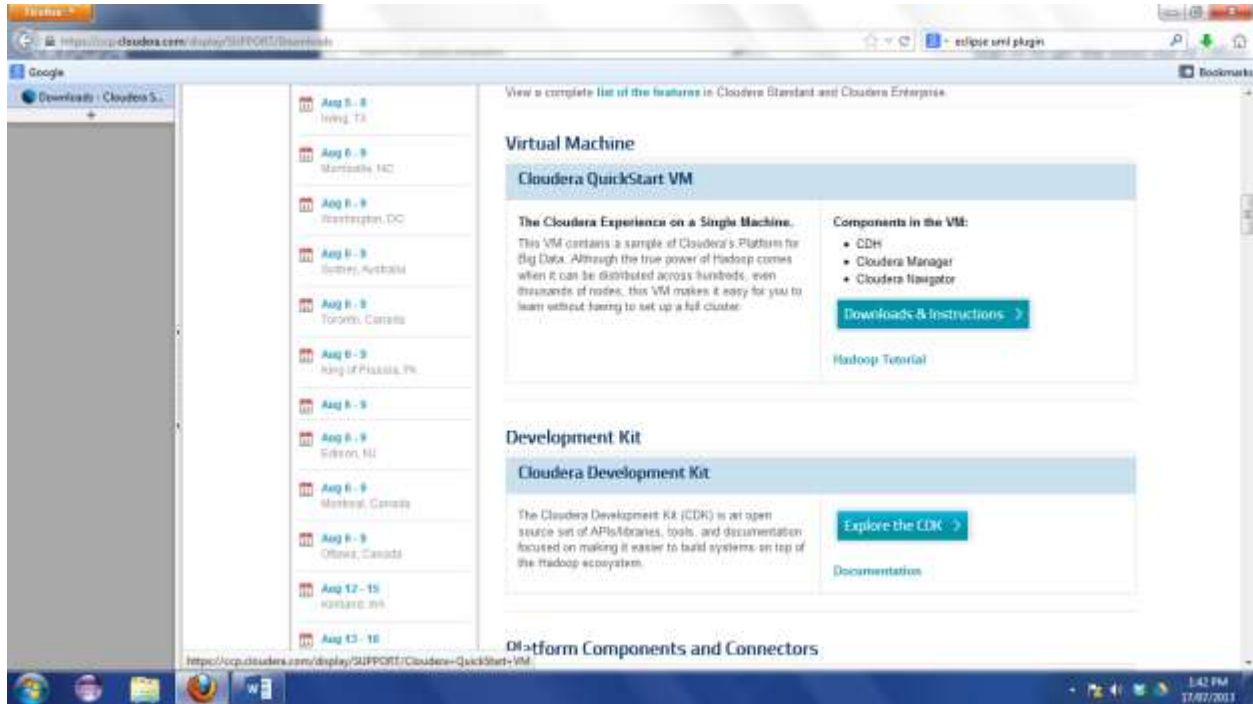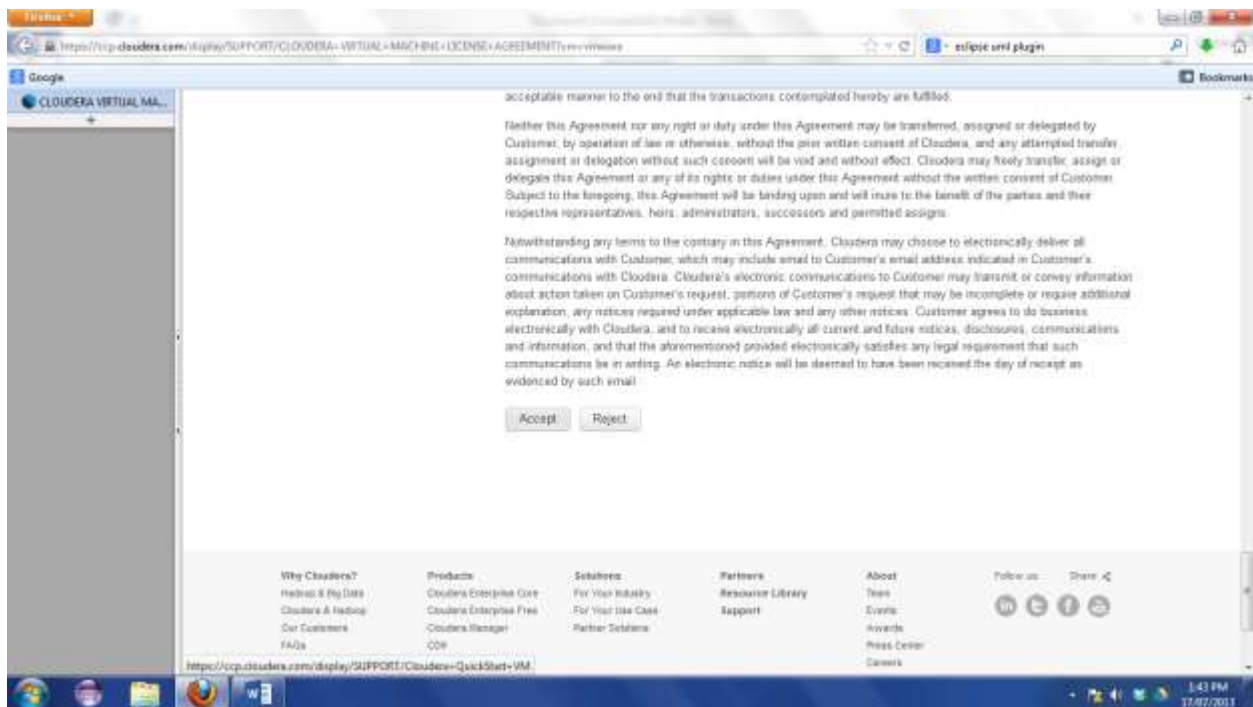# Cloudera Hadoop Installation and Configuration

