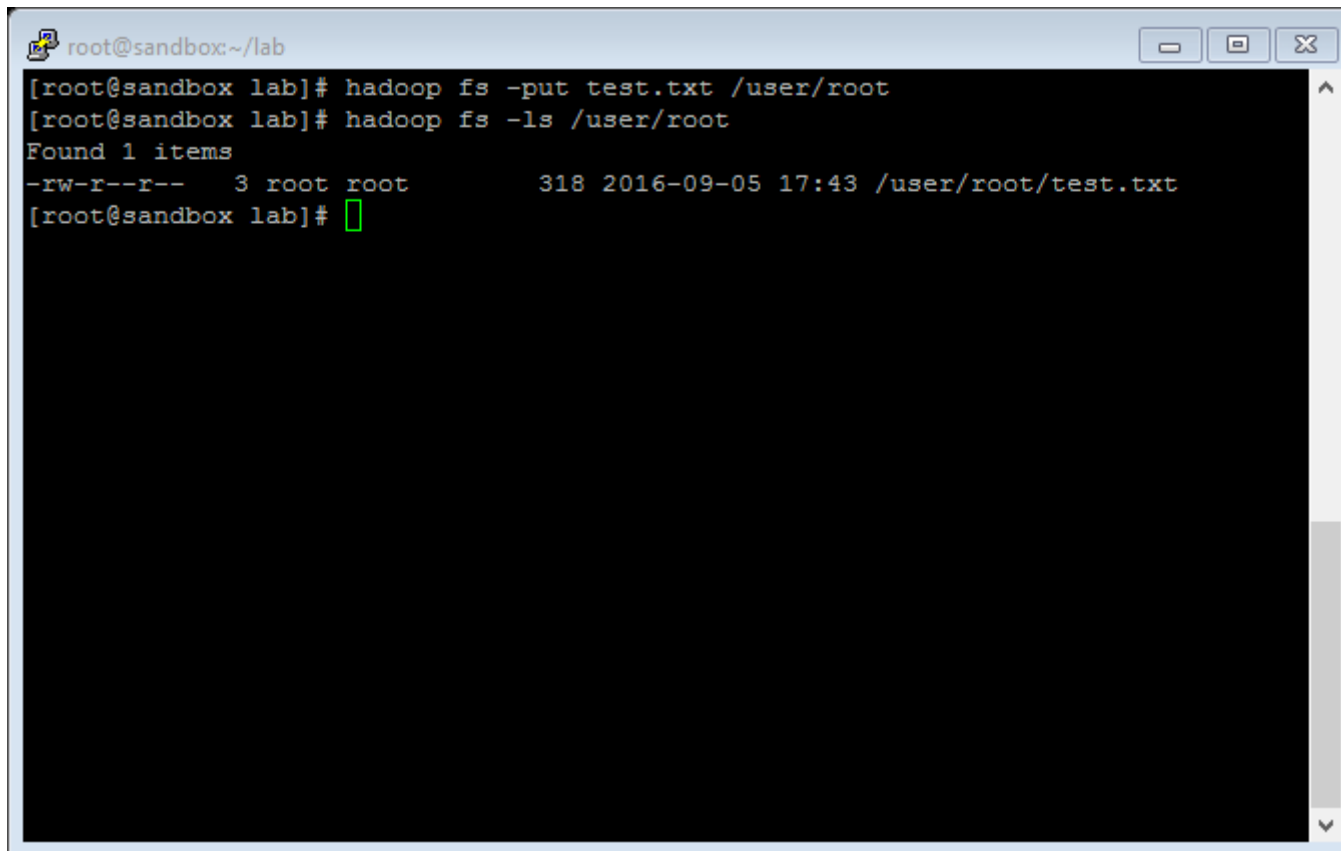
- Let's create a folder on the HDFS that can be accessed as “root” user
- “hadoop fs” command let's us access the hdfs
- “-mkdir” creates a directory.
- Here we are creating a folder in the distributed file system
- “-chown” let's us change the owner of the folder from “hdfs” to “root”

```
kpadmanabhan — root@sandbox:~ — ssh — 80x44

[root@sandbox ~]# sudo -u hdfs hadoop fs -ls /user
Found 11 items
drwxrwx--- - ambari-qa hdfs      0 2015-10-27 12:39 /user/ambari-qa
drwxr-xr-x - guest    guest    0 2015-10-27 12:55 /user/guest
drwxr-xr-x - hcat     hdfs      0 2015-10-27 12:43 /user/hcat
drwx----- - hdfs     hdfs      0 2016-01-14 21:51 /user/hdfs
drwx----- - hive     hdfs      0 2016-01-15 19:51 /user/hive
drwxrwxrwx - hue      hdfs      0 2015-10-27 12:55 /user/hue
drwxrwxr-x - oozie    hdfs      0 2015-10-27 12:44 /user/oozie
drwxr-xr-x - solr     hdfs      0 2015-10-27 12:48 /user/solr
drwxrwxr-x - spark    hdfs      0 2015-10-27 12:41 /user/spark
drwxr-xr-x - unit     hdfs      0 2015-10-27 12:46 /user/unit
drwxr-xr-x - zeppelin zeppelin 0 2015-10-27 12:10 /user/zeppelin
[root@sandbox ~]# sudo -u hdfs hadoop fs -mkdir /user/root
[root@sandbox ~]# sudo -u hdfs hadoop fs -ls /user
Found 12 items
drwxrwx--- - ambari-qa hdfs      0 2015-10-27 12:39 /user/ambari-qa
drwxr-xr-x - guest    guest    0 2015-10-27 12:55 /user/guest
drwxr-xr-x - hcat     hdfs      0 2015-10-27 12:43 /user/hcat
drwx----- - hdfs     hdfs      0 2016-01-14 21:51 /user/hdfs
drwx----- - hive     hdfs      0 2016-01-15 19:51 /user/hive
drwxrwxrwx - hue      hdfs      0 2015-10-27 12:55 /user/hue
drwxrwxr-x - oozie    hdfs      0 2015-10-27 12:44 /user/oozie
drwxr-xr-x - hdfs     hdfs      0 2016-01-22 01:53 /user/root
drwxrwxr-x - solr     hdfs      0 2015-10-27 12:48 /user/solr
drwxrwxr-x - spark    hdfs      0 2015-10-27 12:41 /user/spark
drwxr-xr-x - unit     hdfs      0 2015-10-27 12:46 /user/unit
drwxr-xr-x - zeppelin zeppelin 0 2015-10-27 12:10 /user/zeppelin
[root@sandbox ~]# sudo -u hdfs hadoop fs -chown root:root /user/root
[root@sandbox ~]# sudo -u hdfs hadoop fs -ls /user
Found 12 items
drwxrwx--- - ambari-qa hdfs      0 2015-10-27 12:39 /user/ambari-qa
drwxr-xr-x - guest    guest    0 2015-10-27 12:55 /user/guest
drwxr-xr-x - hcat     hdfs      0 2015-10-27 12:43 /user/hcat
drwx----- - hdfs     hdfs      0 2016-01-14 21:51 /user/hdfs
drwx----- - hive     hdfs      0 2016-01-15 19:51 /user/hive
drwxrwxrwx - hue      hdfs      0 2015-10-27 12:55 /user/hue
drwxrwxr-x - oozie    hdfs      0 2015-10-27 12:44 /user/oozie
drwxr-xr-x - root     root      0 2016-01-22 01:53 /user/root
```

