# Hadoop Installation

## Sandeep Prasad

## 1 Introduction

Hadoop is a system to manage large quantity of data. For this report hadoop-1.0.3 (Released, May 2012) is used and tested on Ubuntu-12.04. The system configuration is Memory (RAM) 4GB, Processor Intel®Core™i3-2120 CPU @3.30GHz X 4, OS Type 32 bit. The installation of hadoop in this [1] installation report is given as

1. Prerequisites
   (a) Java
   (b) Dedicated user
   (c) Configuring ssh
   (d) Disabling ipv6

2. Installation
   (a) .bashrc
   (b) Changes in hadoop-env.sh and *-site.xml file
   (c) Formatting hdfs
   (d) Starting and stopping single-node cluster

## 2 Prerequisites

### 2.1 Java

Hadoop requires Java 1.5 or above but all the Tutorials available on web insist on Java 1.6[1][2] and above. For this installation manual Java 1.7.0_25 is used. Java 1.7 is available in Ubuntu repository and can be installed using command given in Listing 1

```
1  $ sudo apt−get install openjdk−7−jdk
```

Listing 1: Installing Java 1.7

---

[1]http://hadoop.apache.org/docs/stable/single_node_setup.html#Required+Software

Java version can be checked using command in Listing 2 output of command
shown in Figure 1

```
1  $ java −version
```
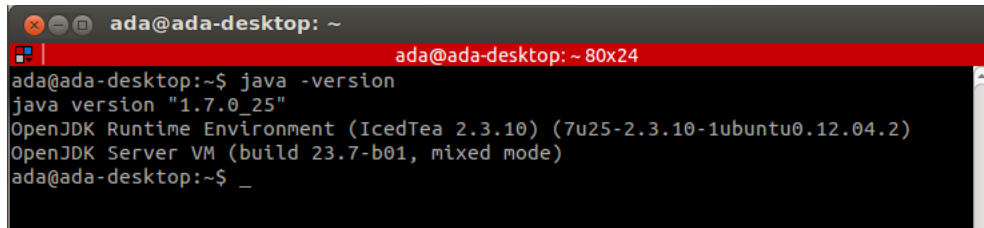Listing 2: Checking Java Version



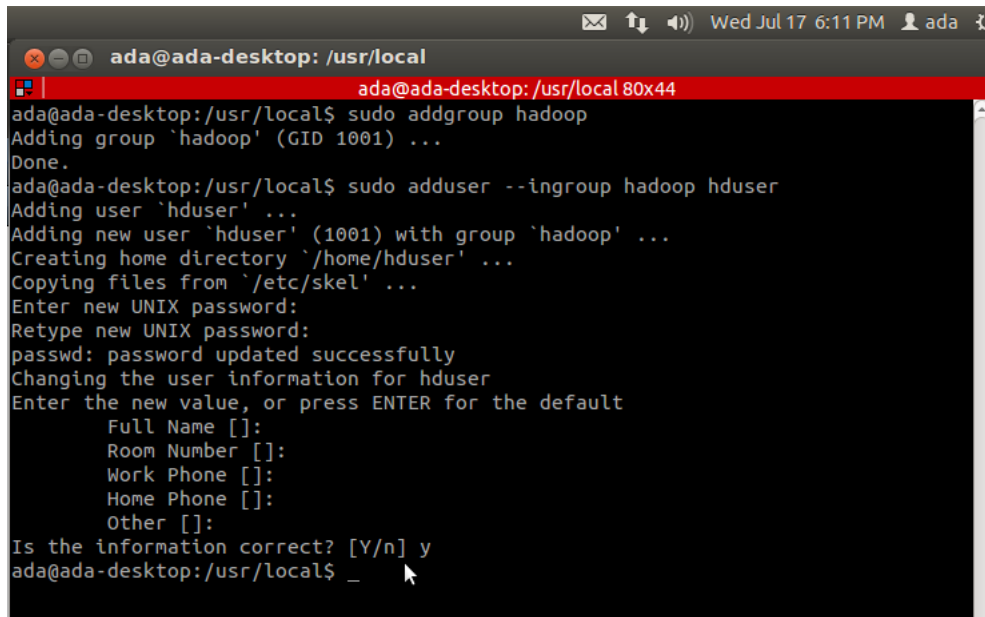Figure 1: Hadoop requires Java 1.5 or higher

## 2.2   Creating dedicated user

Tutorials visited on internet advise creating a new dedicated user for using
hadoop. New group is created (in this installation report new group created is
"hadoop") and user (in this installation report new user added is "hduser") can
be added to the newly created group using commands in Listing 3

```
1  $ sudo addgroup hadoop
2  $ sudo adduser −−ingroup hadoop hduser
```
Listing 3: adding group and user for hadoop

Figure 2 displays the above mentioned commands 'for creating group and
user' executed on my system. When hduser is added it asks for new UNIX
password. This password is password for hduser. Retype the password when
prompted and enter the details asked (details are optional). In the end enter 'y'
to complete the procedure.