1. Go to Cloudera Quickstart VM to download a pre-setup CDH virtual machine.
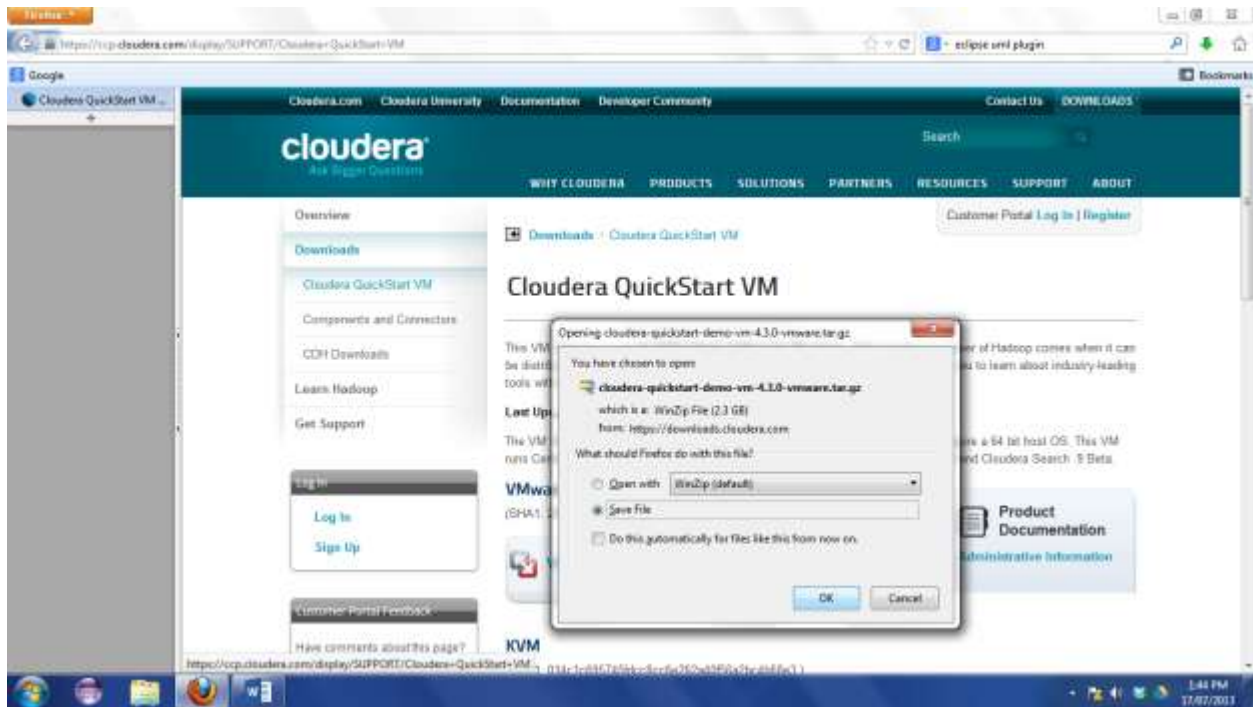


2. Select a VM you wish to download. For purpose of this assignment, I have used VMware Player.
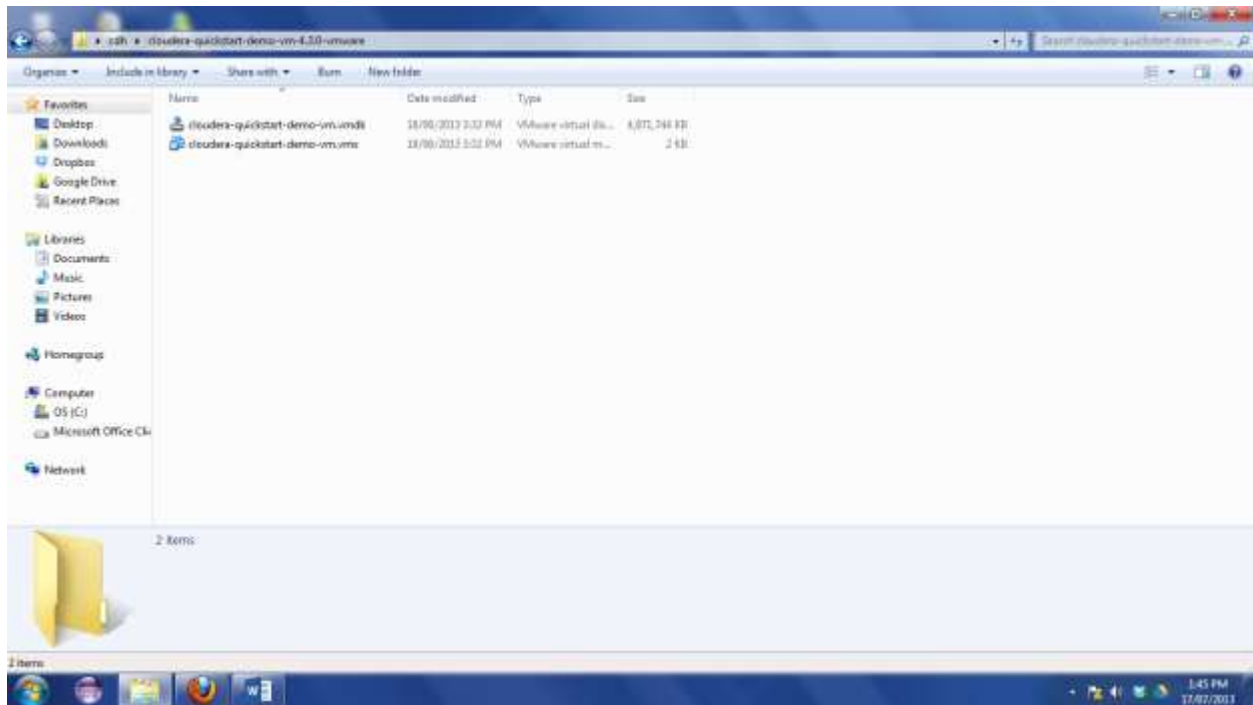
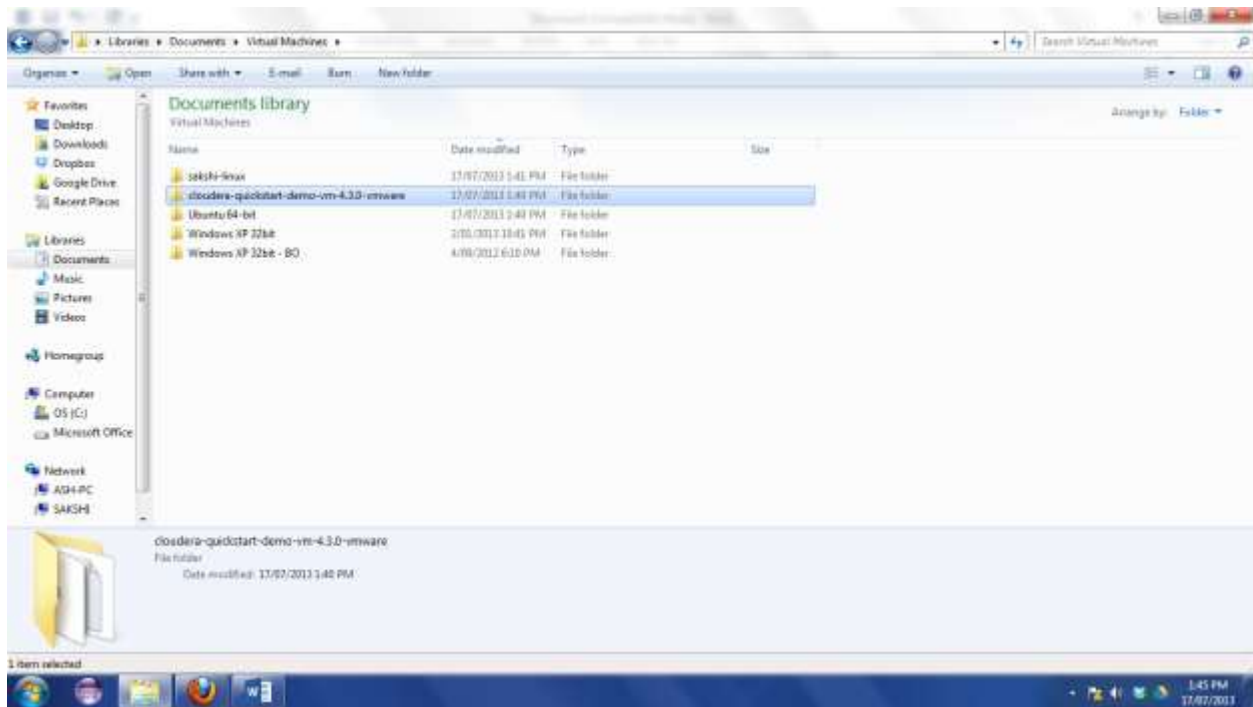3. To get download file, accept the agreement.
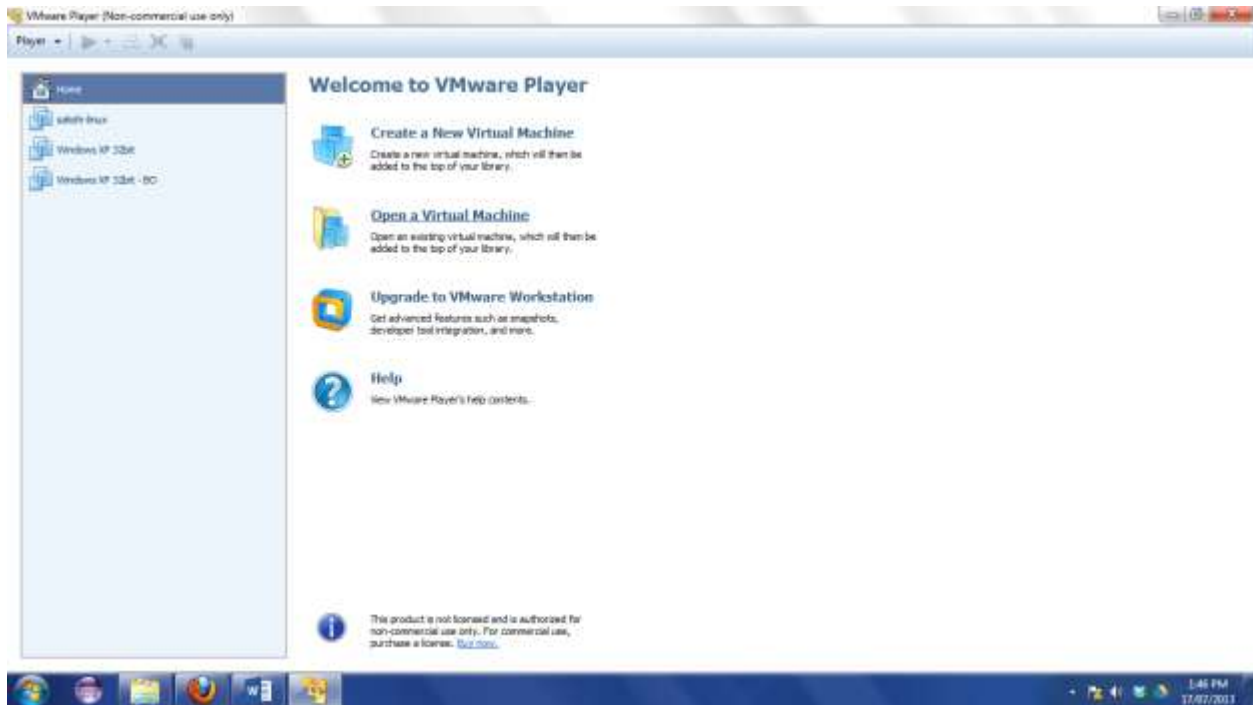
4. Save the downloadable file.



