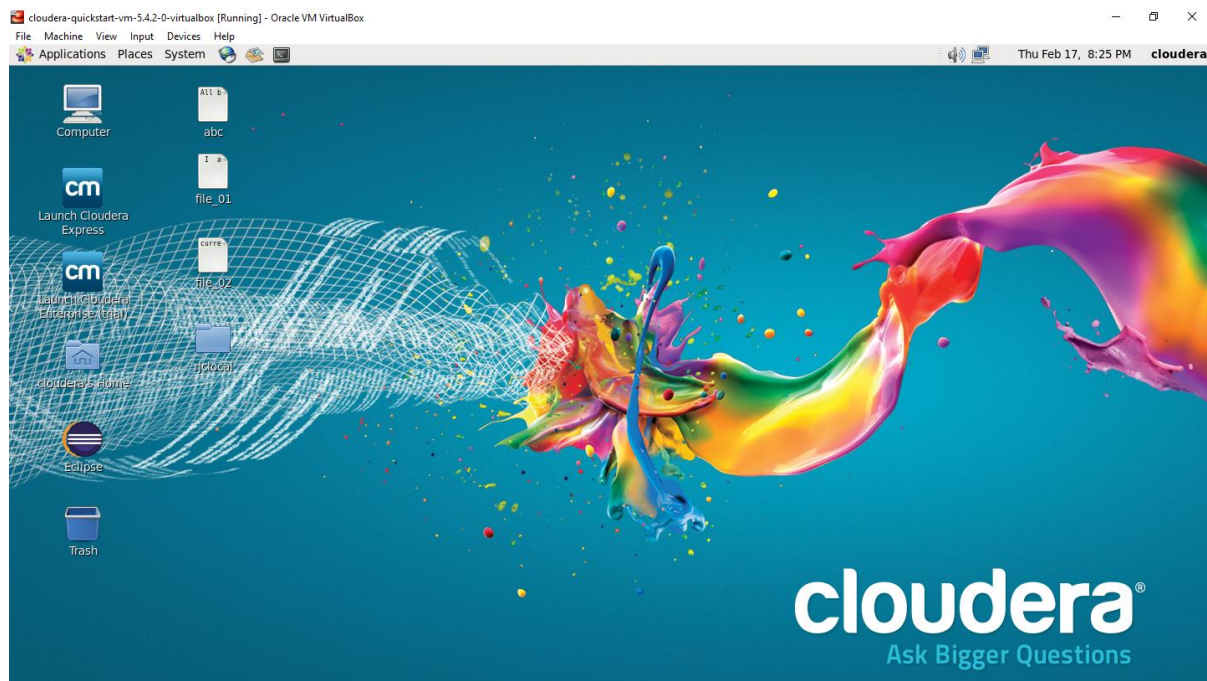# PRACTICAL NO 3

## To Implement Wordcount problem using Hadoop

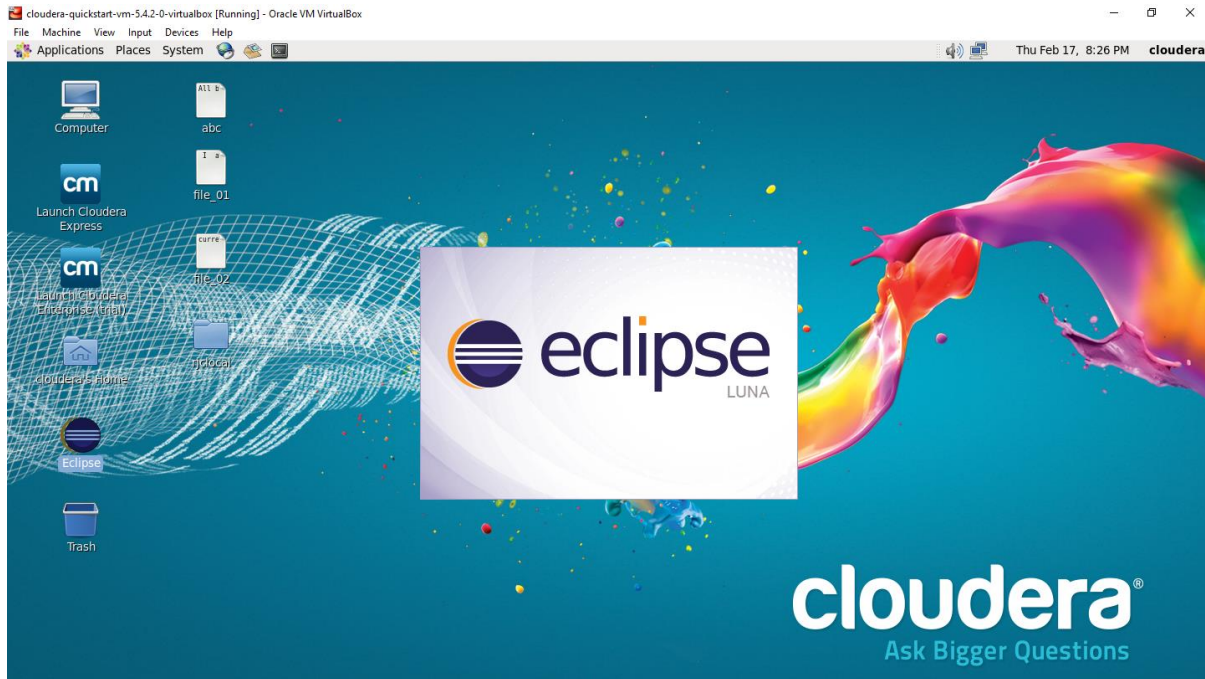## MapReduce in Eclipse: (With Combiner & Without Combiner)