Step 2: Let's copy a file to the HDFS

- “-put” command let's you copy a file from local file system to the HDFS
- “-ls” command let's you copy a list all files in directory

A terminal window titled 'root@sandbox:~/lab' with standard window controls. The terminal shows the execution of two Hadoop commands. The first command, 'hadoop fs -put test.txt /user/root', successfully uploads a file. The second command, 'hadoop fs -ls /user/root', lists the contents of the directory, showing one file with permissions '-rw-r--r--', size '318', and timestamp '2016-09-05 17:43'.

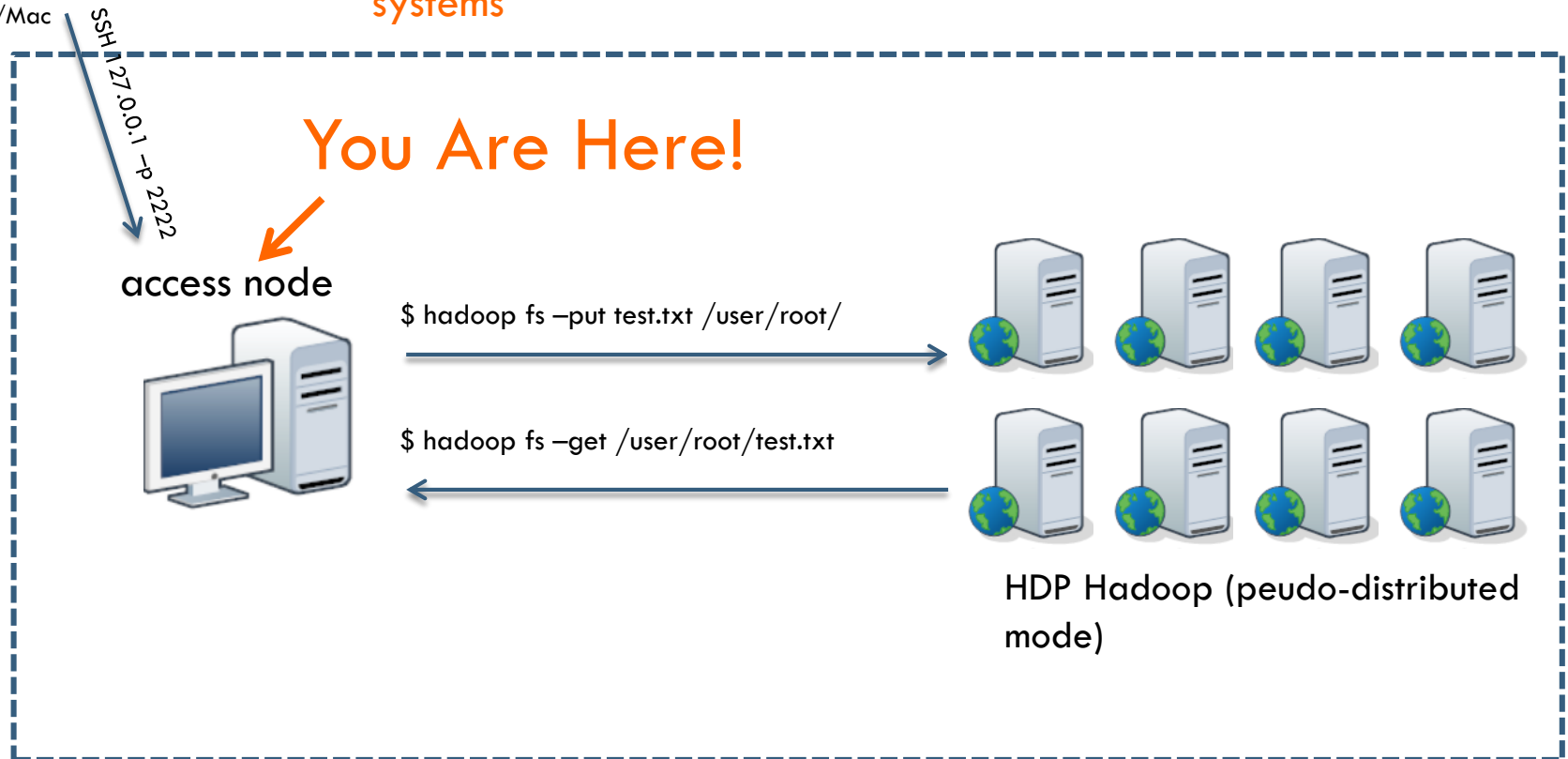
```
root@sandbox:~/lab
[root@sandbox lab]# hadoop fs -put test.txt /user/root
[root@sandbox lab]# hadoop fs -ls /user/root
Found 1 items
-rw-r--r--   3 root root      318 2016-09-05 17:43 /user/root/test.txt
[root@sandbox lab]#
```