5. Unzip the downloaded file. You will get 2 files - .vmx (Virtual machine) and .vmdk
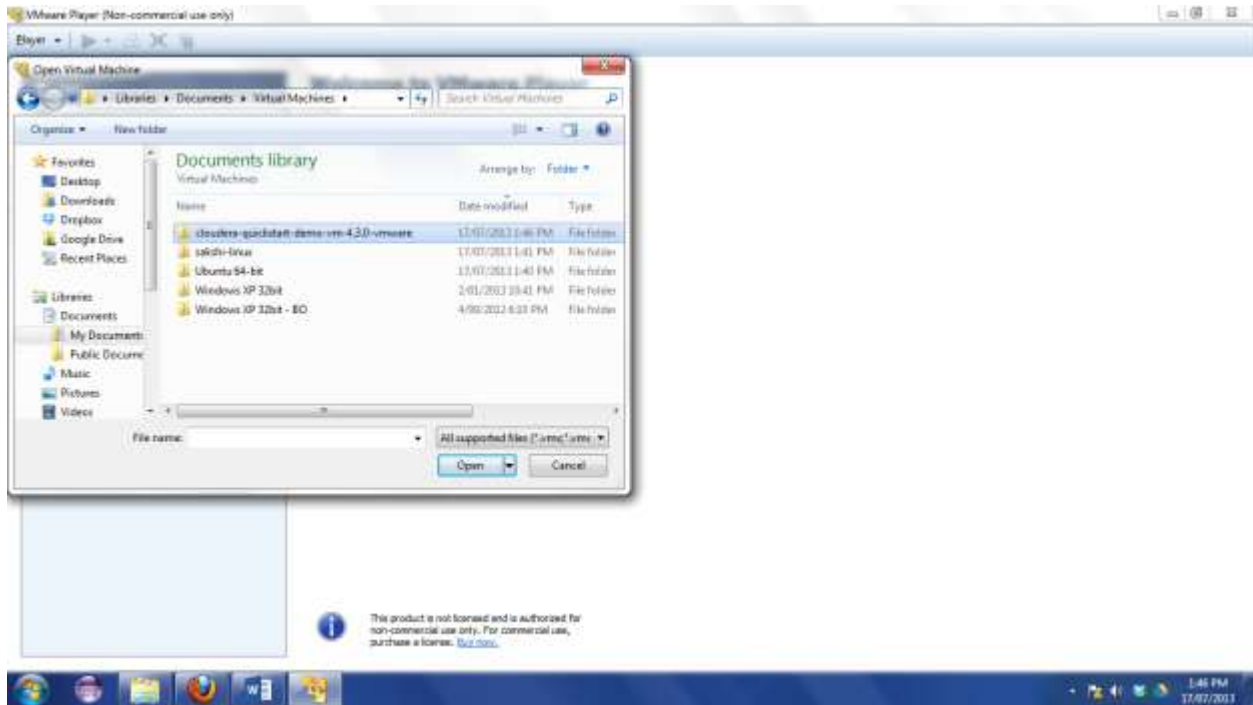
6. Save this folder in the directory where VMPlayer stores all virtual machine files.
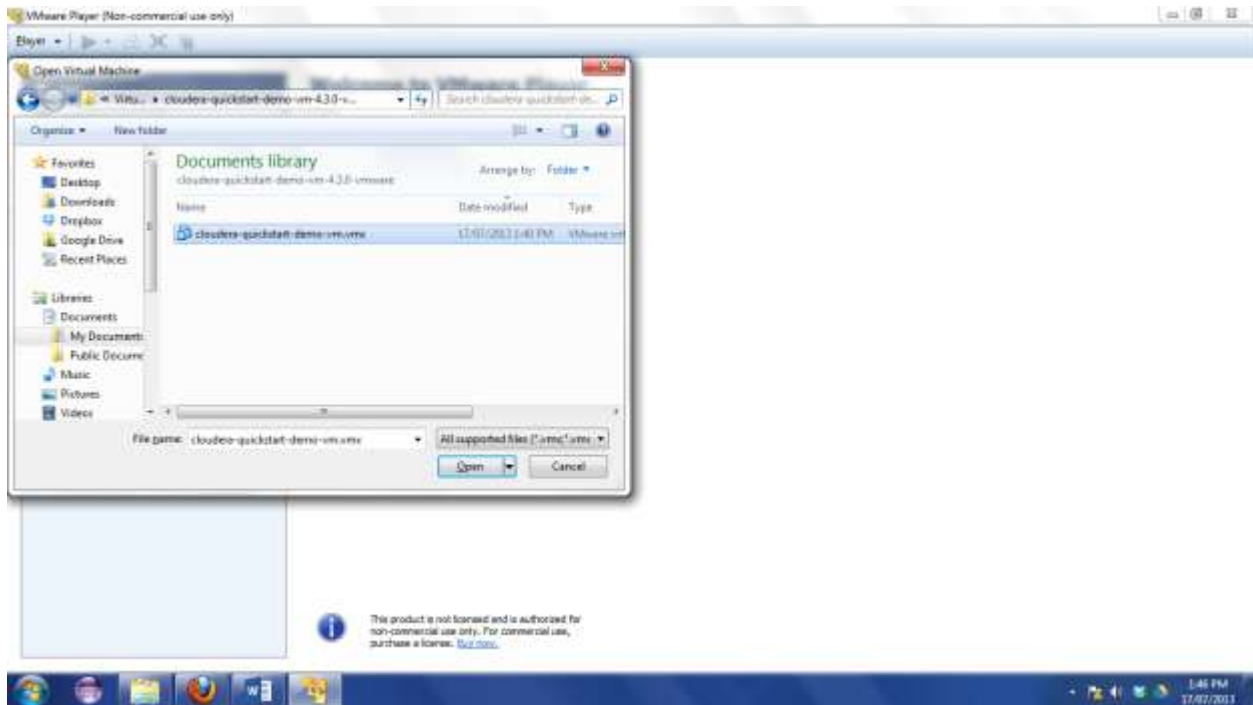