Figure 2: Adding group hadoop and user hduser

Some steps mentioned in this manual require sudo permission. hduser can be added to sudo list using command mentioned in Listing 4.

```
1  $ sudo adduser hduser sudo
```

Listing 4: adding hduser to sudoers list

## 2.3   Configuring ssh

Ssh access is required for hadoop to run. In this installation report we will configure ssh access for localhost to user hduser. If ssh server is not installed on the machine, for Ubuntu it can be installed using command given in Listing 5

```
1  $ sudo apt−get install openssh−server
```

Listing 5: Installing ssh server

To allow ssh access a SSHKey has to be generated which can be generated for user hduser as followed

```
1  $ su − hduser
2  $ ssh−keygen −t rsa −P ""
3  $ cat $HOME/.ssh.id_rsa.pub >> $HOME/.ssh/authorized_keys
4  $ ssh localhost
```

Listing 6: Creating keygen and adding localhost to known hosts

3

The command given above can be explained as

1. Changing from default user to hduser, given in line 1 of Listing 6.

2. Generating keygen, when asked to enter the file to save the key, press enter and key will be saved in default /home/hduser/.ssh/id_rsa file, given in line 2 of Listing 6.

3. Authorizing public key generated as in line 3 of Listing 6.

4. Adding localhost to list of known hosts using ssh, when prompted for 'yes/no', write 'yes' and press enter, given in line 4 of Listing 6.

5. all the above steps is carried out by hduser.

Figure 3 shows the configuration steps for ssh executed on my system.

Figure 3: Configuring ssh on localhost

## 2.4 ipv6

For hadoop to run ipv6 has to be disabled which can be done by editing /etc/sysctl.conf file. Editing sysctl.conf file requires sudo permission. Lines added in /etc/sysctl.conf file is shown in Listing 7

```
1  #disabling ipv6
2  net.ipv6.conf.all.disable_ipv6 = 1
3  net.ipv6.conf.default.disable_ipv6 = 1
4  net.ipv6.conf.lo.disable_ipv6 = 1
```
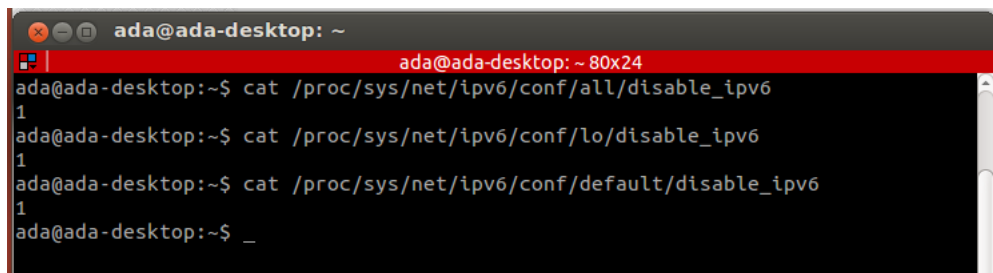
Listing 7: Lines added in /etc/sysctl.conf file

Ipv6 status can be checked using command in listing 8.

```
1  $ cat /proc/sys/net/ipv6/conf/all/disable_ipv6
```

Listing 8: Checking ipv6 status after restarting system

The output will be either 0 or 1. 0 means ipv6 is enabled and 1 means it is disabled as shown in figure 4



Figure 4: Checking status of ipv6 after restarting computer

Alternatively ipv6 can be disabled only for hadoop by adding line below in /usr/local/hadoop/conf/hadoop-env.sh.

```
1  export HADOOP_OPTS=-Djava.net.preferIPv4Stack=true
```