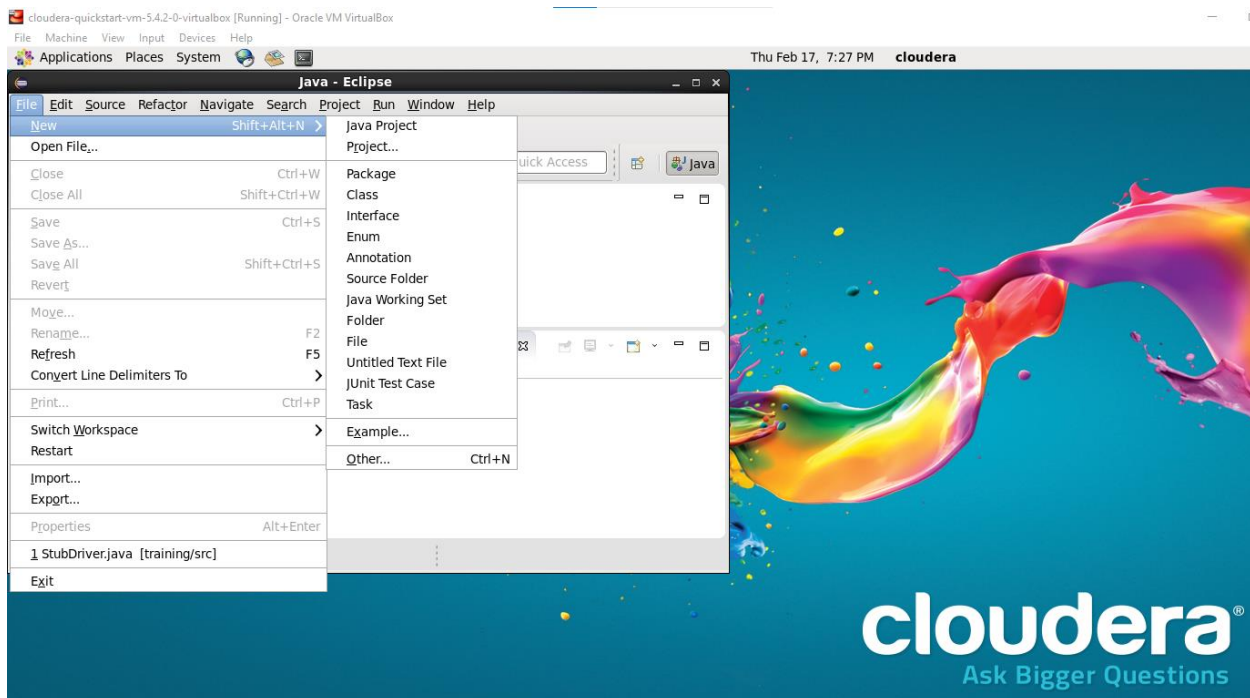- ## Steps for Word Count in Cloudera
  - ## ➢ With Combiner
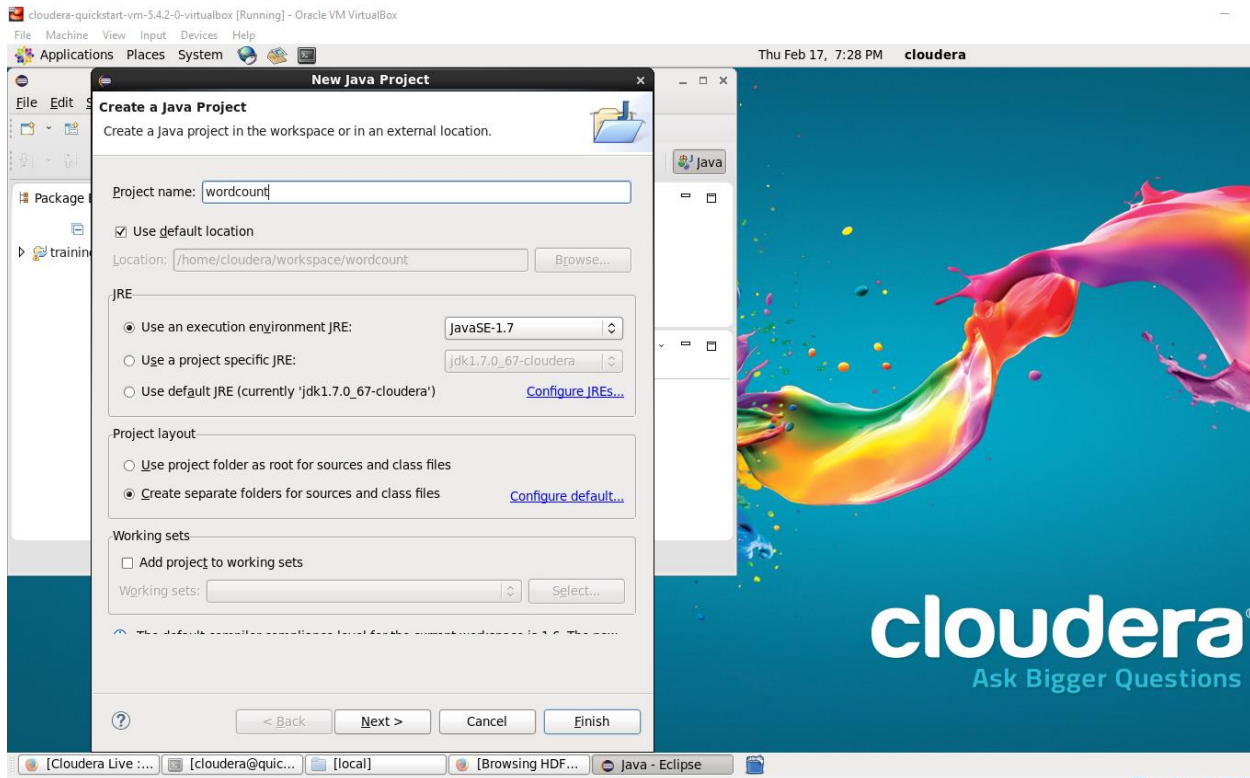
1) Open virtual box and then start cloudera quickstart