Architecture



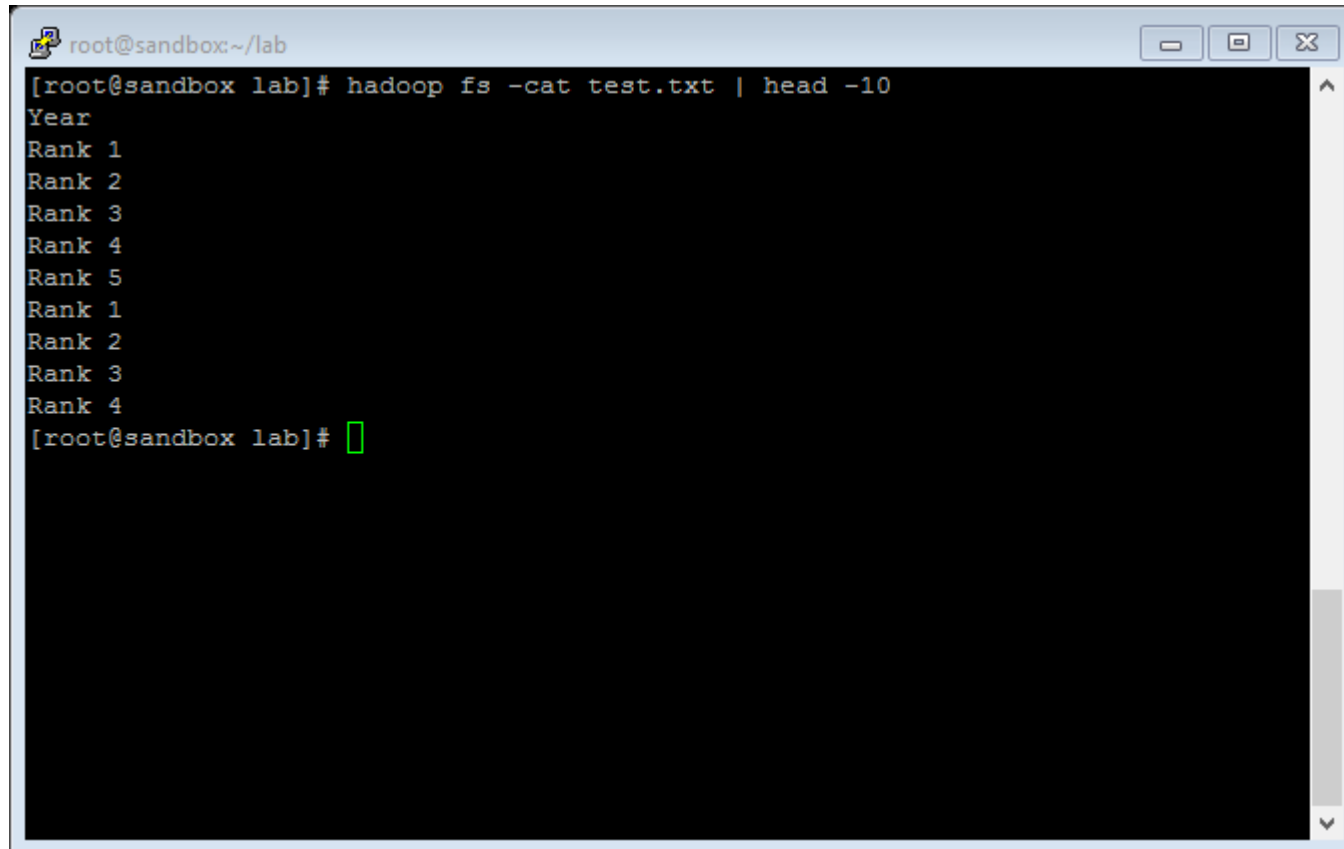
Lab computer
Your PC/Mac

NOTE: Your client linux file system and hadoop filesystem (HDFS) are separate environment. First, you need to learn how to move files between the two file systems



Step 3: Try a few basic Linux commands

- “-cat” concatenates the contents of the file to the screen
- “head” to print the first few lines of an hdfs file

A terminal window titled 'root@sandbox:~/lab' with standard window controls. The command '[root@sandbox lab]# hadoop fs -cat test.txt | head -10' has been executed. The output is displayed on a black background with white text, showing 'Year' followed by eight lines of 'Rank' values (Rank 