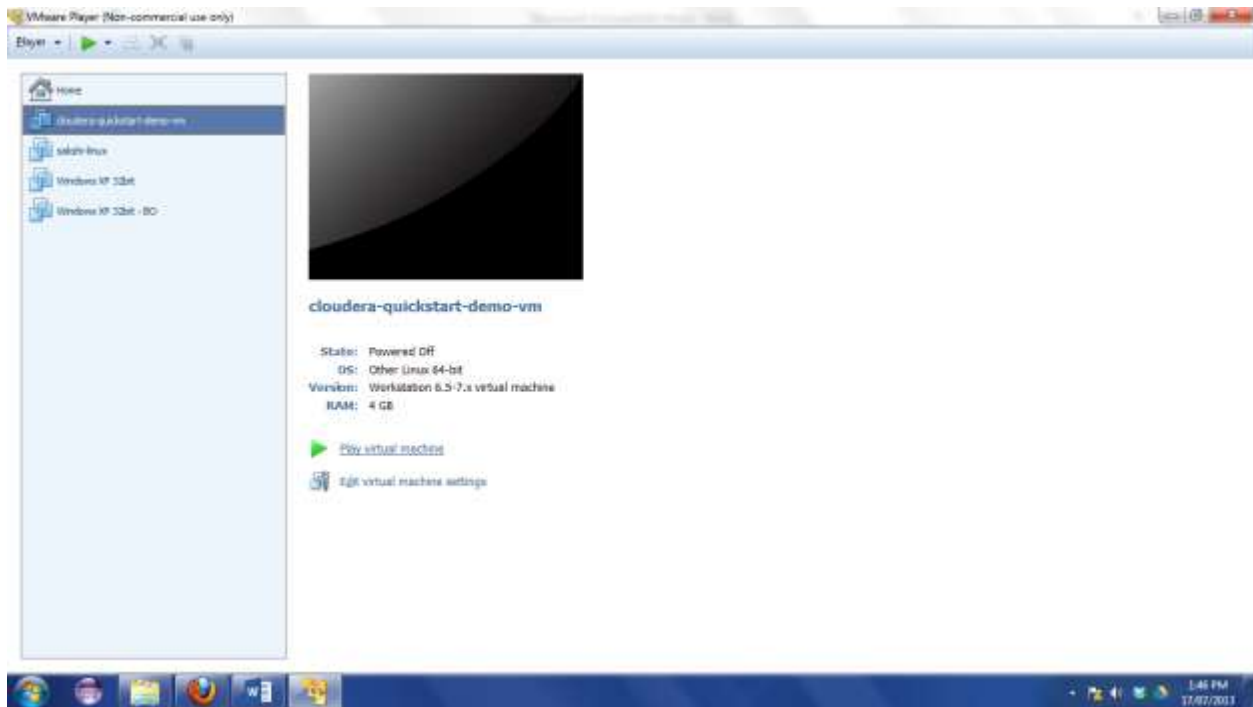7. Go to VMWare player and click on Open a Virtual Machine.
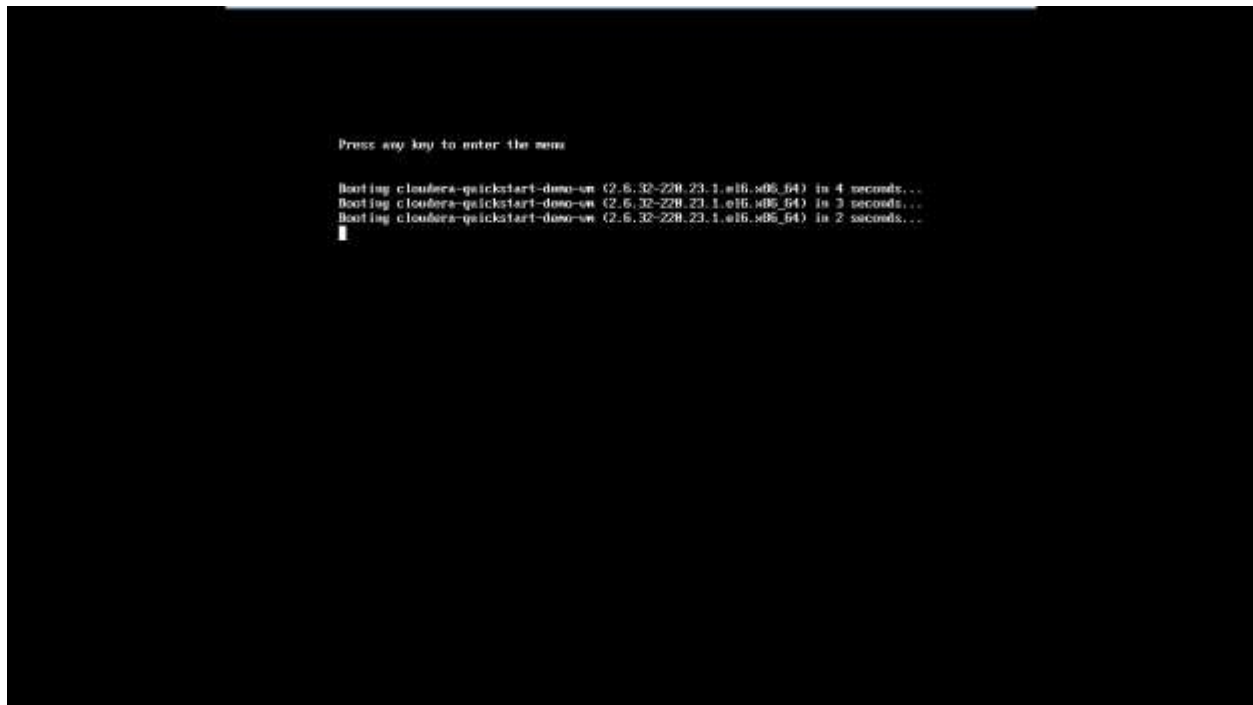
8. Select the newly copied virtual machine.



9. Select .vmx file which is the virtual machine file.

10. Once the VM is available, click on Play virtual machine.



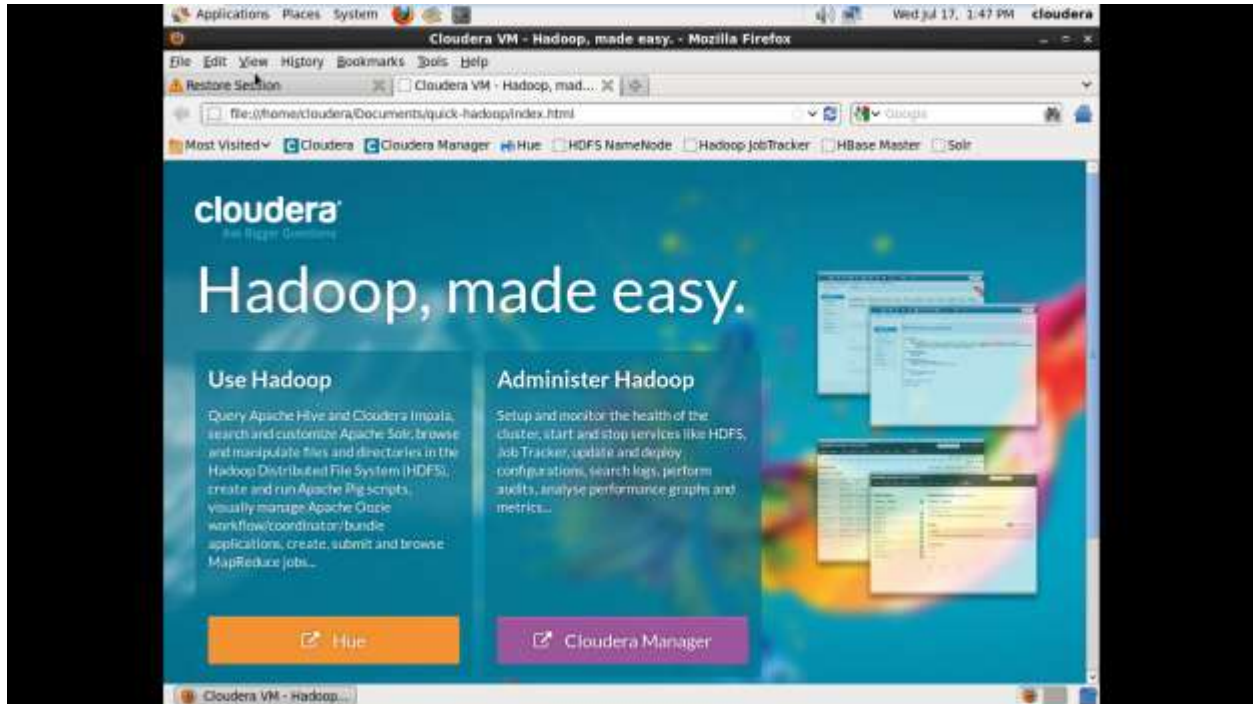11. As this is a new virtual machine, It will start deploying.