2) Open Eclipse present on the cloudera desktop

3) Create a new Java project clicking: File -> New -> Project -> Java Project -> Next ("WordCount" is the project name).
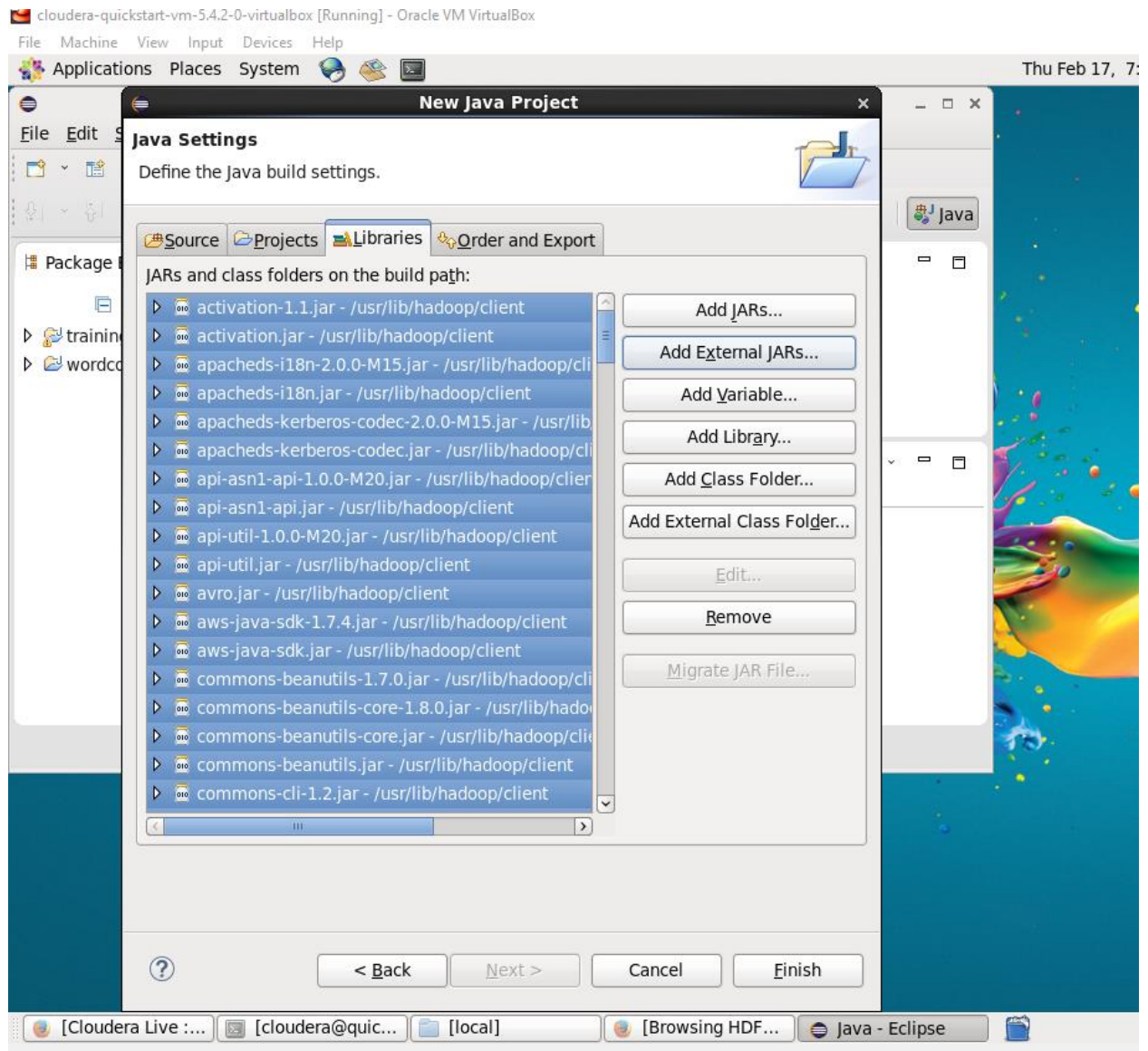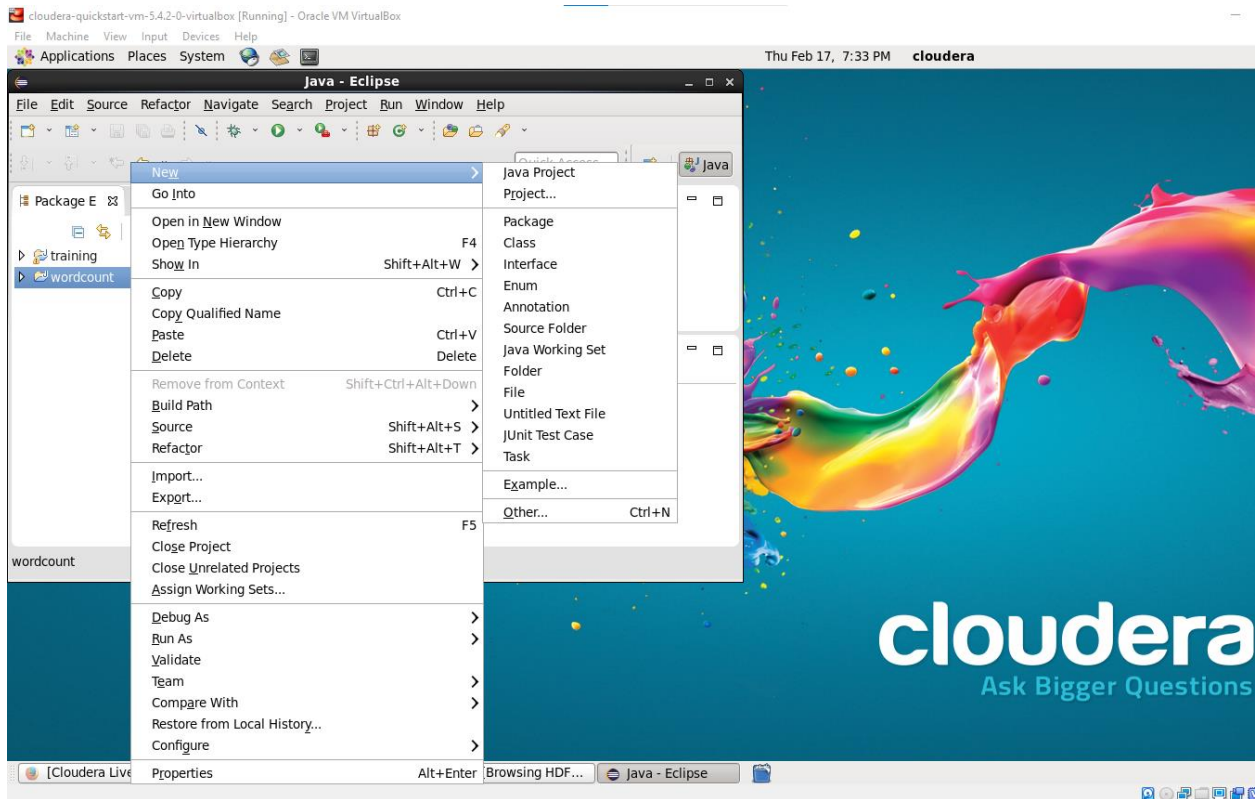
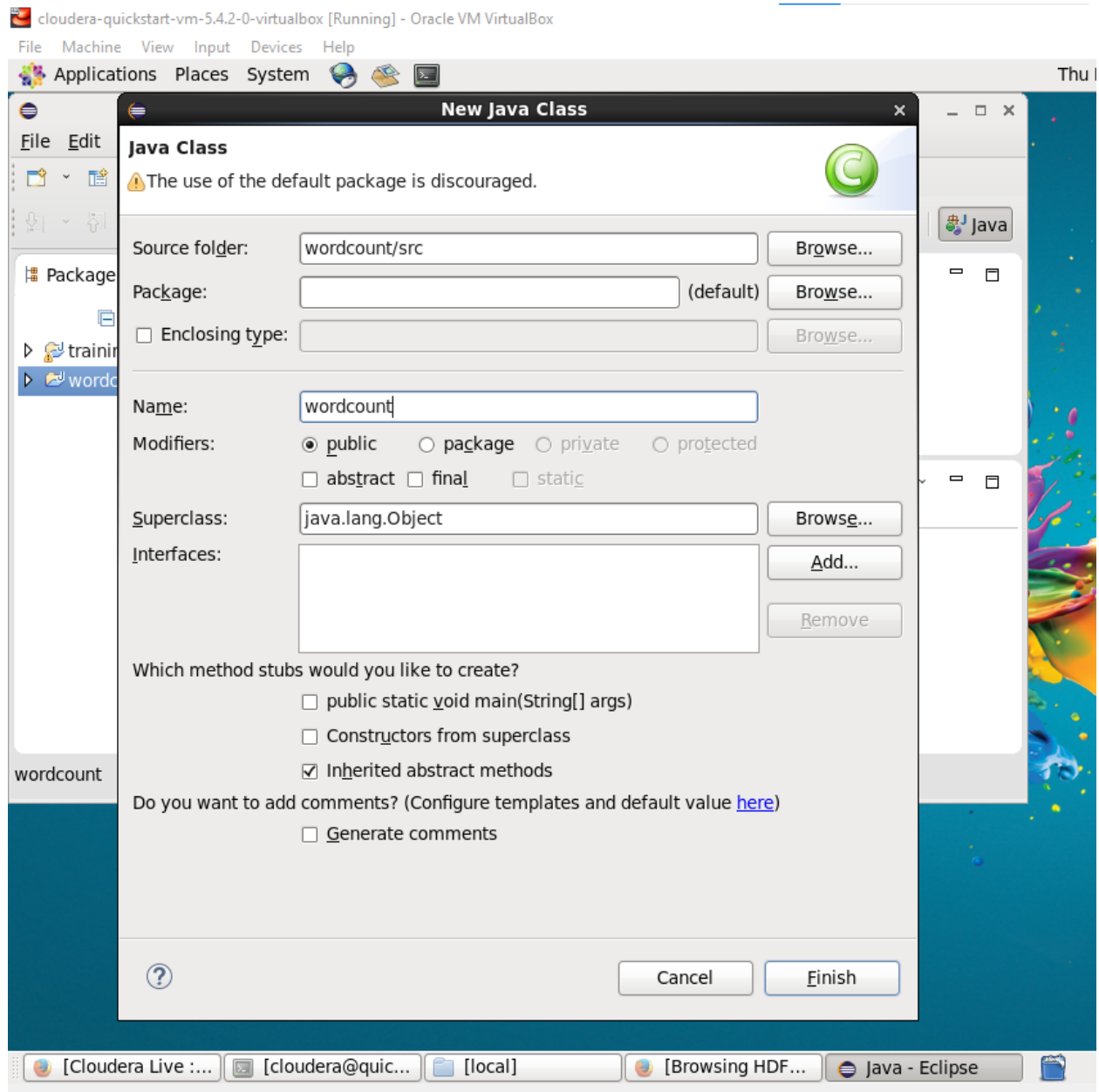4) Adding the Hadoop libraries to the project Click on Libraries -> Add External JARs Click on