1, Rank 2, Rank 3, Rank 4, Rank 5, Rank 1, Rank 2, Rank 3, Rank 4). The prompt '[root@sandbox lab]# ' is visible at the bottom with a green cursor.

```
root@sandbox:~/lab
[root@sandbox lab]# hadoop fs -cat test.txt | head -10
Year
Rank 1
Rank 2
Rank 3
Rank 4
Rank 5
Rank 1
Rank 2
Rank 3
Rank 4
[root@sandbox lab]#
```

Recall



Lab computer
Your PC/Mac

NOTE: Your client linux file system and hadoop filesystem (HDFS) are separate environment. First, you need to learn how to move files between the two file systems

You Are Here!

ssh 127.0.0.1 -p 2222

access node



\$ `hadoop fs -put myFiles.txt /user/lab`

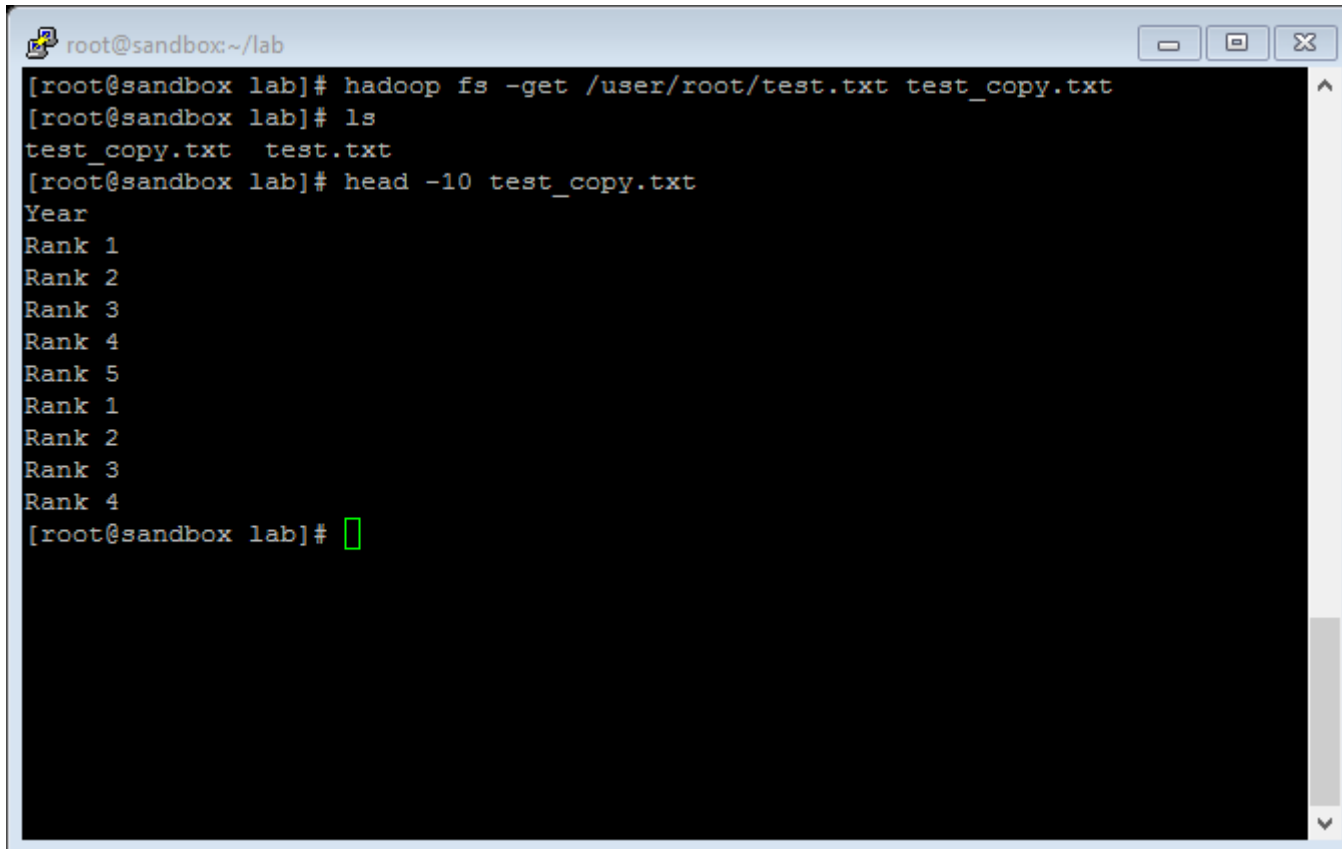
\$ `hadoop fs -get /user/lab/myFile.txt`