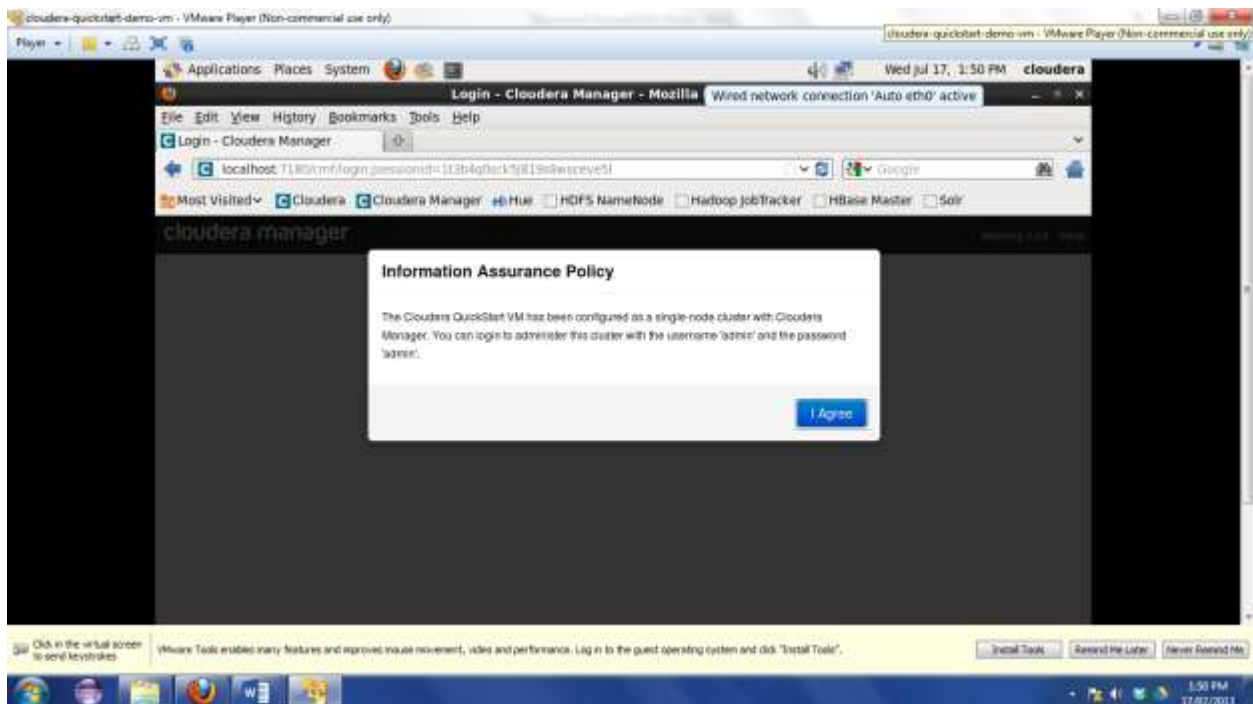12. After the installation/ deployment is complete, you get the virtual machine desktop.
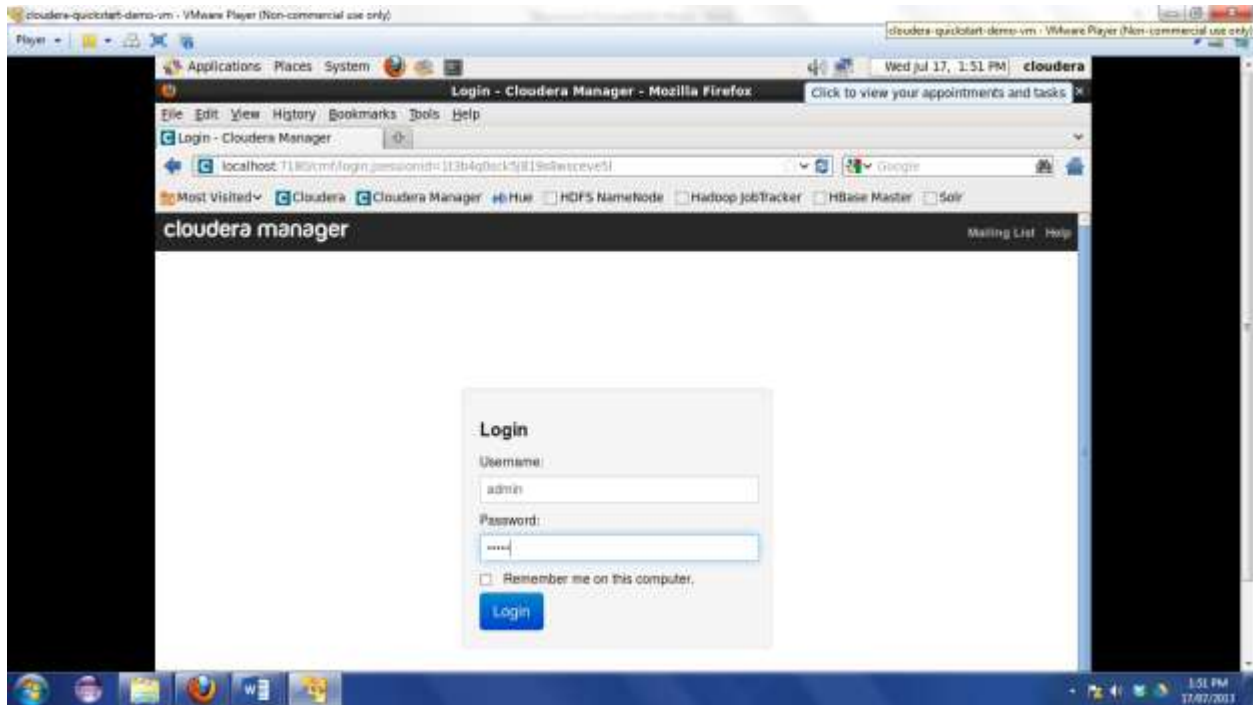
# Hadoop Configuration

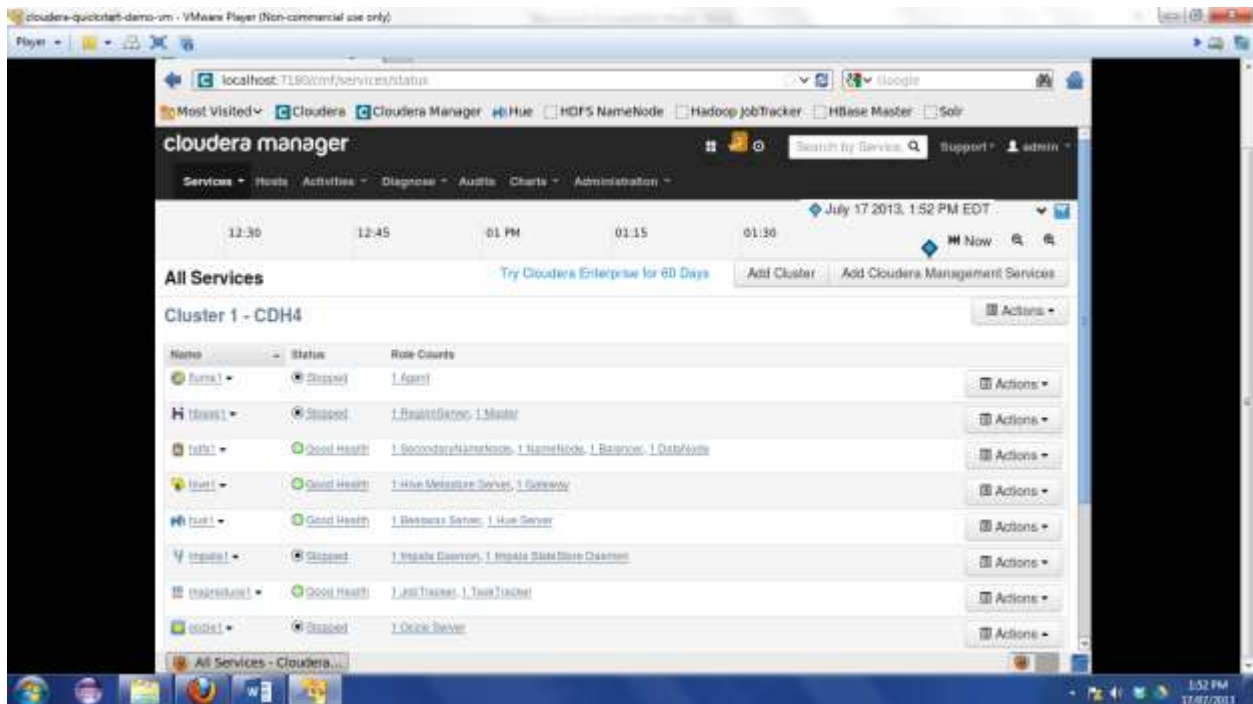1. Once the installation is complete, you will get below welcome page.