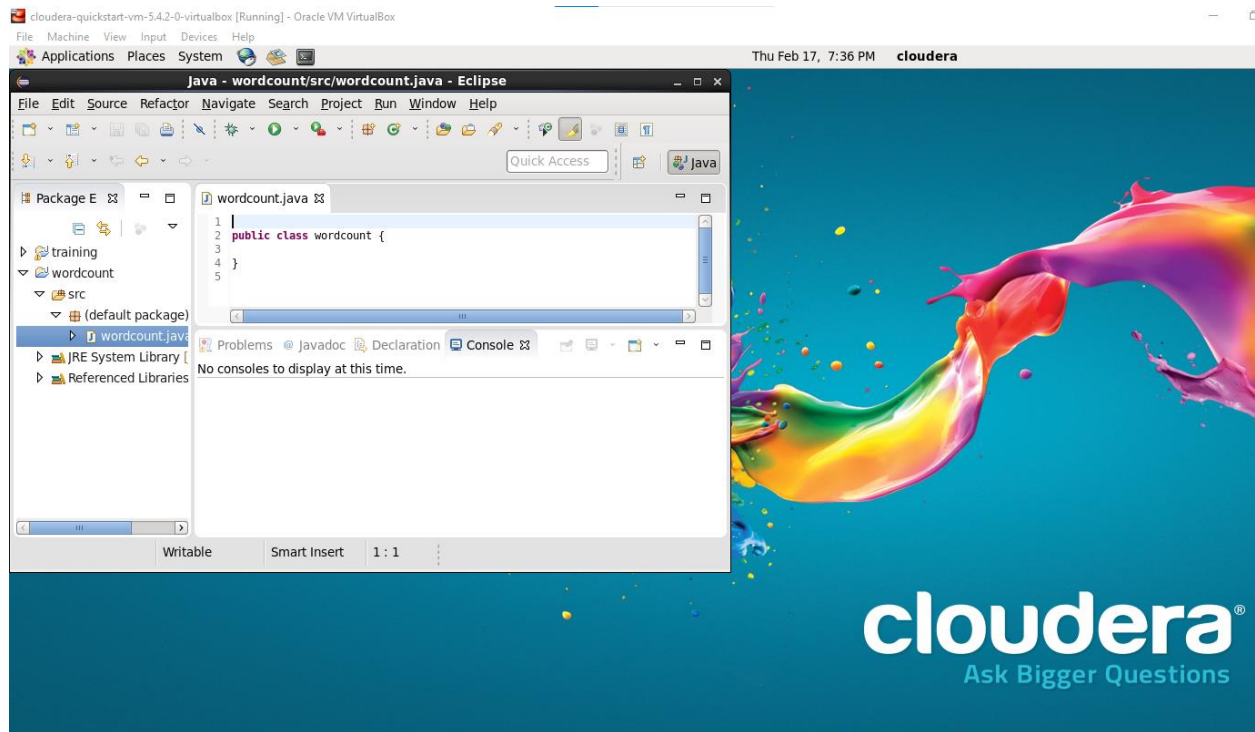File System -> usr -> lib -> hadoop Select all the libraries (JAR Files) -> click OK Click on

Add External jars, -> client -> select all jar files -> ok -> Finish

5) Right Click on the name of Project "WordCount" -> New -> class Don't write anything for package Write Name Textbox write "WordCount" -> Finish Then WordCount.java window will pop up
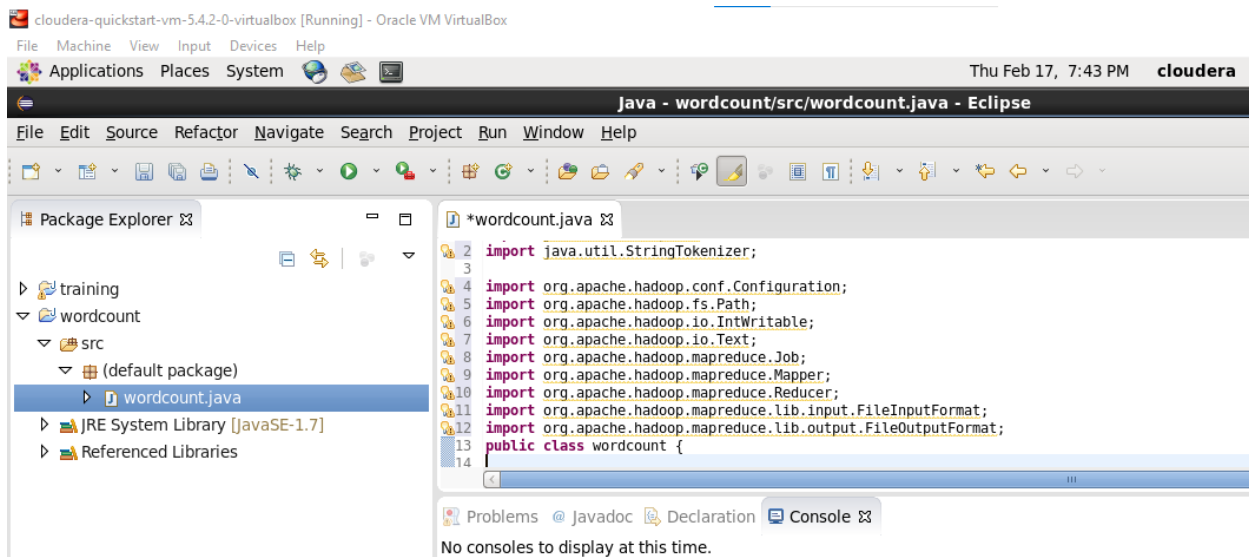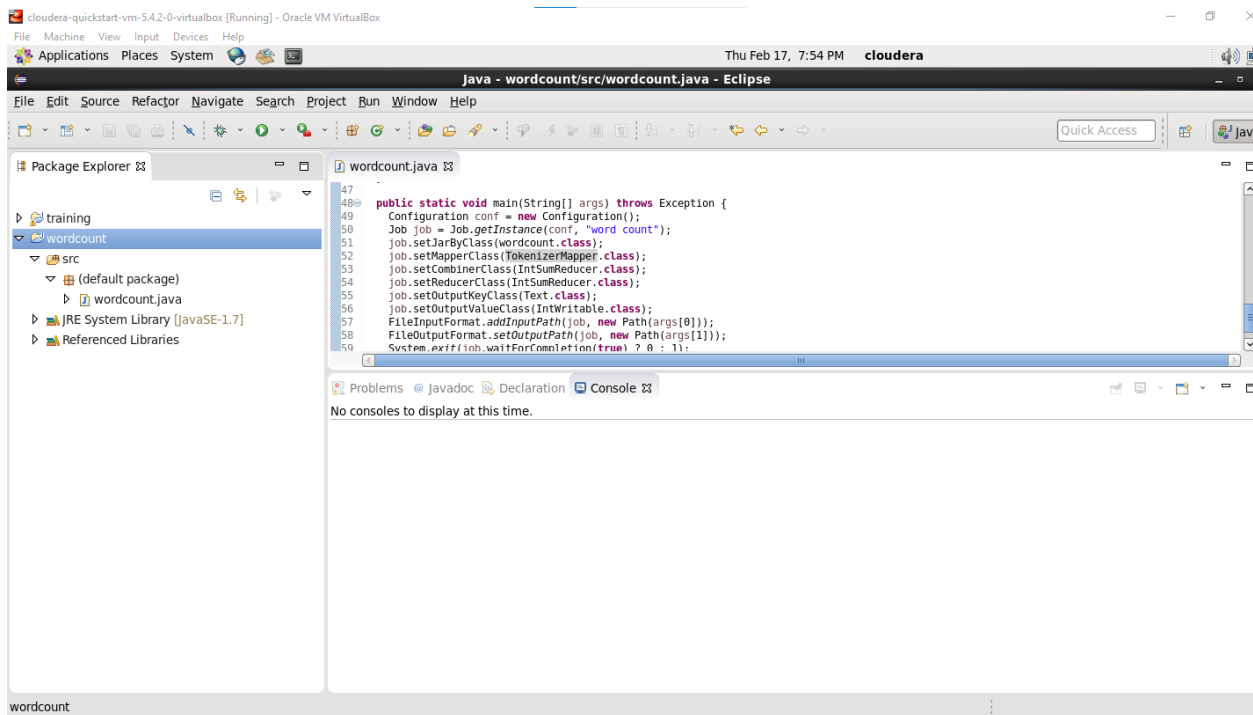
cloudera-quickstart-vm-5.4.2-0-virtualbox [Running] - Oracle VM VirtualBox

File   Machine   View   Input   Devices   Help

Applications   Places   System                                                    Thu

File   Edit                                                                                    Java

**New Java Class**                                                        ×

**Java Class**

⚠ The use of the default package is discouraged.

Source folder:   wordcount/src                          Browse...

Package:                                        (default)   Browse...

☐ Enclosing type:                                          Browse...

