

Assignment – 3

Experiment 3 - To Implement WordCount problem using Hadoop MapReduce in Eclipse: (Without Combiner & With Combiner)

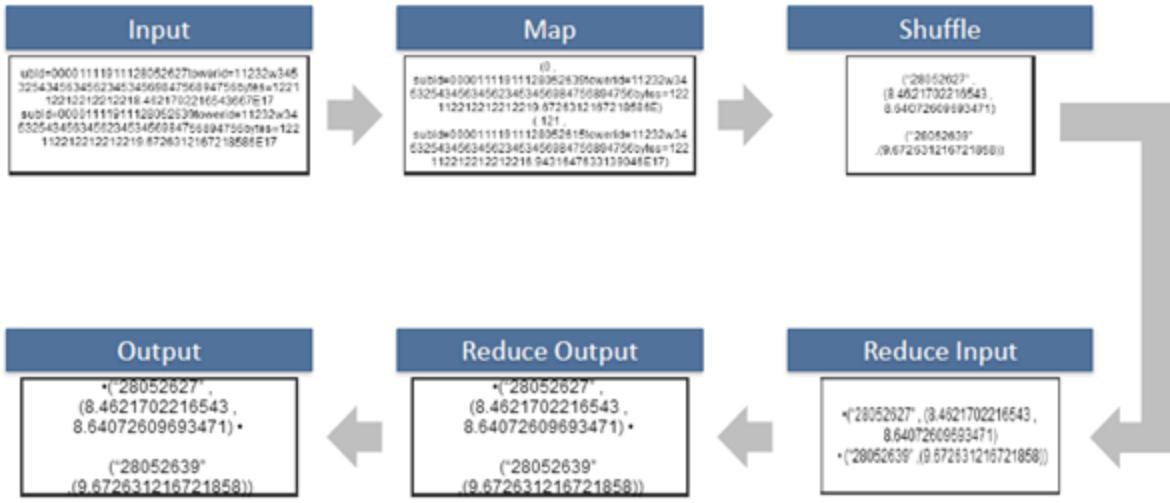
What is MapReduce ?

A MapReduce is a data processing tool which is used to process the data parallelly in a distributed form. It was developed in 2004, on the basis of paper titled as "MapReduce: Simplified Data Processing on Large Clusters," published by Google.

The MapReduce is a paradigm which has two phases, the mapper phase, and the reducer phase. In the Mapper, the input is given in the form of a key-value pair. The output of the Mapper is fed to the reducer as input. The reducer runs only after the Mapper is over. The reducer too takes input in key-value format, and the output of reducer is the final output.

Steps in Map Reduce

- The map takes data in the form of pairs and returns a list of <key, value> pairs. The keys will not be unique in this case.
- Using the output of Map, sort and shuffle are applied by the Hadoop architecture. This sort and shuffle acts on these list of <key, value> pairs and sends out unique keys and a list of values associated with this unique key <key, list(values)>.
- An output of sort and shuffle sent to the reducer phase. The reducer performs a defined function on a list of values for unique keys, and Final output <key, value> will be stored/displayed.

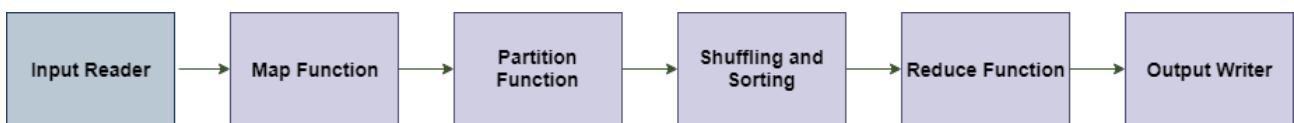


Sort and Shuffle

The sort and shuffle occur on the output of Mapper and before the reducer. When the Mapper task is complete, the results are sorted by key, partitioned if there are multiple reducers, and then written to disk. Using the input from each Mapper $\langle k2, v2 \rangle$, we collect all the values for each unique key $k2$. This output from the shuffle phase in the form of $\langle k2, \text{list}(v2) \rangle$ is sent as input to reducer phase.

Data Flow In MapReduce

MapReduce is used to compute the huge amount of data . To handle the upcoming data in a parallel and distributed form, the data has to flow from various phases.



