# Hadoop Introduction - Part III

**Goal:** In this tutorial you will set up a distributed, multi-node Hadoop cluster and run a Hadoop MapReduce [WordCount](#) example job.

## Set up two single-node Hadoop clusters

- Log into [NCSU V](#)CL using your Unity ID, and reserve two "**CentOS 5.9 Base (64 bit VM)**" virtual machines. Use a SSH client (e.g., [Putty](#)) to connect to the virtual machines.

- Download and install JDK in both machines as described in Part I of the tutorial.

- Download and install Hadoop and configure a single-node Hadoop cluster in both machines as described in Part I of the tutorial. <u>Do not start the single-node Hadoop clusters</u>.

  - **NOTE:** it is best to NOT reuse the nodes you used in Part 1/Part 2 and instead use/configure two new machines. Otherwise, to resuse the same machine, stop the currently running cluster (`bin/stop-mapred.sh and bin/stop-dfs.sh`) and issue the following commands:

    - `rm -Rf hadoop/tmp/*`

    - `bin/hadoop namenode -format`

## Set up a multi-node Hadoop cluster

- After setting up the two single-node Hadoop clusters, you will modify the Hadoop configuration to make one cluster (e.g., `152.xxx.xxx.xxx`) the master and the other cluster (e.g., `152.yyy.yyy.yyy`) a slave.

- Disable the `iptables` firewall on both machines.

```
[user@vcl-master ~]$ sudo /etc/init.d/iptables stop
Flushing firewall rules:                                [  OK  ]
Setting chains to policy ACCEPT: filter                 [  OK  ]
Unloading iptables modules:                             [  OK  ]
```

```
[user@vcl-slave ~]$ sudo /etc/init.d/iptables stop
Flushing firewall rules:                                [  OK  ]
Setting chains to policy ACCEPT: filter                 [  OK  ]
Unloading iptables modules:                             [  OK  ]
```

- Set up the SSH RSA keys, so that the master machine can communicate with the slave machine without entering a password. Substitute `user` for your Unity ID and `152.yyy.yyy.yyy` for the IP address of the slave machine. Enter your password when prompted.

```
[user@vcl-master ~]$ ssh-copy-id -i $HOME/.ssh/id_rsa.pub
user@152.yyy.yyy.yyy

30

The authenticity of host '152.yyy.yyy.yyy (152.yyy.yyy.yyy)'
can't be established.

RSA key fingerprint is

[…]

Are you sure you want to continue connecting (yes/no)? yes

Warning: Permanently added '152.yyy.yyy.yyy' (RSA) to the list of
known hosts.

user@152.yyy.yyy.yyy's password:

Now try logging into the machine, with "ssh
'user@152.yyy.yyy.yyy'", and check in:



  .ssh/authorized_keys



to make sure we haven't added extra keys that you weren't
expecting.
```

- Test the SSH setup by connecting from the master machine (`152.xxx.xxx.xxx`) to the slave machine (`152.yyy.yyy.yyy`).

```
[user@vcl-master ~]$ ssh 152.yyy.yyy.yyy
```

- Connect again to the master machine using a SSH client.

- On the master machine, open file "**conf/slaves**." If you followed Part I of the tutorial, this file should contain the IP address of the master machine (`152.xxx.xxx.xxx`). Add a new line with the IP address of the slave machine (`152.yyy.yyy.yyy`).

```
152.xxx.xxx.xxx
152.yyy.yyy.yyy
```

- On both machines, open file "**conf/core-site.xml**" and change the `fs.default.name` parameter, which specifies the NameNode host and port, by substituting `localhost` for the IP address of the master machine (`152.xxx.xxx.xxx`).

```
<property>
    <name>fs.default.name</name>
    <value>hdfs://152.xxx.xxx.xxx:54310</value>
    <description>The name of the default file system.  A URI
    whose scheme and authority determine the FileSystem
    implementation. The uri's scheme determines the config
    property (fs.SCHEME.impl) naming the FileSystem
    implementation class.  The uri's authority is used to
    determine the host, port, etc. for a filesystem.
    </description>
</property>
```