HDP Hadoop (pseudo-distributed mode)

Step 4: Download Files From HDFS to Local File System

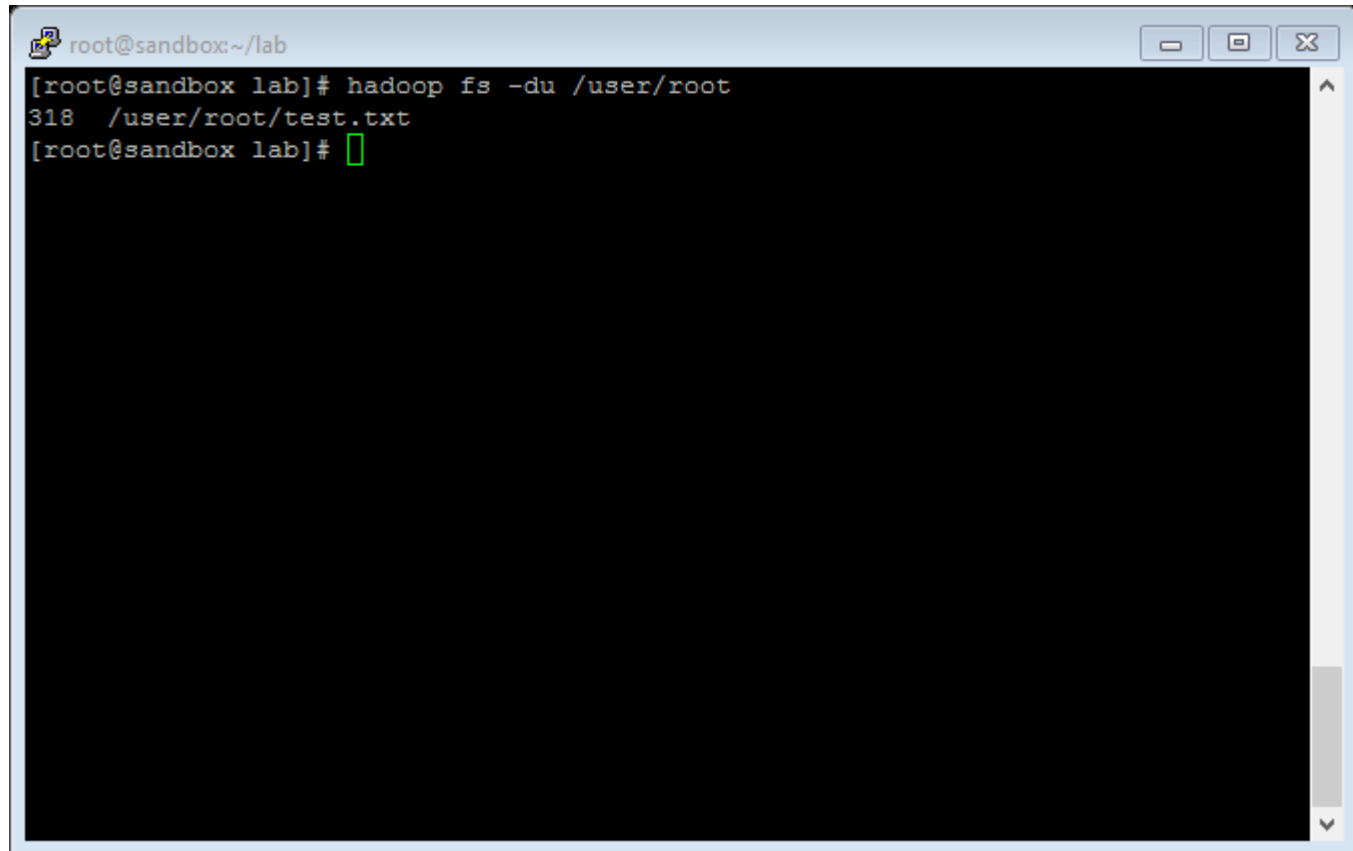
- ❑ “-get” Copies/Downloads files from HDFS to the local file system