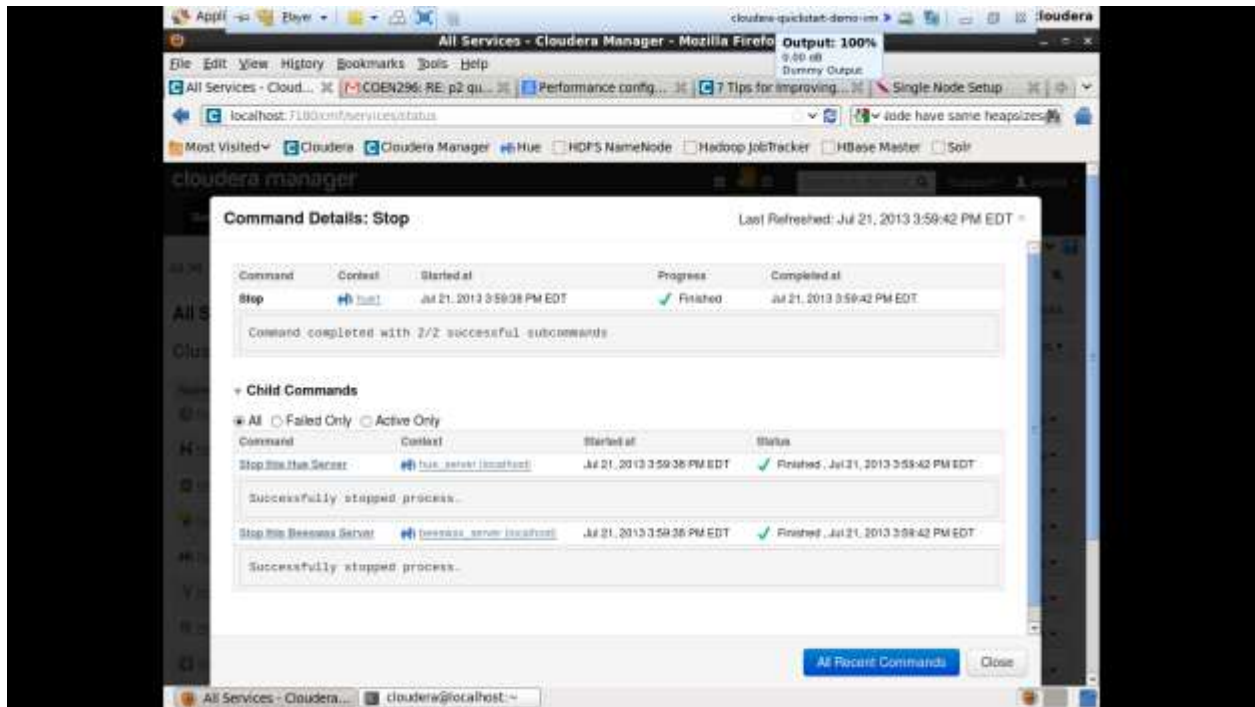2. Click on Cloudera Manager, as this is your first login, agree to Information Assurance policy.

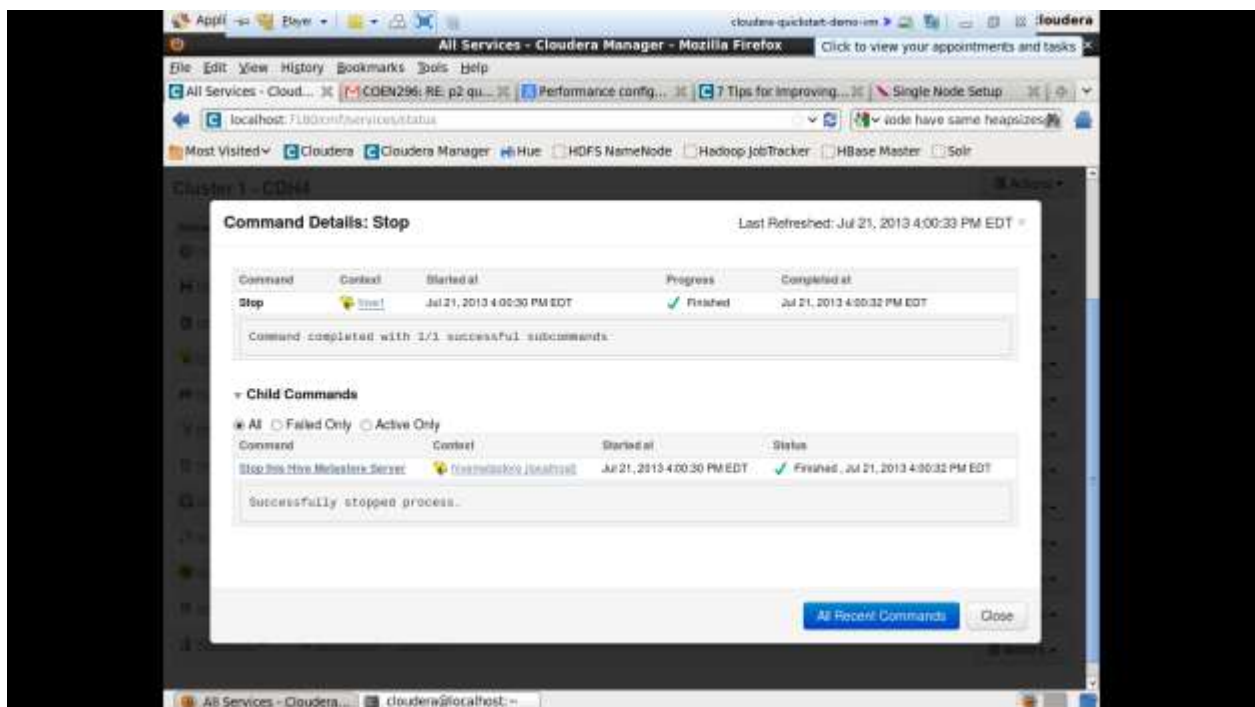3. Login to CDH using username: admin, password: admin.



4. Once you login, you will get list of all services available. Check and confirm there is no unhealthy service that needs to be fixed.