- On both machines, open file "**conf/mapred-site.xml**" and change the `mapred.job.tracker` parameter, which specifies the JobTracker host and port, by substituting `localhost` for the IP address of the master machine (`152.xxx.xxx.xxx`).

```
<property>
    <name>mapred.job.tracker</name>
    <value>152.xxx.xxx.xxx:54311</value>
    <description>The host and port that the MapReduce job
    tracker runs at. If "local", then jobs are run in-process as
    a single map and reduce task.
    </description>
</property>
```

- On both machines, open file "**conf/hdfs-site.xml**" and change the value of the `dfs.replication` parameter to `2`.

```
<property>
    <name>dfs.replication</name>
    <value>2</value>
    <description>Default block replication. The actual number of
    replications can be specified when the file is created. The
    default is used if replication is not specified in create
    time.
    </description>
</property>
```

## Start the multi-node Hadoop cluster

- On the master machine, clear the "**tmp**" directory and then use command `namenode` to format the HDFS filesystem.

```
[user@vcl-master ~]$ rm -Rf hadoop/tmp/*
[user@vcl-master ~]$ cd $HOME/hadoop/hadoop-1.2.1
[user@vcl-master hadoop-1.2.1]$ bin/hadoop namenode -format
[…] INFO namenode.NameNode: STARTUP_MSG:
/************************************************************
STARTUP_MSG: Starting NameNode
STARTUP_MSG:   host = bnxxx-xxx.dcs.mcnc.org/152.xxx.xxx.xxx
STARTUP_MSG:   args = [-format]
STARTUP_MSG:   version = 1.2.1
STARTUP_MSG:   build =
https://svn.apache.org/repos/asf/hadoop/common/branches/branch-
1.2 -r 1503152; compiled by 'mattf' on Mon Jul 22 15:23:09 PDT
2013
STARTUP_MSG:   java = 1.7.0_25
************************************************************/
[…]
/************************************************************
SHUTDOWN_MSG: Shutting down NameNode at bnxxx-
xxx.dcs.mcnc.org/152.xxx.xxx.xxx
************************************************************/
```

- On the master machine, use command `start-dfs.sh` to start the HDFS daemons. This command will start up the `NameNode` on the master machine and the `DataNode` on the master and the slave machines.

```
[user@vcl-master hadoop-1.2.1]$ bin/start-dfs.sh
starting namenode, logging to /home/user/hadoop/hadoop-
1.2.1/libexec/../logs/hadoop-user-namenode-bnxxx-
xxx.dcs.mcnc.org.out
152.xxx.xxx.xxx: starting datanode, logging to
/home/user/hadoop/hadoop-1.2.1/libexec/../logs/hadoop-user-
datanode-bnxxx-xxx.dcs.mcnc.org.out
152.yyy.yyy.yyy: starting datanode, logging to
/home/user/hadoop/hadoop-1.2.1/libexec/../logs/hadoop-user-
datanode-bnyyy-yyy.dcs.mcnc.org.out
152.xxx.xxx.xxx: starting secondarynamenode, logging to
/home/user/hadoop/hadoop-1.2.1/libexec/../logs/hadoop-user-
secondarynamenode-bnxxx-xxx.dcs.mcnc.org.out
```

- Verify that the `NameNode` and the `DataNode` are running on the corresponding machines using command `jps`.

```
[user@vcl-master hadoop-1.2.1]$ jps
24564 DataNode
24440 NameNode
24759 Jps
24704 SecondaryNameNode
```

```
[user@vcl-slave ~]$ jps
24309 DataNode
24389 Jps
```

- On the master machine, use command `start-mapred.sh` to start the MapReduce daemons. This command will start up the `JobTracker` on the master machine and the `TaskTracker` on the master and the slave machines.

```
[user@vcl-master hadoop-1.2.1]$ bin/start-mapred.sh
starting jobtracker, logging to /home/user/hadoop/hadoop-
1.2.1/libexec/../logs/hadoop-user-jobtracker-bnxxx-
xxx.dcs.mcnc.org.out
152.yyy.yyy.yyy: starting tasktracker, logging to
/home/user/hadoop/hadoop-1.2.1/libexec/../logs/hadoop-user-
tasktracker-bnyyy-yyy.dcs.mcnc.org.out
152.xxx.xxx.xxx: starting tasktracker, logging to
/home/user/hadoop/hadoop-1.2.1/libexec/../logs/hadoop-user-
tasktracker-bnxxx-xxx.dcs.mcnc.org.out
```