# 3   Installation

## 3.1   Hadoop's folder

Copy hadoop-1.0.3.tar.gz file in /usr/local directory and untar it. Also create temporary folder which will be used by hadoop's hdfs file system (in my case temporary folder created is 'tmp' in /usr/local folder). After that we have to change ownership of hadoop and temporary directory just created. Copying the file in /usr/local, untaring, creating temporary folder and changing owner requires sudo permission. The commands executed is given in figure 5

```
1  $ sudo tar −xzf hadoop −1.0.3.tar.gz
2  $ sudo mv hadoop −1.0.3  hadoop
3  $ sudo mkdir tmp
4  $ sudo chown −R hduser:hadoop hadoop
5  $ sudo chown −R hduser:hadoop tmp
```

Listing 9: Steps to be followed before using hadoop

The steps mentioned in Listing 9 assumes that hadoop's tar file has been copied in /usr/local folder and user with sudo permission is in /usr/local folder (check the working folder using 'pwd' command on terminal, now the steps can be explained as below

1. untar hadoop-1.0.3.tar.gz file line 1 of Listing 9. It will create a folder called hadoop-1.0.3.

2. line 2 of Listing 9 changes the name of hadoop folder from hadoop-1.0.3 to hadoop. This step is not required but is carried out as convenience.

3. Line 3 of Listing 9 makes 'tmp' directory that will be used by hdfs as it's temporary folder and it's location will be mentioned in core-site.xml file.

4. Line 4 and line 5 of Listing 9 changes the ownership of hadoop and tmp folder from root:root to hduser:hadoop.

Figure 5: Steps executed with sudo user

## 3.2 Updating .bashrc for hduser

We can edit .bashrc file for hduser. The edited .bashrc file is shown in Listing 10

```
1  #Set Hadoop−related environment variables
2  export HADOOP_HOME=/usr/local/hadoop
3
4  #Set JAVA_HOME= (we will also configure JAVA_HOME directly for
       Hadoop later on)
5  export JAVA_HOME=/usr/lib/jvm/java−7−openjdk−i386
6
7  #Some convenient aliases and functions for running Hadoop−related
       commands
8  unalias fs &> /dev/null
9  alias fs="hadoop fs"
10 unalias hls &> /dev/null
11 alias hls="fs −ls"
12
13 #If you have LZO compression enabled in your hadoop cluster and
14 #compress job outputs with LZOP (not covered in this tutorial)
15 #Conveniently inspect an LZOP compressed file from the command
16 #line: run via:
17 #
18 # $ lzohead /hdfs/path/to/lzop/compressed/file.zo
19 #
20 #Requires installed 'lzop' command
21 lzohead () {
22         hadoop fs −cat $1 | lzop −dc | head −1000 | less
```

8

```
23  }
24
25  #Add Hadoop bin/ directory to PATH
26  export PATH=$PATH:$HADOOP_HOME/bin
```

<div align="center">Listing 10: Changes made in .bashrc file for hduser</div>

Hadoop uses lzop which is a compression tool. In Ubuntu lzop can be installed using command in Listing 11

```
1  $ sudo apt−get install lzop
```

<div align="center">Listing 11: Installing lzop in Ubuntu</div>

## 3.3 Changes in Hadoop folder

In Hadoop's folder we have to edit few files for hadoop to run. The files can be found in /usr/local/hadoop/conf directory. The files are hadoop-env.sh, core-site.xml, hdfs-site.xml and mapred-site.xml. This changes can be done using user 'hduser'.

### 3.3.1 hadoop-env.sh

In hadoop-env.sh we have to define path for JAVA_HOME. By default it will be commented and it's value will be set to j2sdk1.5-sun as shown in Listing 12, un-comment it and change it's value to the Java to be used. Original and edited hadoop-env.sh files are given in Listing 12 and Listing 13 respectively.

```
1  # The java implementation to use.   Required.
2  # export JAVA_HOME=/usr/lib/j2sdk1.5−sun
```

<div align="center">Listing 12: Java path in original hadoop-env.sh</div>

```
1  # The java implementation to use.   Required.
2    export JAVA_HOME=/usr/lib/jvm/java−7−openjdk−i386
```

<div align="center">Listing 13: Java path provided in hdfs-env.sh</div>

### 3.3.2 core-site.xml

Edited core-site.xml is given in Listing 14. Notice the <value> field in first <property> tag, the value points to 'tmp' folder we created earlier as mentioned in Line 3 of Listing 9.

```
1  <?xml version="1.0"?>
2  <?xml−stylesheet type="text/xsl" href="configuration.xsl"?>
3
4  <!−− Put site−specific property overrides in this file. −−>
5
6  <configuration>
7          <property>
8                  <name>hadoop.tmp.dir</name>
9                  <value>/usr/local/tmp</value>
10                 <description>A base for other temporary directories
                         .</description>
```

```
11            </property>
12
13            <property>
14                    <name>fs.default.name</name>
15                    <value>hdfs://localhost:54310</value>
16                    <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
                        >
17            </property>
18 </configuration>
```