Phases of MapReduce data flow

Input reader

The input reader reads the upcoming data and splits it into the data blocks of the appropriate size (64 MB to 128 MB). Each data block is associated with a Map function.

Once input reads the data, it generates the corresponding key-value pairs. The input files reside in HDFS.

Map function

The map function process the upcoming key-value pairs and generated the corresponding output key-value pairs. The map input and output type may be different from each other.

Partition function

The partition function assigns the output of each Map function to the appropriate reducer. The available key and value provide this function. It returns the index of reducers.

Shuffling and Sorting

The data are shuffled between/within nodes so that it moves out from the map and get ready to process for reduce function. Sometimes, the shuffling of data can take much computation time.

The sorting operation is performed on input data for Reduce function. Here, the data is compared using comparison function and arranged in a sorted form.

Reduce function

The Reduce function is assigned to each unique key. These keys are already arranged in sorted order. The values associated with the keys can iterate the Reduce and generates the corresponding output.

Output writer

Once the data flow from all the above phases, Output writer executes. The role of Output writer is to write the Reduce output to the stable storage.

To Implement WordCount problem using Hadoop MapReduce in Eclipse:

Hadoop WordCount operation occurs in 3 stages –

- Mapper Phase
- Shuffle Phase
- Reducer Phase

Hadoop WordCount - Mapper Phase Execution

- The text from the input text file is tokenized into words to form a key value pair with all the words present in the input text file. The key is the word from the input file and value is ‘1’.

- For instance if you consider the sentence “An elephant is an animal”.
- The mapper phase in the WordCount example will split the string into individual tokens i.e. words. In this case, the entire sentence will be split into 5 tokens (one for each word) with a value 1 as shown below –

Key-Value pairs from Hadoop Map Phase Execution-

(an,1)

(elephant,1)

(is,1)

(an,1)

(animal,1)

Hadoop WordCount Example- Shuffle Phase Execution

- After the map phase execution is completed successfully, shuffle phase is executed automatically wherein the key-value pairs generated in the map phase are taken as input and then sorted in alphabetical order.
- After the shuffle phase is executed from the WordCount example code, the output will look like this -

(an,1)

(an,1)

(animal,1)

(elephant,1)

(is,1)

Hadoop WordCount Example- Reducer Phase Execution

- In the reduce phase, all the keys are grouped together and the values for similar keys are added up to find the occurrences for a particular word.

- It is like an aggregation phase for the keys generated by the map phase. The reducer phase takes the output of shuffle phase as input and then reduces the key-value pairs to unique keys with values added up.
- In our example “An elephant is an animal.” is the only word that appears twice in the sentence.
- After the execution of the reduce phase of MapReduce WordCount example program, appears as a key only once but with a count of 2 as shown below -

(an,2)

(animal,1)