- Verify that the `JobTracker` and the `TaskTracker` are running on the corresponding machines using command `jps`.

```
[user@vcl-master hadoop-1.2.1]$ jps
25049 TaskTracker
24564 DataNode
24928 JobTracker
24440 NameNode
24704 SecondaryNameNode
25150 Jps
```

```
[user@vcl-slave ~]$ jps
24309 DataNode
24545 TaskTracker
24616 Jps
```

## Download example input data

- Download the following ebooks from Project Gutenberg as text files in "**Plain Text UTF-8**" encoding:

  o [A Tale of Two Cities by Charles Dickens](#).

  o [Les Misérables by Victor Hugo](#).

  o [Pride and Prejudice by Jane Austen](#).

  o [Ulysses by James Joyce](#).

  Create a directory named "**gutenberg**" in the home directory of the master machine and store these files.

```
[user@vcl-master hadoop-1.2.1]$ cd $HOME
[user@vcl-master ~]$ mkdir gutenberg
```

- On the master machine, copy the files to the HDFS.

```
[user@vcl-master ~]$ cd $HOME/hadoop/hadoop-1.2.1/

[user@vcl-master hadoop-1.2.1]$ bin/hadoop dfs -copyFromLocal
$HOME/gutenberg /gutenberg
[user@vcl-master hadoop-1.2.1]$ bin/hadoop dfs -ls /gutenberg
Found 4 items
-rw-r--r--   2 user supergroup 717569  […] /gutenberg/pg1342.txt
```

```
-rw-r--r--    2 user supergroup 3322647 […] /gutenberg/pg135.txt
-rw-r--r--    2 user supergroup 1573150 […] /gutenberg/pg4300.txt
-rw-r--r--    2 user supergroup 792927  […] /gutenberg/pg98.txt
```

### Run WordCount job

- On the master machine, run the WordCount example job.

```
[user@vcl-master hadoop-1.2.1]$ bin/hadoop jar
hadoop*examples*.jar wordcount /gutenberg /gutenberg-output

[…] INFO input.FileInputFormat: Total input paths to process : 4

[…] INFO util.NativeCodeLoader: Loaded the native-hadoop library

[…] WARN snappy.LoadSnappy: Snappy native library not loaded

[…] INFO mapred.JobClient: Running job: job_xxx

[…] INFO mapred.JobClient:  map 0% reduce 0%

[…] INFO mapred.JobClient:  map 50% reduce 0%

[…] INFO mapred.JobClient:  map 98% reduce 0%

[…] INFO mapred.JobClient:  map 100% reduce 0%

[…] INFO mapred.JobClient:  map 100% reduce 33%

[…] INFO mapred.JobClient:  map 100% reduce 100%

[…] INFO mapred.JobClient: Job complete: job_xxx

[…]
```

- On the master machine, copy the output file of the WordCount example job from the HDFS to the local file system.

```
[user@vcl-master hadoop-1.2.1]$ bin/hadoop dfs -getmerge
/gutenberg-output $HOME/gutenberg-output
```

Open file "**gutenberg-output**" in folder "**~/gutenberg-output**" on the master machine. Output file must look as follows:

```
"            27

"'A        2

"'After    1

"'At       1

"'Do       2

[…]
```

- Submit file "**gutenberg-output**" through Moodle (rename the file "**hadoop-multi-node-user**," where **user** is your Unity ID).

## Stop the multi-node Hadoop cluster

- On the master machine, use commands `stop-mapred.sh` and `stop-dfs.sh` to stop the MapReduce and the HDFS daemons, respectively.

```
[user@vcl-master hadoop-1.2.1]$ bin/stop-mapred.sh

stopping jobtracker

152.yyy.yyy.yyy: stopping tasktracker

152.xxx.xxx.xxx: stopping tasktracker

[user@vcl-master hadoop-1.2.1]$ bin/stop-dfs.sh

stopping namenode

152.xxx.xxx.xxx: stopping datanode
```

```
152.yyy.yyy.yyy: stopping datanode

152.xxx.xxx.xxx: stopping secondarynamenode
```

## References

- [http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-single-node-cluster/](http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-single-node-cluster/)

- [http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-multi-node-cluster/](http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-multi-node-cluster/)