A terminal window titled 'root@sandbox ~/lab' with standard window controls. The terminal shows the execution of 'hadoop fs -get /user/root/test.txt test_copy.txt', followed by 'ls' which lists 'test_copy.txt' and 'test.txt'. Then 'head -10 test_copy.txt' is run, displaying the first 10 lines of the file: 'Year', 'Rank 1', 'Rank 2', 'Rank 3', 'Rank 4', 'Rank 5', 'Rank 1', 'Rank 2', 'Rank 3', and 'Rank 4'. The prompt is currently at '[root@sandbox lab]#'.

```
root@sandbox ~/lab
[root@sandbox lab]# hadoop fs -get /user/root/test.txt test_copy.txt
[root@sandbox lab]# ls
test_copy.txt  test.txt
[root@sandbox lab]# head -10 test_copy.txt
Year
Rank 1
Rank 2
Rank 3
Rank 4
Rank 5
Rank 1
Rank 2
Rank 3
Rank 4
[root@sandbox lab]#
```

Step 5: Find Out Space Utilization in a HDFS Directory

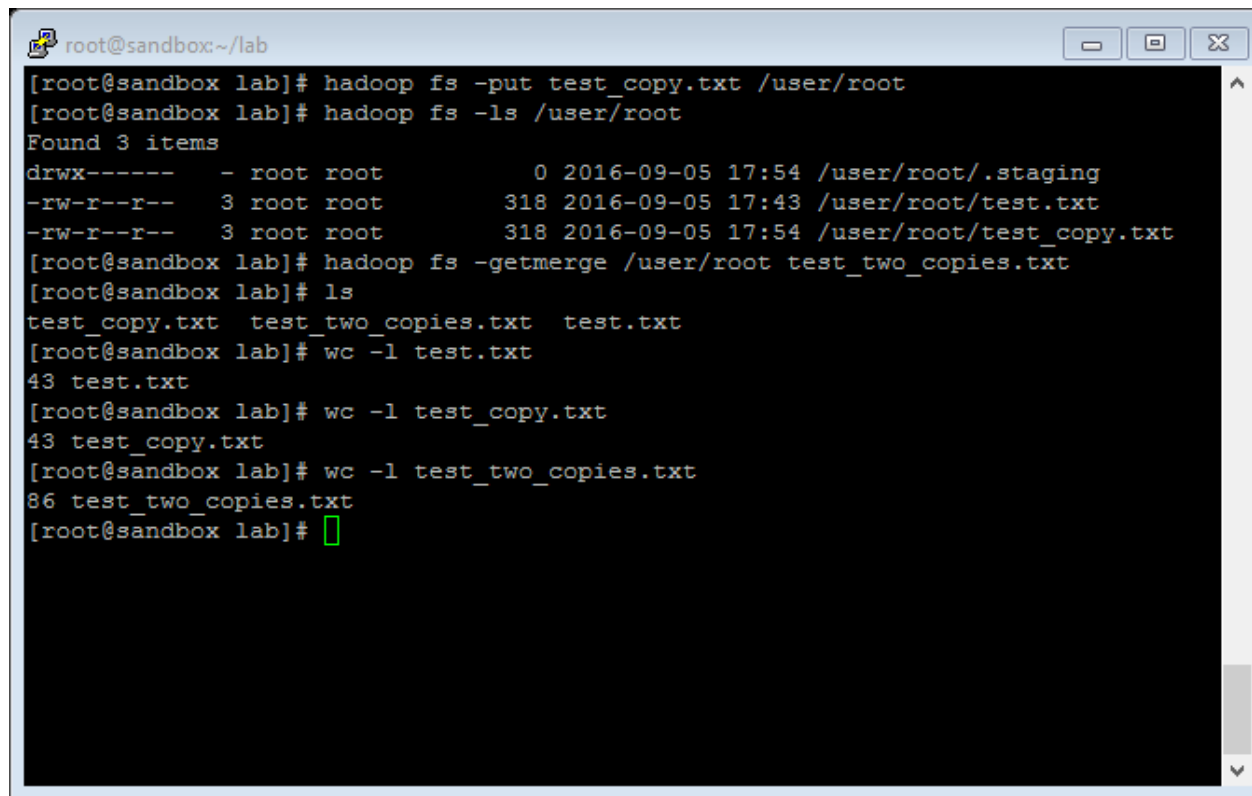
- ❑ `-du` Displays size of files and directories contained in the given directory or the size of a file if its just a file