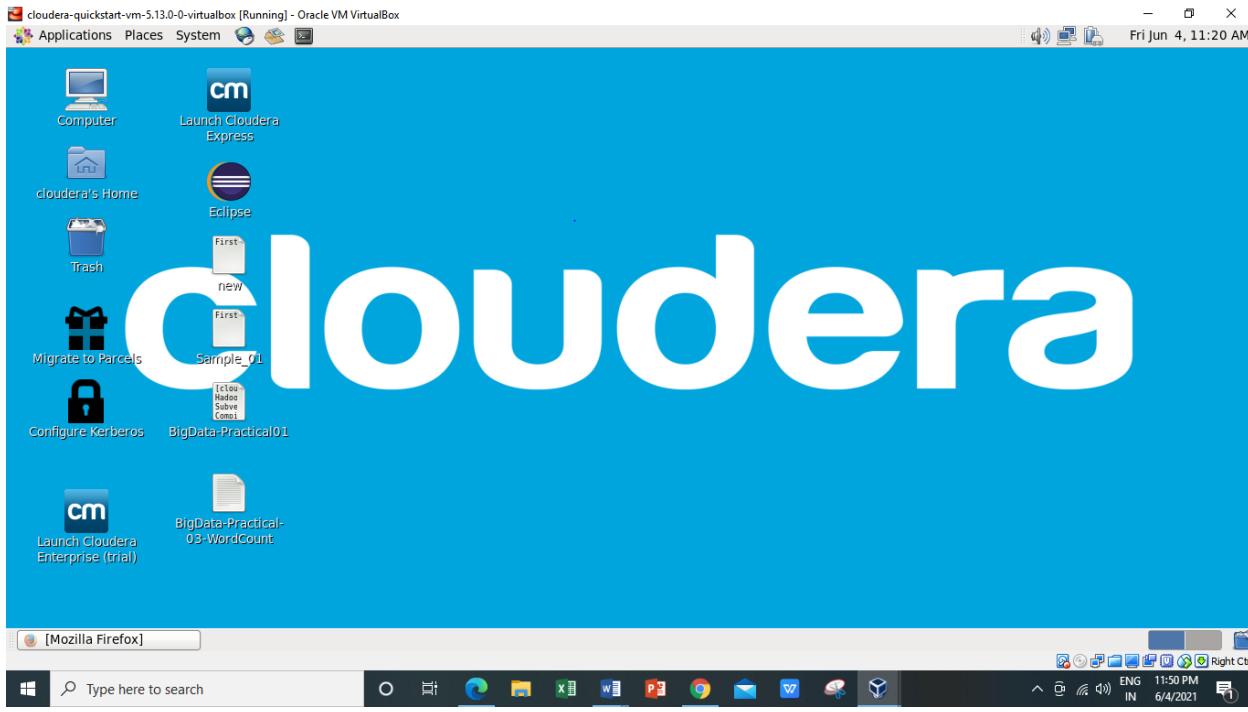
(elephant,1)

(is,1)

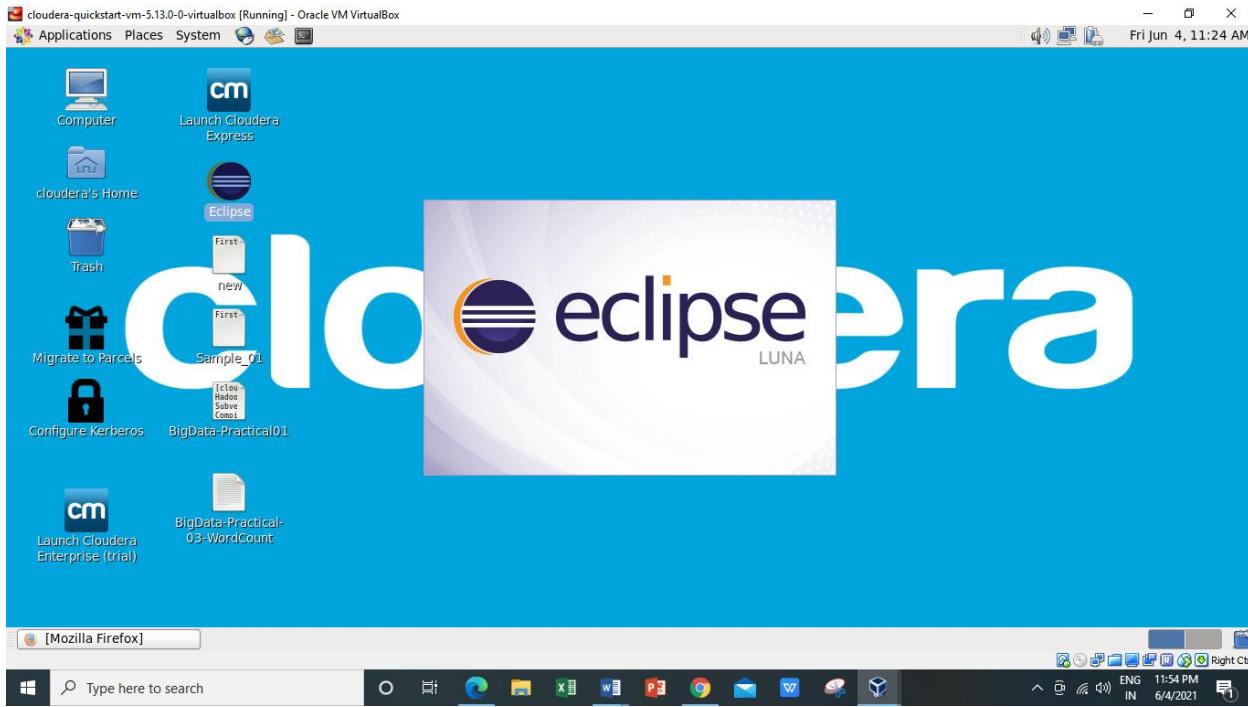
- This is how the MapReduce word count program executes and outputs the number of occurrences of a word in any given input file.
- An important point to note during the execution of the WordCount example is that the mapper class in the WordCount program will execute completely on the entire input file and not just a single sentence.
- Suppose if the input file has 15 lines then the mapper class will split the words of all the 15 lines and form initial key value pairs for the entire dataset.
- The reducer execution will begin only after the mapper phase is executed successfully.