Name:   wordcount

Modifiers:   ● public   ○ package   ○ private   ○ protected

☐ abstract   ☐ final   ☐ static

Superclass:   java.lang.Object                          Browse...

Interfaces:                                               Add...

Remove

Which method stubs would you like to create?

☐ public static void main(String[] args)

☐ Constructors from superclass

☑ Inherited abstract methods

Do you want to add comments? (Configure templates and default value here)

☐ Generate comments

wordcount

?                                    Cancel        Finish

[Cloudera Live :...]   [cloudera@quic...]   [local]   [Browsing HDF...]   Java - Eclipse
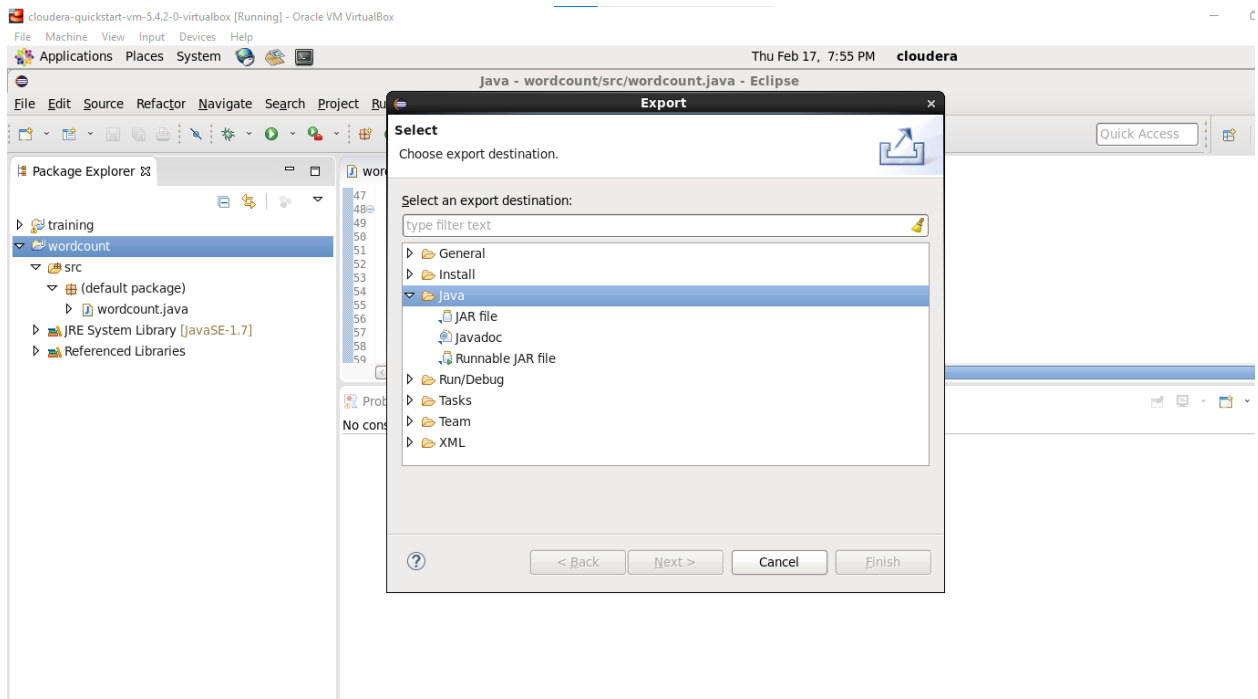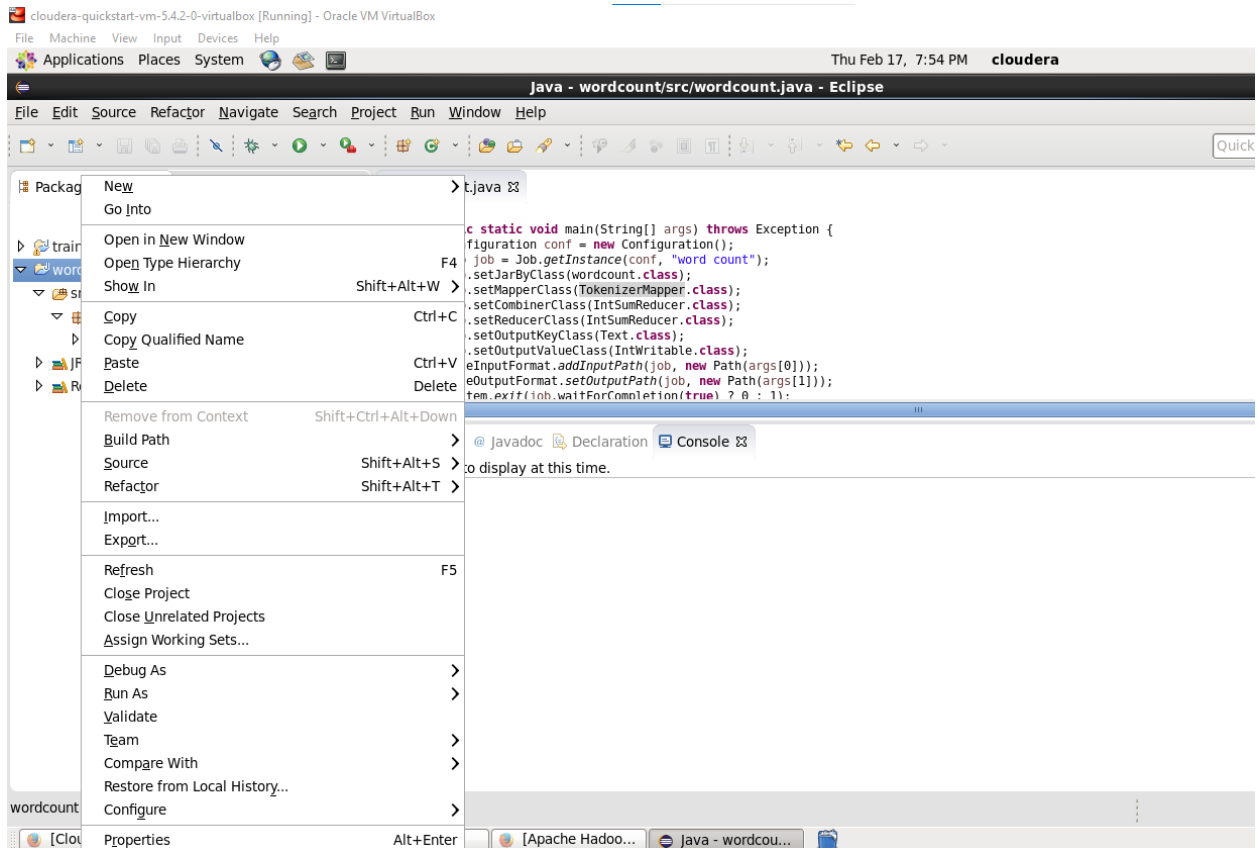
**Source code:**