Listing 14: Edited core-site.xml

### 3.3.3   hdfs-site.xml

Edited hdfs-site.xml file is given in Listing 15

```
1 <?xml version="1.0"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3
4 <!-- Put site-specific property overrides in this file. -->
5
6 <configuration>
7         <property>
8                 <name>dfs.replication</name>
9                 <value>1</value>
10                <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>
11        </property>
12 </configuration>
```

Listing 15: Edited hdfs-site.xml

### 3.3.4   mapred-site.xml

Edited mapred-site.xml file is given in Listing 16

```
1 <?xml version="1.0"?>
2 <?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
3
4 <!-- Put site-specific property overrides in this file. -->
5
6 <configuration>
7         <property>
8                 <name>mapred.job.tracker</name>
9                 <value>localhost:54311</value>
10                <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>
11        </property>
12 </configuration>
```
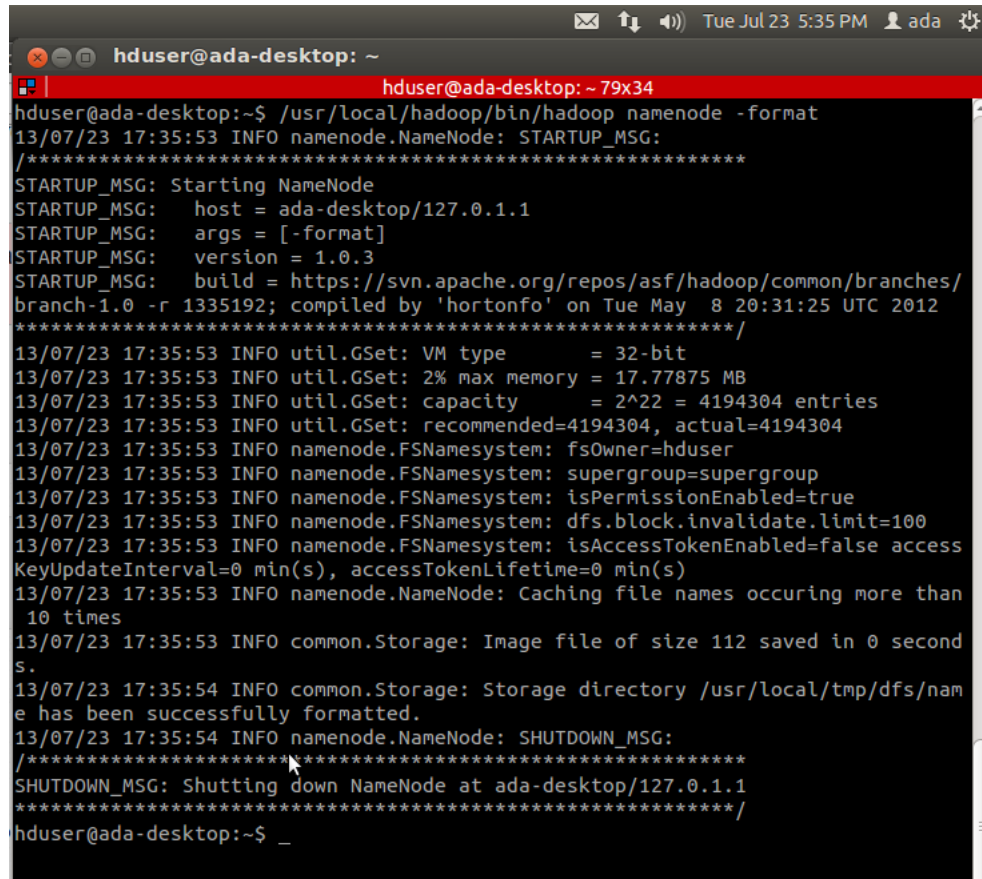
## 3.4   Formating hdfs FileSystem