There are some services that are not required to be running to run the wordcount program. These services can be stopped for better performance.
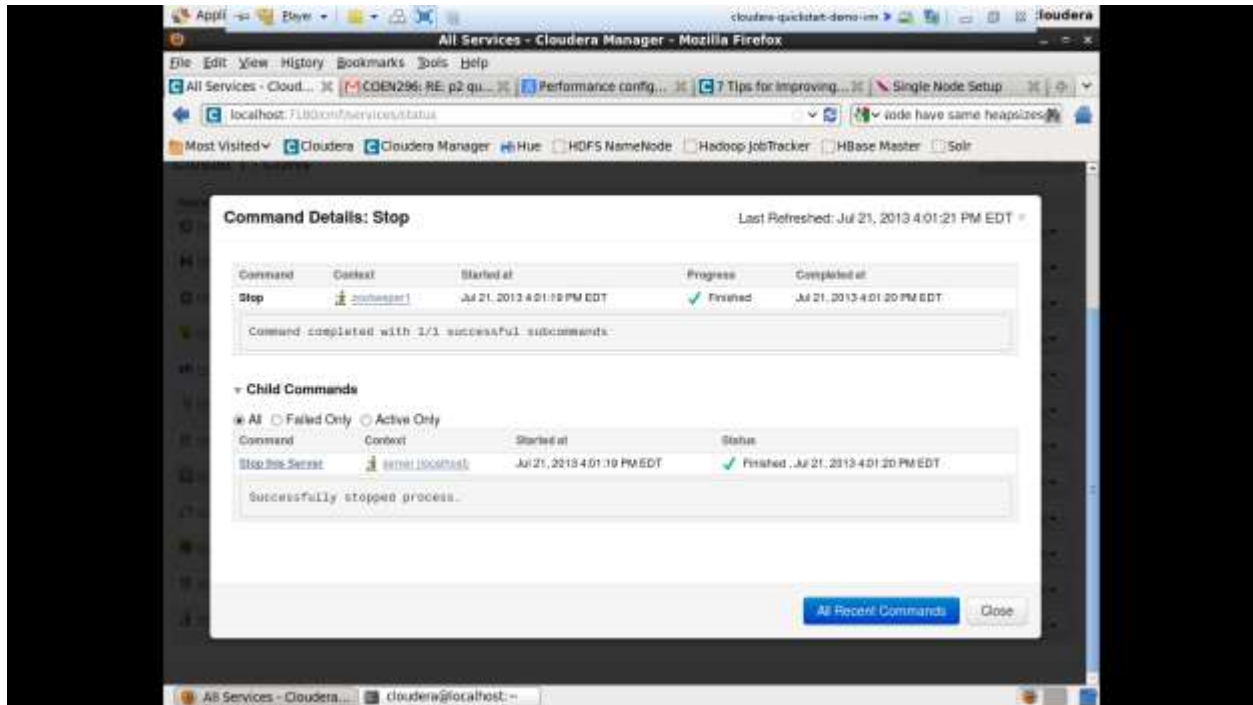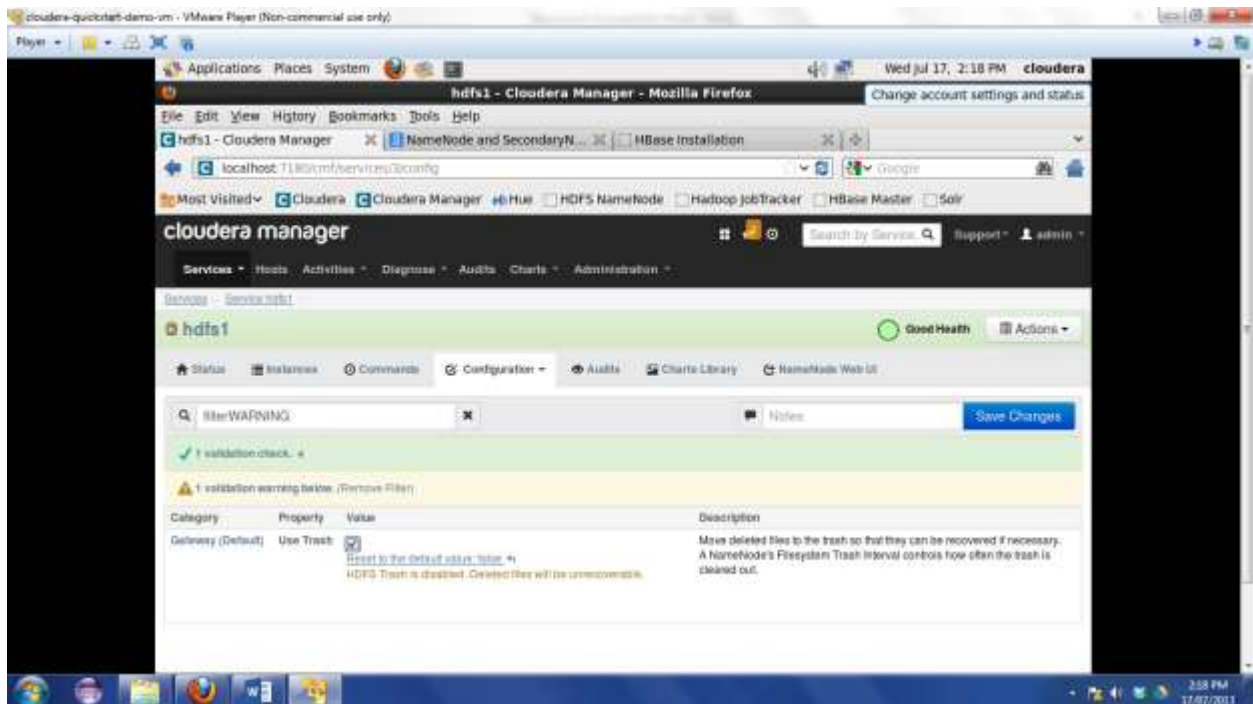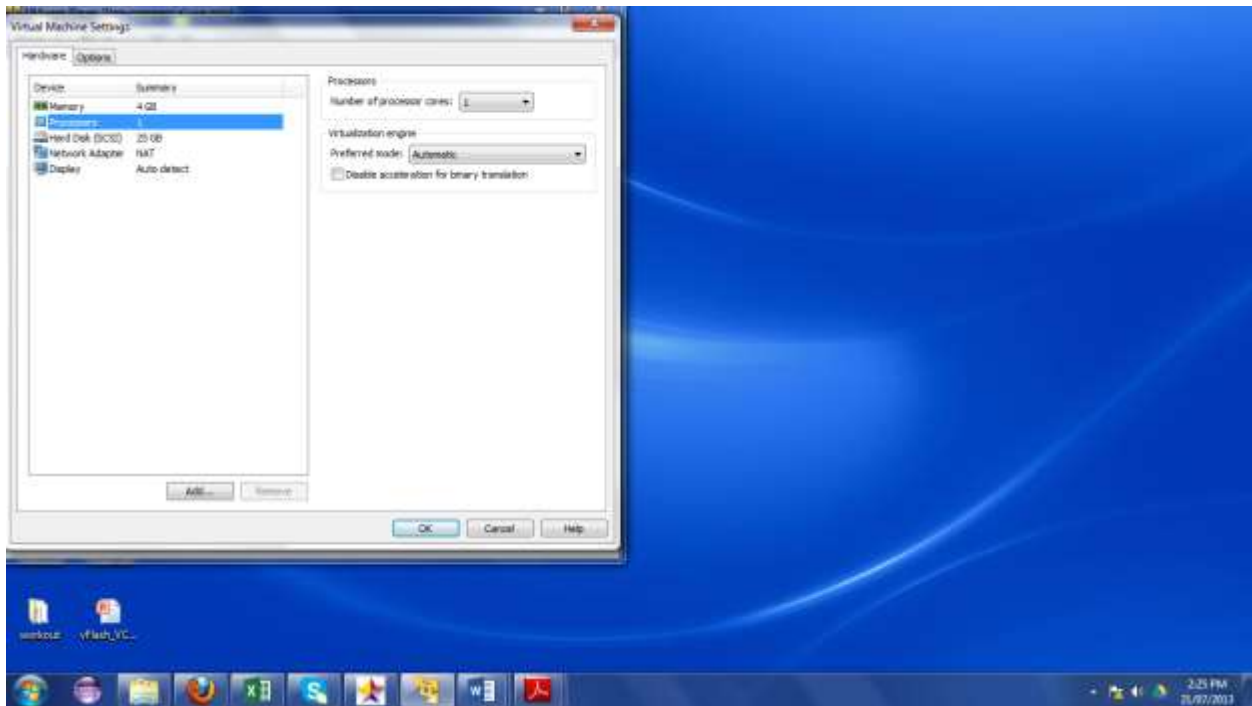
5. Stop Hue as it is not required to run word count program.