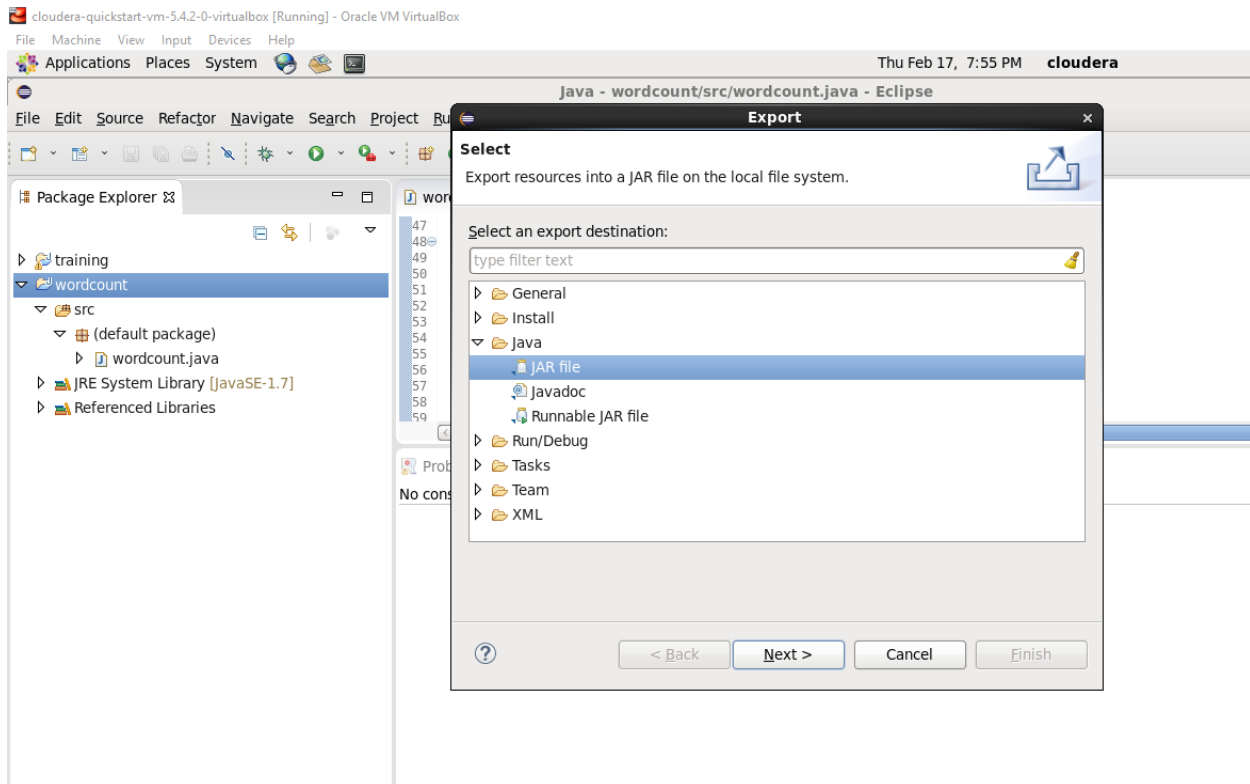**Packages**

## Mapper Logic



6) Right Click on the project name WordCount -> Export -> Java -> JAR File -> Next -> for

select the export destination for JAR file: browse -> Name : WordCount.jar -> save in folder

-> cloudera -> Finish -> OK

7) Verify jar file from terminal by using Open terminal & type " ls " There it will show WordCount.jar

Check current working directory

->pwd

->ls

8)   We need to create an input file in local file system

Creating an input file named as "abc".

9)   Now we have to move this input file to hdfs. For this we create a direcory on hdfs using

command hdfs dfs -mkdir /inputdir.

Then we can verify whether this directory is created or not using ls command hdfs dfs -ls /

Move the input file to this directory created in hdfs by using either put command or copyFromLocal command.

Now checking whether the "abc" present in /inputdir directory of hdfs or not using hdfs dfs -ls /inputdir command



As we can see "abc" file is present in /inputdir directory of hdfs. Now we will see the content of this file using hdfs dfs –cat /inputdir/abc command



10) Running Mapreduce Program on Hadoop, syntax is hadoop jar jarFileName.jar ClassName /InputFileAddress /outputdir

**i.e. hadoop jar /home/cloudera/WordCount.jar WordCount /inputdir/abc /outputdir**

```
[cloudera@quickstart ~]$ hadoop jar /home/cloudera/WordCount.jar WordCount /inpu
tdir/abc /outpurdir
22/02/17 20:09:44 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0
:8032
22/02/17 20:09:45 WARN mapreduce.JobSubmitter: Hadoop command-line option parsin
g not performed. Implement the Tool interface and execute your application with
ToolRunner to remedy this.
22/02/17 20:09:45 INFO input.FileInputFormat: Total input paths to process : 1
22/02/17 20:09:45 INFO mapreduce.JobSubmitter: number of splits:1
22/02/17 20:09:46 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_16
44894610889_0001
22/02/17 20:09:47 INFO impl.YarnClientImpl: Submitted application application_16
44894610889_0001
22/02/17 20:09:47 INFO mapreduce.Job: The url to track the job: http://quickstar
t.cloudera:8088/proxy/application_1644894610889_0001/
22/02/17 20:09:47 INFO mapreduce.Job: Running job: job_1644894610889_0001
22/02/17 20:10:02 INFO mapreduce.Job: Job job_1644894610889_0001 running in uber
 mode : false
22/02/17 20:10:02 INFO mapreduce.Job:  map 0% reduce 0%
22/02/17 20:10:17 INFO mapreduce.Job:  map 100% reduce 0%
22/02/17 20:10:28 INFO mapreduce.Job:  map 100% reduce 100%
22/02/17 20:10:28 INFO mapreduce.Job: Job job_1644894610889_0001 completed succe
ssfully
22/02/17 20:10:28 INFO mapreduce.Job: Counters: 49
        File System Counters
                FILE: Number of bytes read=1156
                FILE: Number of bytes written=222899
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
```

| [part-r-00000 (~/Down... | [Java - WordCount/src/... | [Cloudera Live : Welco... | cloudera@quickstart:~ | [Microsoft Word - Dow... | Browsing HDFS - Mozil... |

## Map-Reduce Framework

```
cloudera-quickstart-vm-5.4.2-0-virtualbox [Running] - Oracle VM VirtualBox
File  Machine  View  Input  Devices  Help
Applications  Places  System                                     Thu Feb 17, 8:42 PM   cloudera
                              cloudera@quickstart:~
File  Edit  View  Search  Terminal  Help
22/02/17 20:10:28 INFO mapreduce.Job: Job job_1644894610889_0001 completed succe
ssfully
22/02/17 20:10:28 INFO mapreduce.Job: Counters: 49
        File System Counters
                FILE: Number of bytes read=1156
                FILE: Number of bytes written=222899
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=914
                HDFS: Number of bytes written=799
                HDFS: Number of read operations=6
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
        Job Counters
                Launched map tasks=1
                Launched reduce tasks=1
                Data-local map tasks=1
                Total time spent by all maps in occupied slots (ms)=11842
                Total time spent by all reduces in occupied slots (ms)=8821
                Total time spent by all map tasks (ms)=11842
                Total time spent by all reduce tasks (ms)=8821
                Total vcore-seconds taken by all map tasks=11842
                Total vcore-seconds taken by all reduce tasks=8821
                Total megabyte-seconds taken by all map tasks=12126208
                Total megabyte-seconds taken by all reduce tasks=9032704
        Map-Reduce Framework
                Map input records=1
                Map output records=132
                Map output bytes=1333
                Map output materialized bytes=1156
                Input split bytes=109
                Combine input records=132
                Combine output records=88
                Reduce input groups=88
                Reduce shuffle bytes=1156
                Reduce input records=88
                Reduce output records=88
                Spilled Records=176
                Shuffled Maps =1
```

| [part-r-00000 (~/Down... | [Java - WordCount/src/... | [Cloudera Live : Welco... | cloudera@quickstart:~ | [Microsoft Word - Dow... | Browsing HDFS - Mozil... |

**As we can see in the above output,**

**Combine input records=132**

**Combine output records=88**

**And Reduce shuffle bytes coming as,**

**Reduce shuffle bytes=1876**

11) Then we can verify the content of outputdir directory and in that part-r file has the actual

output by using the command Hdfs dfs -cat /outputdir/part-r-00000 This will give us final output.
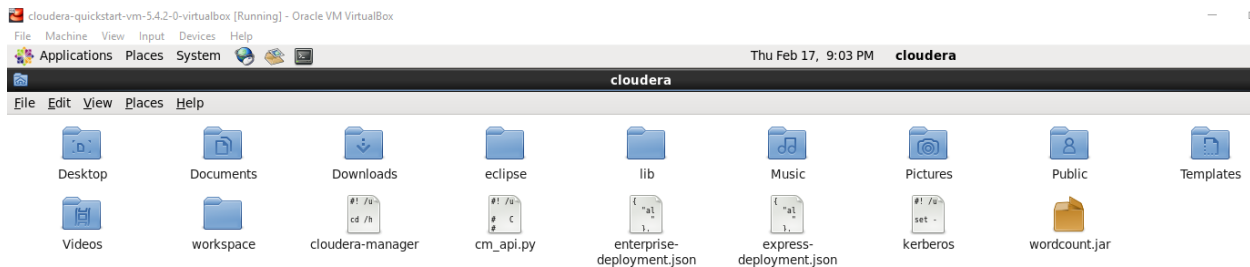
The same file can also be accessed using a browser. For every execution of this program we need

to delete the output directory or give a new name to the output directory every time.