A terminal window titled 'root@sandbox ~/lab' with standard window controls. The terminal shows the command 'hadoop fs -du /user/root' being executed, which returns '318 /user/root/test.txt'. The prompt is then ready for the next command.

```
root@sandbox ~/lab
[root@sandbox lab]# hadoop fs -du /user/root
318 /user/root/test.txt
[root@sandbox lab]#
```

Step 6: Advanced HDFS feature

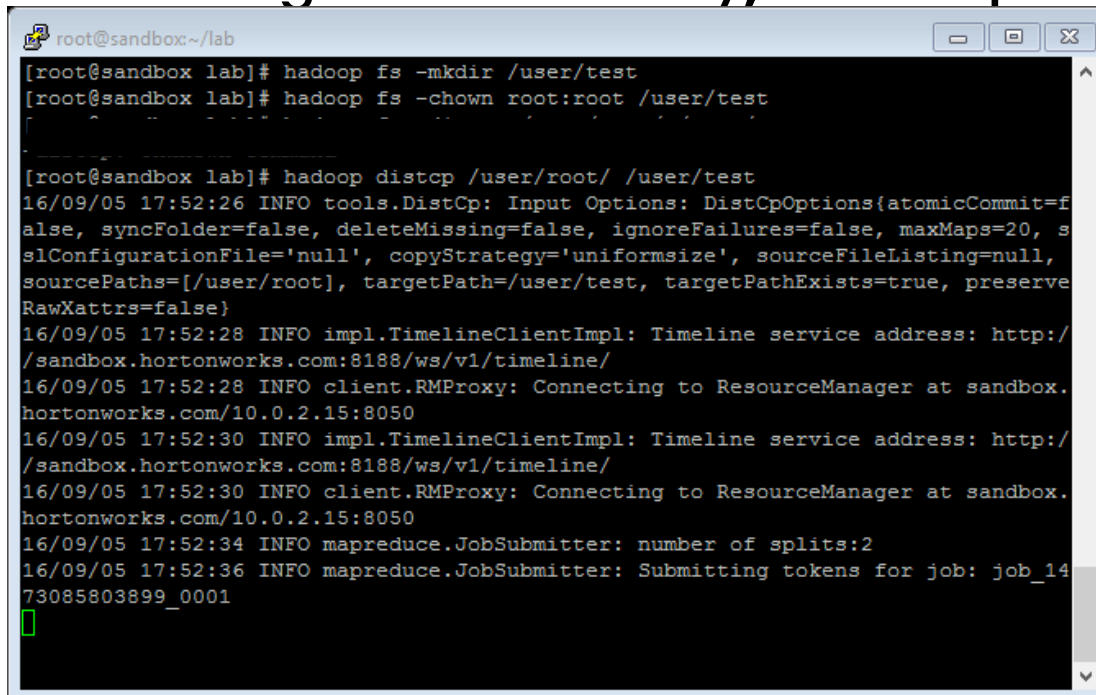
- ❑ “-getmerge” Takes a source directory file or files as input and concatenates files in src into the local destination file



```
root@sandbox:~/lab
[root@sandbox lab]# hadoop fs -put test_copy.txt /user/root
[root@sandbox lab]# hadoop fs -ls /user/root
Found 3 items
drwx----- - root root          0 2016-09-05 17:54 /user/root/.staging
-rw-r--r--  3 root root        318 2016-09-05 17:43 /user/root/test.txt
-rw-r--r--  3 root root        318 2016-09-05 17:54 /user/root/test_copy.txt
[root@sandbox lab]# hadoop fs -getmerge /user/root test_two_copies.txt
[root@sandbox lab]# ls
test_copy.txt  test_two_copies.txt  test.txt
[root@sandbox lab]# wc -l test.txt
43 test.txt
[root@sandbox lab]# wc -l test_copy.txt
43 test_copy.txt
[root@sandbox lab]# wc -l test_two_copies.txt
86 test_two_copies.txt
[root@sandbox lab]#
```

Step 6: Advanced HDFS feature

- ❑ “-distcp”
- ❑ It is a tool used for large inter/intra-cluster copying
- ❑ It uses MapReduce to effect its distribution copy, error handling and recovery, and reporting