Steps for Word Count in Cloudera: (Without Combiner)

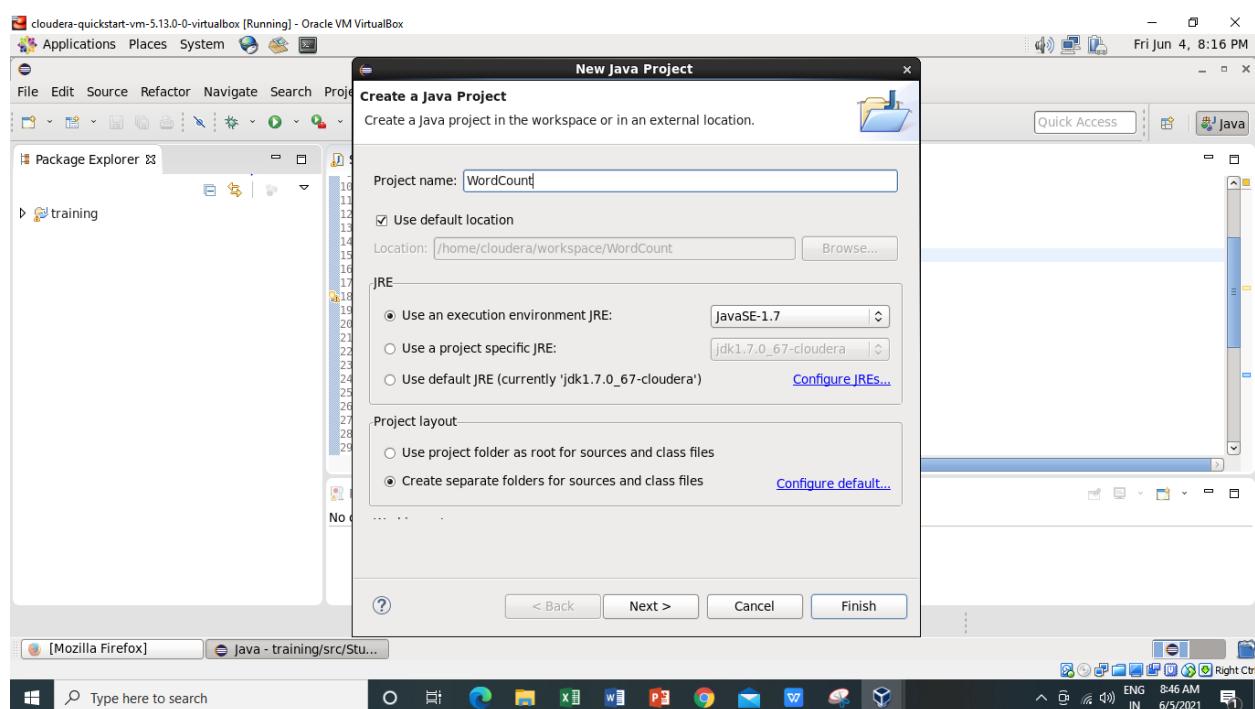
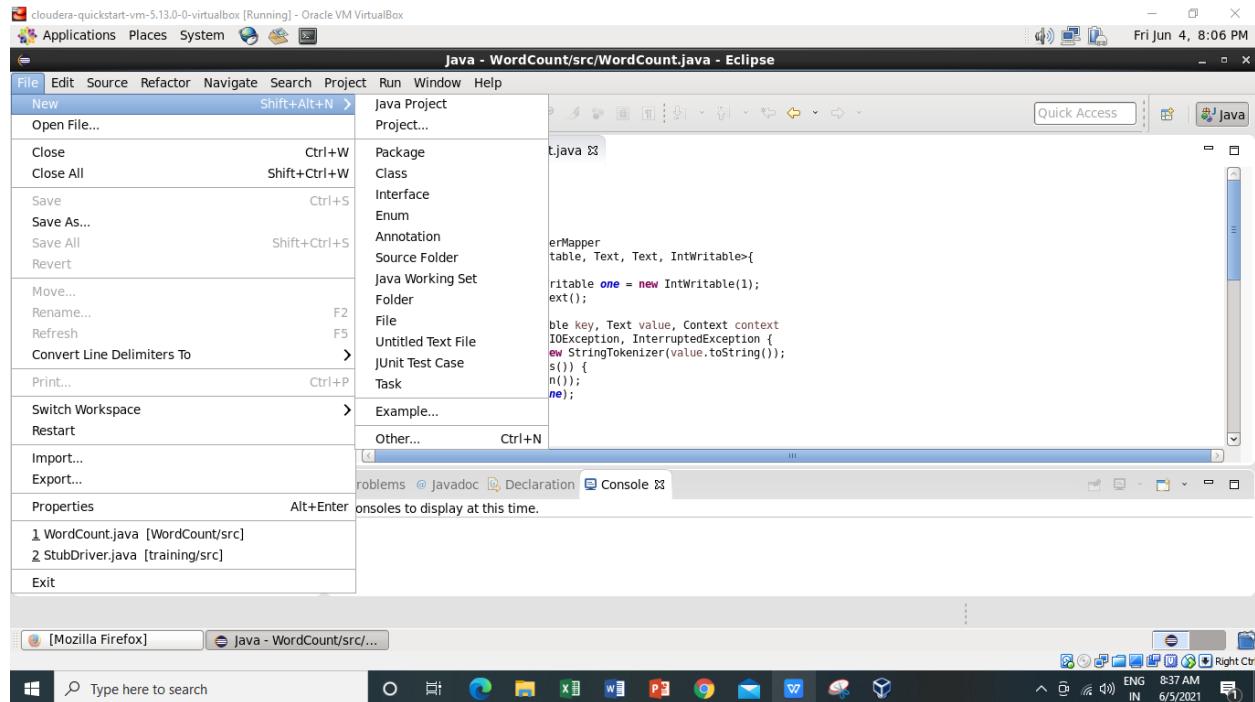
- 1) Open virtual box and then start cloudera quickstart