6. Stop Hive
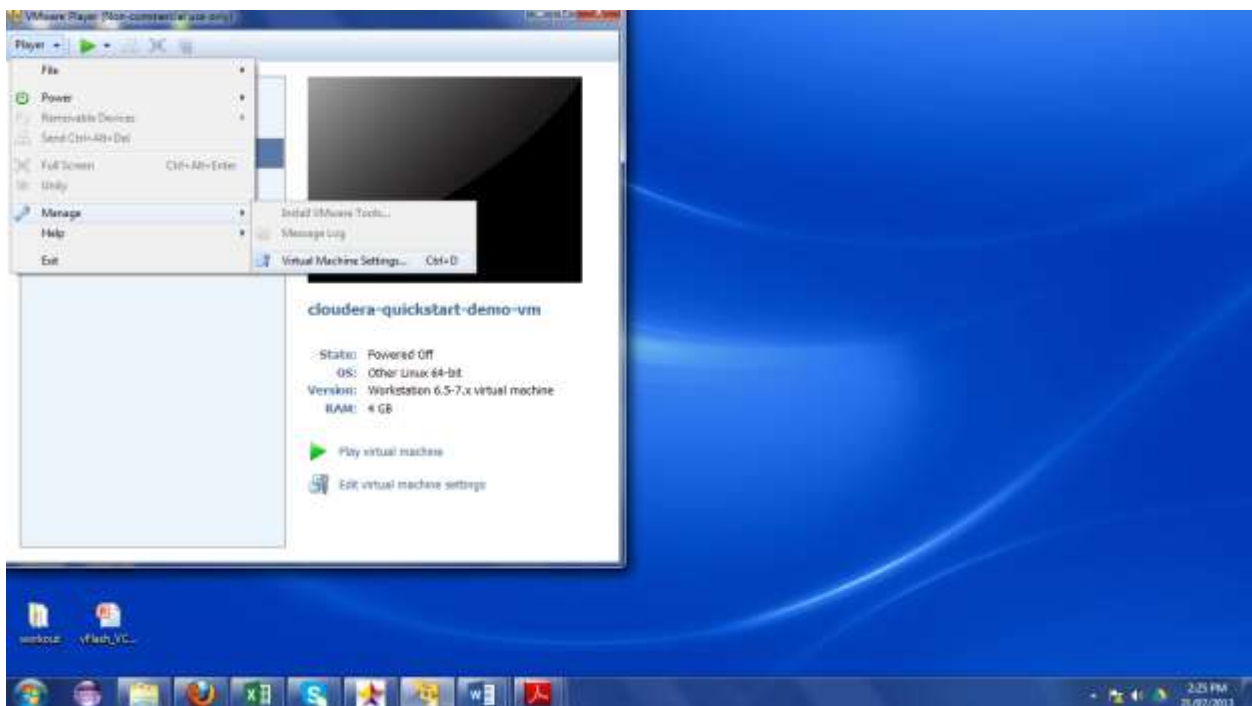
7. Stop service Zookeeper



8. Check Use trash property. This keeps a copy of every file that is deleted in _trash folder. Setting this property helps in better monitoring.
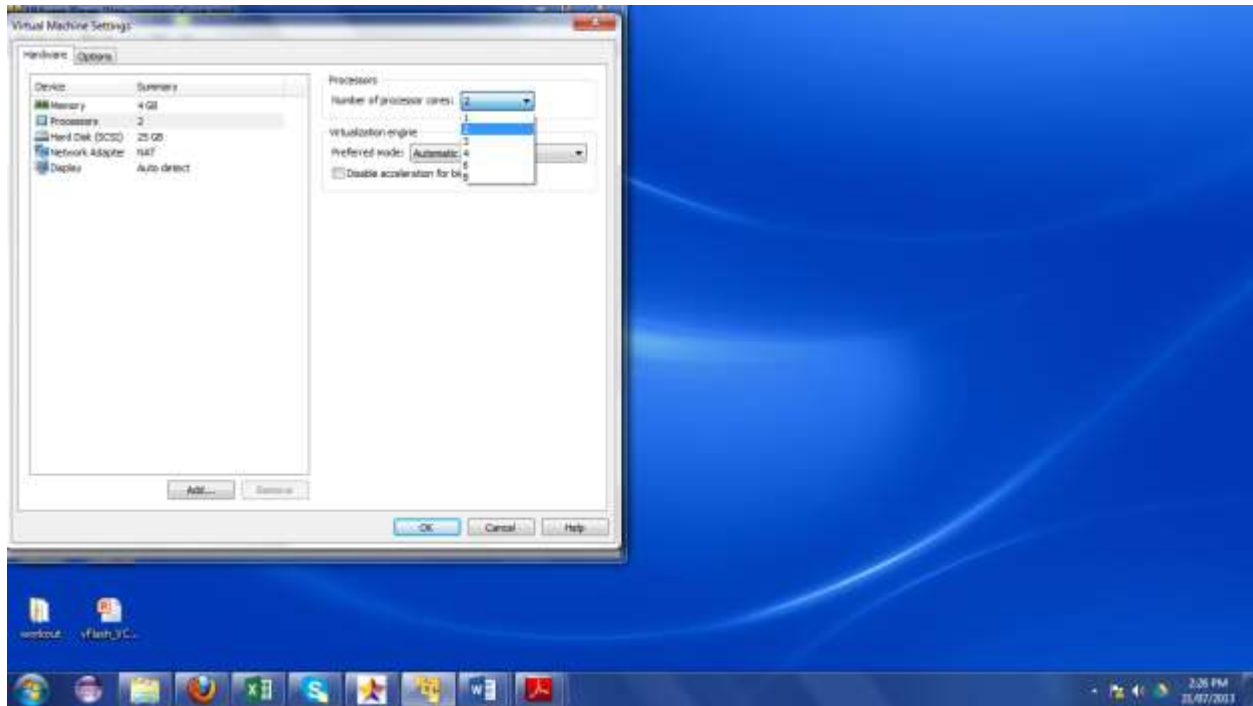
9. Also, By default each VM has a single processor, as shown in the below screenshot.



10. To enhance performance VM can be configured to be multi CPU. TO change the setting, Go to Player → Manage → Virtual Machine Settings…

11. Click on Processors and change the number to desired number. For our purposes, I have changed it to 2.

# Running a prgoram on Hadoop

1. Create a directory P2 and copy the sample program files to current working directory.

Command: mkdir P2

cd P2

ls –ltr

scp ssingh@linux.dc.engr/scu/edu:/home/mwang2/test/coen296/t2*.dat ./

ls -ltr



2. Create a directory wordcount_classes that will contain .jar file of the program.

Command: cd ..

mkdir wordcount_classes

ls -ltr



3. Copy one of the program files to WordCount.java file

Command: cp P2/t21.dat WordCount.java

ls -ltr

4. Set classpath variable to access Hadoop libraries.

Command:

Export classpath= "/usr/lib/hadoop/*:/usr/lib/hadoop/client-0.20/*"

Echo $classpath



5. Compile WordCount.java and save class files in wordcount_classes directory. Create a JAR file for all classes.

Command:

javac -cp $*classpath* -d wordcount_classes WordCount.java

jar -cvf wordcount.jar -C wordcount_classes/ .

ls -lrt

6. Create 2 sample files – file0 and file1.

Command:

echo "Hello World Bye World" > file0
echo "Hello Hadoop Goodbye Hadoop" > file1
ls –ltr file*
cat file0
cat file1



7. Create a wordcount directory and wordcount/input directory at Hadoop server.

Command:

hadoop fs -mkdir /user/cloudera/wordcount /user/cloudera/wordcount/input
hadoop fs -ls /user/cloudera/wordcount
hadoop fs -ls /user/cloudera/wordcount/input

8.  Copy the files file0 and file1 to the Hadoop input files directory
    /user/cloudera/wordcount/input.

Command:

hadoop fs -put file* /user/cloudera/wordcount/input

hadoop fs -ls /user/cloudera/wordcount/input

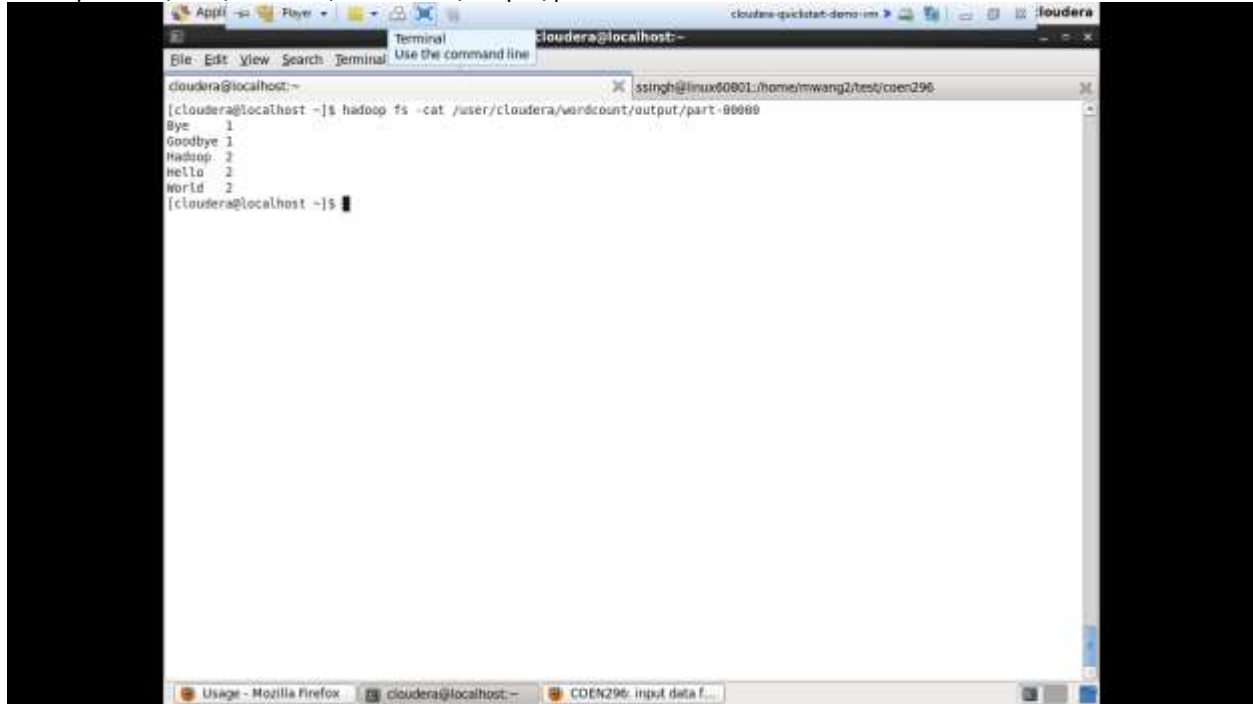9. Run the jar file wordcount.jar, reading input files from /user/cloudera/wordcount/input and providing output at /user/cloudera/wordcount/output.

Command:

hadoop jar wordcount.jar org.myorg.WordCount /user/cloudera/wordcount/input /user/cloudera/wordcount/output

10. See the output byh displaying contents of file part-00000 in output directory
    /user/cloudera/wordcount/output

Command:

hadoop fs -cat /user/cloudera/wordcount/output/part-00000