Formatting hdfs FileSystem will format the virtually created File System. Anything stored in the cluster will be lost. hdfs can be formatted using command given in Listing 17. Figure 6 shows the the output obtained by formatting hdfs on my system.

```
1 $ /usr/local/hadoop/bin/hadoop namenode −format
```
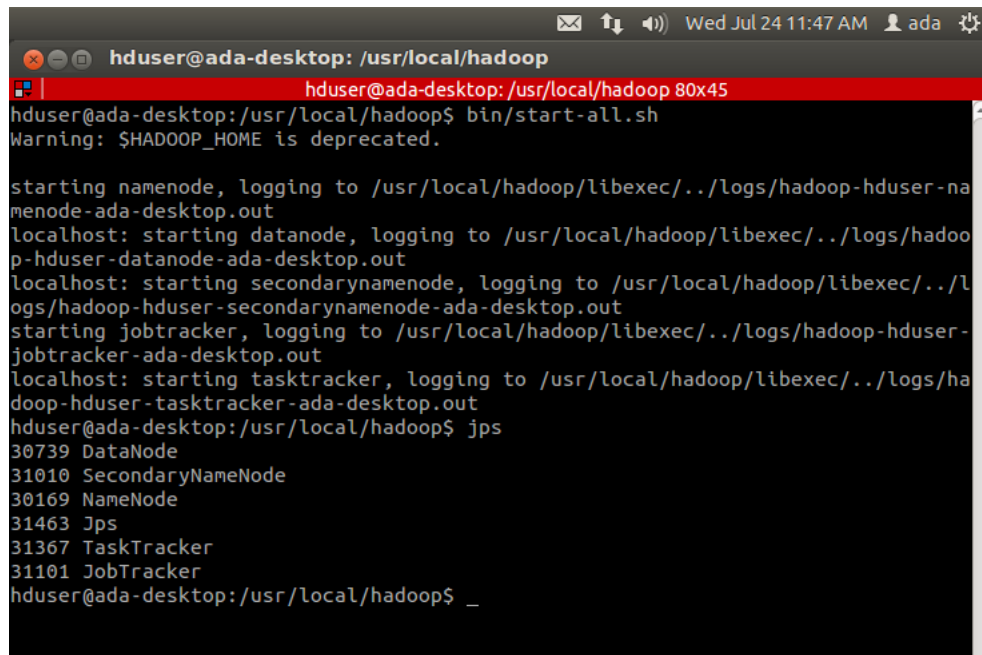
Listing 17: formatting hdfs



Figure 6: Output when hdfs is formatted

## 3.5   Starting and stopping hdfs

After completing all prerequisites, installation steps mentioned and formatting hdfs, hadoop is ready for use. Hadoop can be started and stopped using the start and stop script available in bin directory (done using hduser). Script to
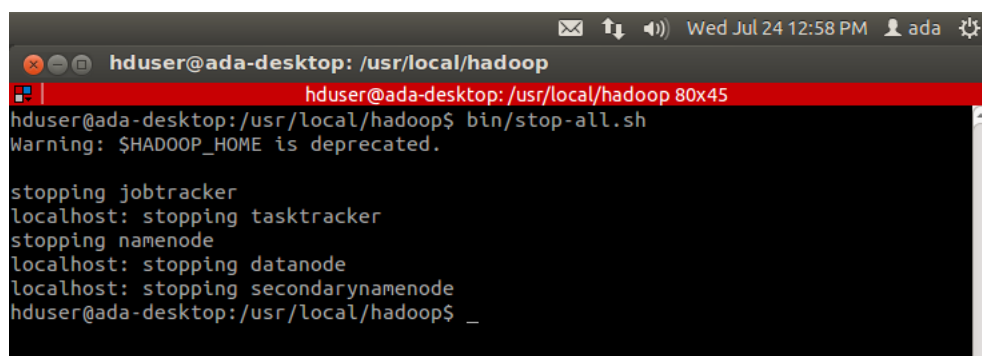
start and stop hadoop when run on my system are given in figure 7 and figure 8 respectively. The command to start hadoop services is (it is assumed you are in /usr/local/hadoop directory). Figure 7 also mentions jps, jps is a tool available in Java used to check the services started. When start script is executed the services started are DataNode, SecondaryNameNode, NameNode, TaskTracker and JobTracker.



Figure 7: Starting hadoop and checking the status of started processes using jps



Figure 8: Stopping hadoop processes

# References

[1] Michael G. Noll. Running hadoop on ubuntu linux (single-node cluster) - michael g. noll. http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-single-node-cluster/.

[2] Hadoop 1.1.2 Documentation. Single node setup. http://hadoop.apache.org/docs/stable/index.html.