1st we are checking whether the outputdir directory is created in hdfs or not using command

**hdfs dfs -ls /**



Now let's check what we have inside this **outputdir** directory using command as **hdfs dfs -ls**

**/outputdir**



Now we want to read the content of the **part-r-00000 file** which present inside the **outputdir**

using command **hdfs dfs -cat /outputdir/part-r-00000**

**It will give the count of number of times each word has occurred as output.**

**12) The same file can also be accessed using a browser.**

Browse the Directory by

**Hadoop->HDFS Namenode->Ultilities ->Browse the file system**



Now downloading the **part-r-00000** file.



Inside the **part-r-00000** file it will have the same output as we are getting after executing using

command **hadoop jar /home/cloudera/WordCount.jar WordCount /inputdir/abc /op1**

**For every execution of this program we need to delete the output directory or give a new**

**name to the output directory every time.**

## ➢ **Implementation of WordCount problem using Hadoop MapReduce (Without Combiner) in Eclipse:**

1) We will perform the same steps as we have done above for WordCount (without using

combiner) in that we just commenting the combiner line in main function**.**



**2)** And will delete the WordCount.jar file in which all jar files are present from **/home/cloudera.**

**We have successfully deleted the WordCount.jar file.**

3)  Now exporting the jar files Right Click on the project name WordCount -> Export -> Java ->JAR
    File -> Next -> for select the export destination for JAR file: browse -> Name : WordCount.jar ->
    save in folder -> cloudera -> Finish -> OK

4)  Now checking the WordCount.jar file is created or not using –ls command



5)  Running Mapreduce Program on Hadoop, syntax is hadoop jar jarFileName.jar ClassName

/InputFileAddress /outputdir

**i.e. hadoop jar /home/cloudera/WordCount.jar WordCount /inputdir/abc /op1**

here I am using the same input file 'abc' which I have created earlier for WordCount

example (Without Combiner**). For every execution of this program we need to delete the**

**output directory or give a new name to the output directory every time.** So here I am

giving the new name to the output directory as **'op1'.**

- As we can see from above image the the combiner input and output records coming out as,

   **Combine input records=0**

   **Combine output records=0**

- Earlier it was coming out as "zero" while executing WordCount (without combiner).

   **Combine input records=132**

**Combine output records=88**

- And also here we are getting the Reduce Shuffle bytes as,

  **Reduce shuffle bytes=942**

   Earlier while executing WordCount (without combiner) it is coming out as,
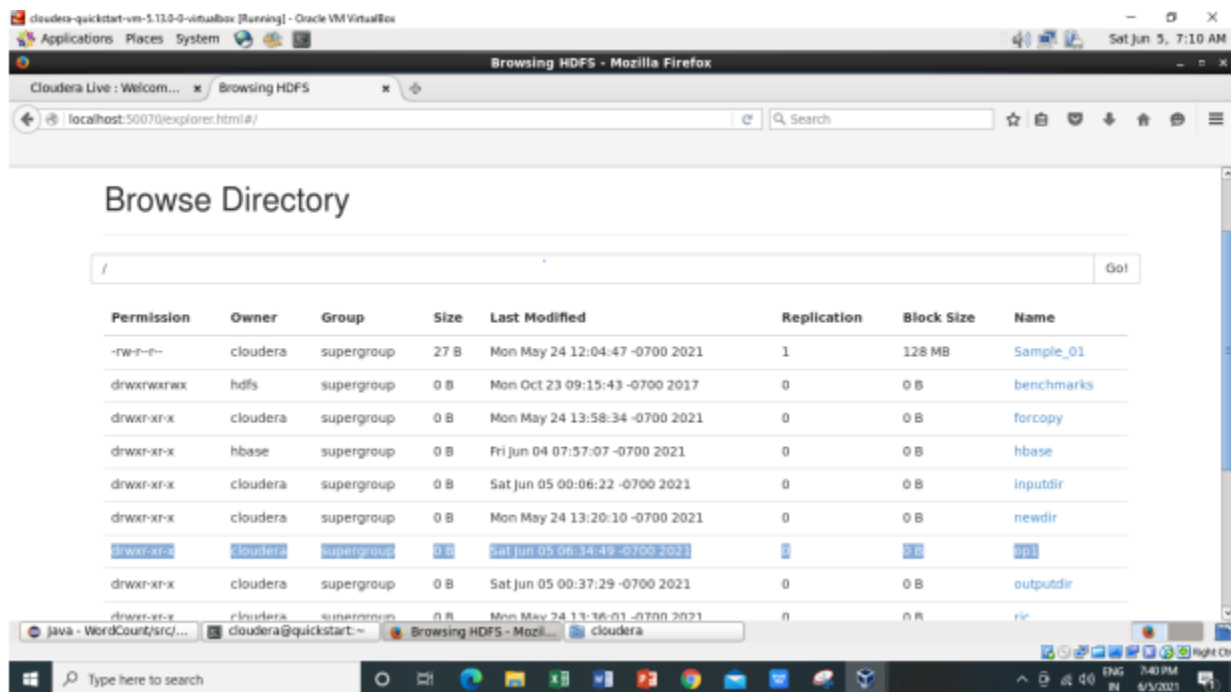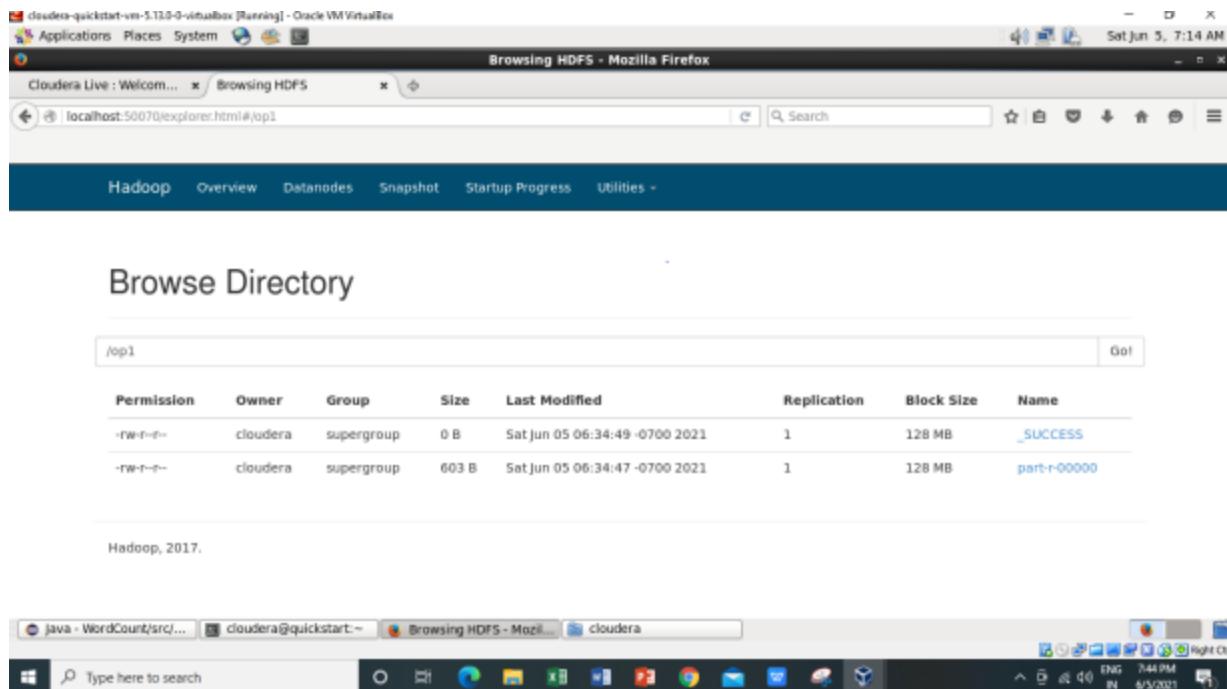
  **Reduce shuffle bytes=1876**

- So Combiner is used to save the Network Bandwidth. So for saving the Network bandwidth we make use of combiner. So instead of sending every word over the network what we do is we incorporate the logic of the reducer at the combiner side so that the less amount of information can be transmitted over the network.
- So when we are not using combiner 1876 bytes acting as an input for the reducer. And when we are making use of the combiner so 942 bytes acting as input for the reducer.