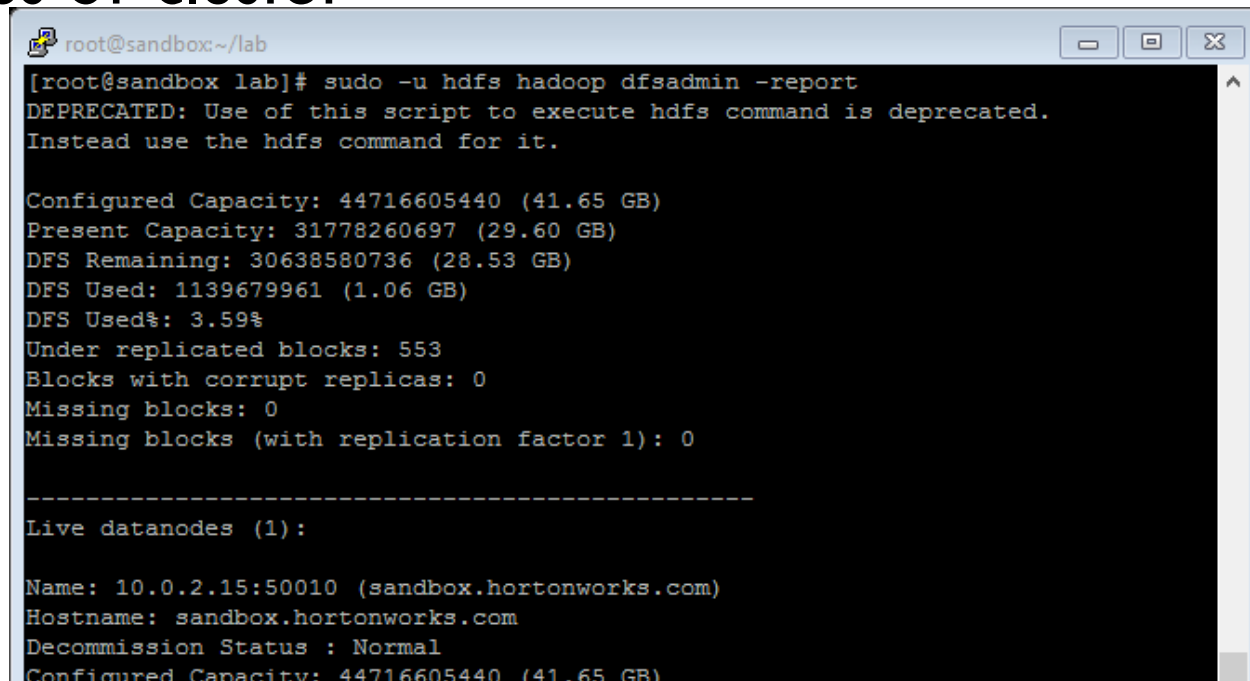


```
root@sandbox ~/lab
[root@sandbox lab]# hadoop fs -mkdir /user/test
[root@sandbox lab]# hadoop fs -chown root:root /user/test

[root@sandbox lab]# hadoop distcp /user/root/ /user/test
16/09/05 17:52:26 INFO tools.DistCp: Input Options: DistCpOptions{atomicCommit=false, syncFolder=false, deleteMissing=false, ignoreFailures=false, maxMaps=20, sslConfigurationFile='null', copyStrategy='uniformsize', sourceFileListing=null, sourcePaths=[/user/root], targetPath=/user/test, targetPathExists=true, preserveRawXattrs=false}
16/09/05 17:52:28 INFO impl.TimelineClientImpl: Timeline service address: http://sandbox.hortonworks.com:8188/ws/v1/timeline/
16/09/05 17:52:28 INFO client.RMProxy: Connecting to ResourceManager at sandbox.hortonworks.com/10.0.2.15:8050
16/09/05 17:52:30 INFO impl.TimelineClientImpl: Timeline service address: http://sandbox.hortonworks.com:8188/ws/v1/timeline/
16/09/05 17:52:30 INFO client.RMProxy: Connecting to ResourceManager at sandbox.hortonworks.com/10.0.2.15:8050
16/09/05 17:52:34 INFO mapreduce.JobSubmitter: number of splits:2
16/09/05 17:52:36 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1473085803899_0001
```

Step 6: Advanced HDFS feature

- ❑ “hadoop dfsadmin -report”
- ❑ To understand storage availability in the cluster
- ❑ Remember: user “hdfs” can view status of storage status of cluster



```
root@sandbox:~/lab
[root@sandbox lab]# sudo -u hdfs hadoop dfsadmin -report
DEPRECATED: Use of this script to execute hdfs command is deprecated.
Instead use the hdfs command for it.

Configured Capacity: 44716605440 (41.65 GB)
Present Capacity: 31778260697 (29.60 GB)
DFS Remaining: 30638580736 (28.53 GB)
DFS Used: 1139679961 (1.06 GB)
DFS Used%: 3.59%
Under replicated blocks: 553
Blocks with corrupt replicas: 0
Missing blocks: 0
Missing blocks (with replication factor 1): 0

-----
Live datanodes (1):

Name: 10.0.2.15:50010 (sandbox.hortonworks.com)
Hostname: sandbox.hortonworks.com
Decommission Status : Normal
Configured Capacity: 44716605440 (41.65 GB)
```


Try yourself

- Print the last 2 lines of the Shakespeare_100.txt file
- Create another HDFS directory called “labs_temp”
- Move the “Shakespeare_100.txt” to labs_temp

Note: Look at the “Hadoop HDFS Shell Commands” link on the D2L