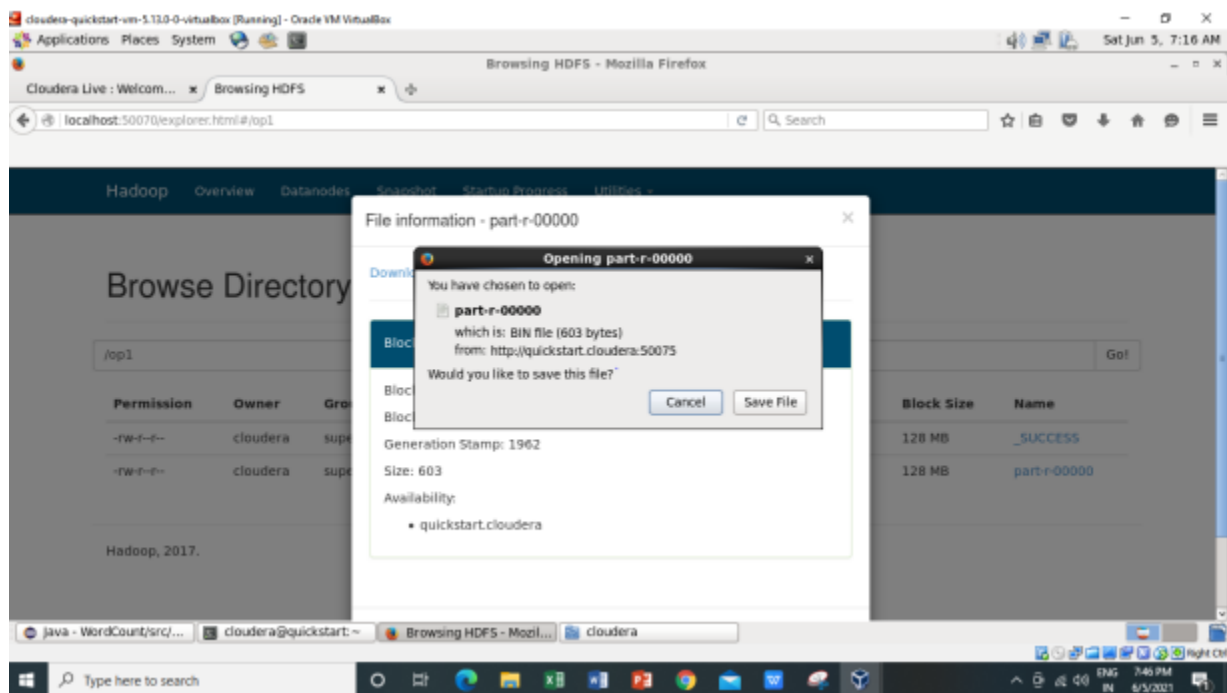6) The same file can also be accessed using a browser.

Browse the Directory by

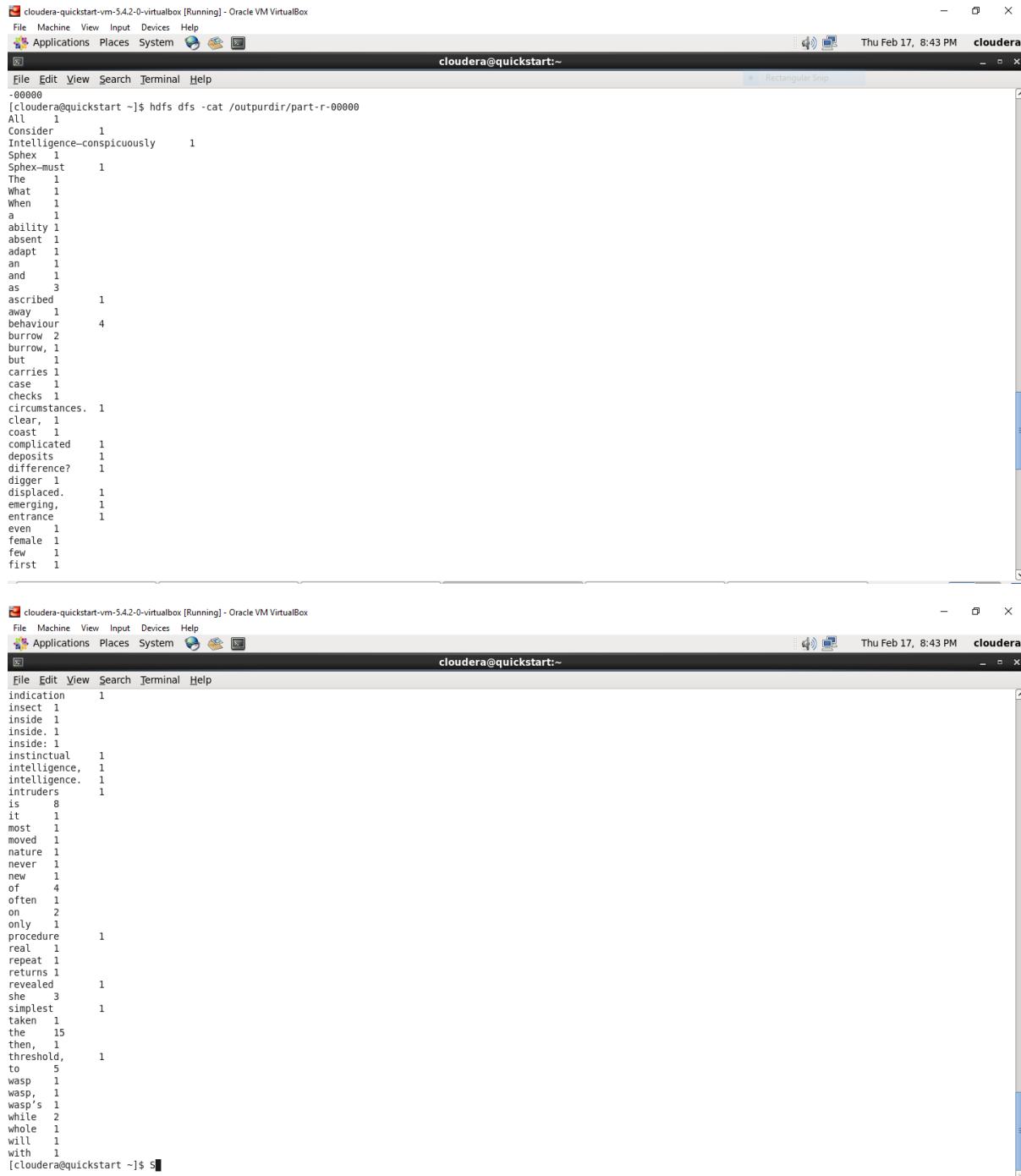**Hadoop->HDFS Namenode->Ultilities ->Browse the file system**

Now downloading the **part-r-00000** file.



Inside the **part-r-00000** file it will have the same output as we are getting after executing using

command **hadoop jar /home/cloudera/WordCount.jar WordCount /inputdir/abc /op1**