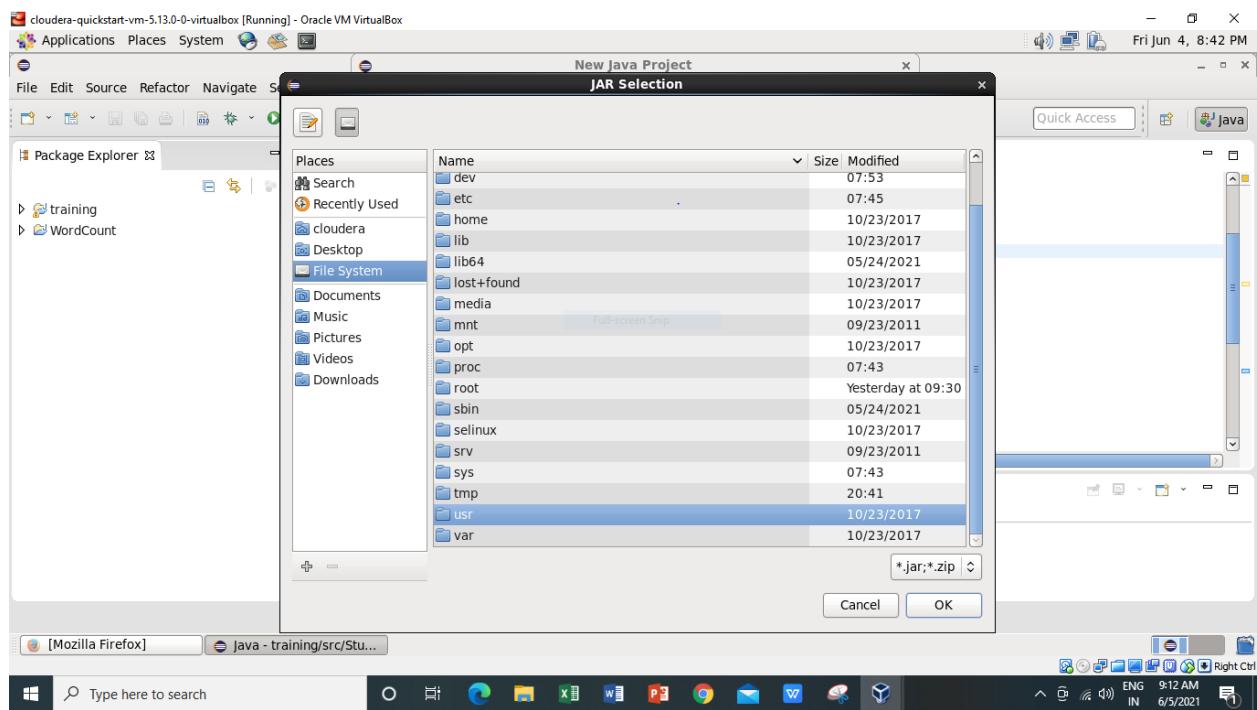
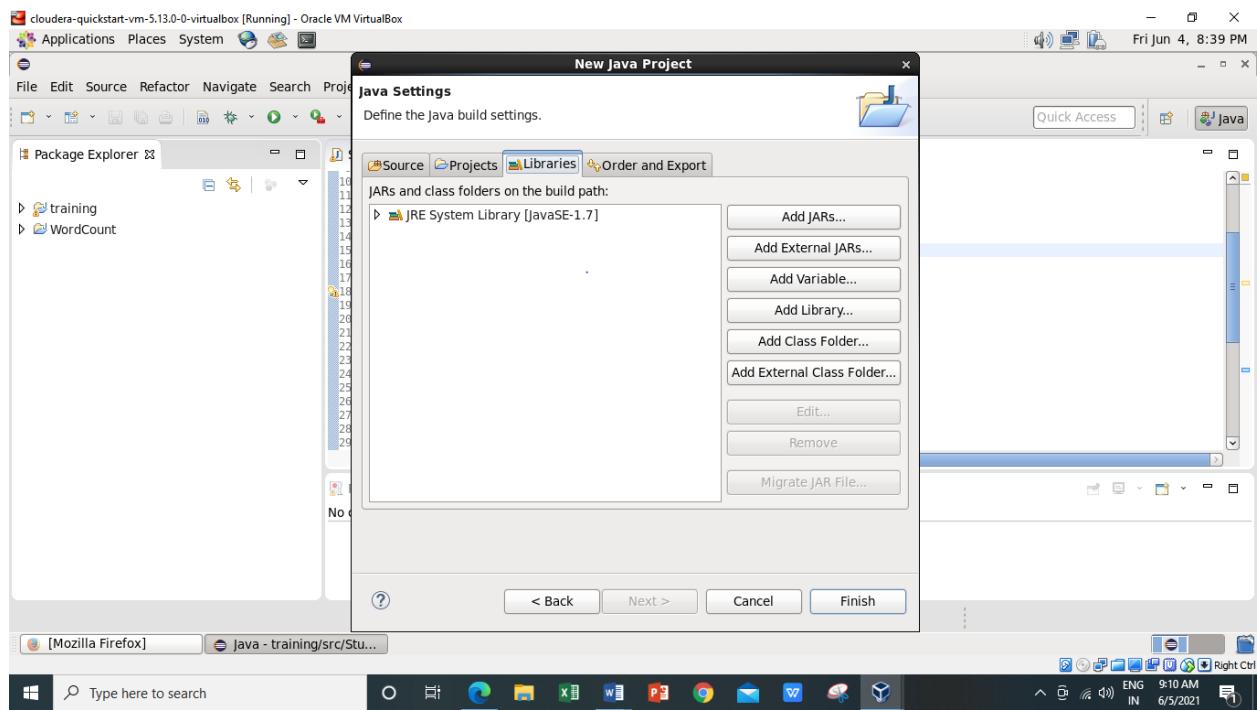
2) Open Eclipse present on the cloudera desktop

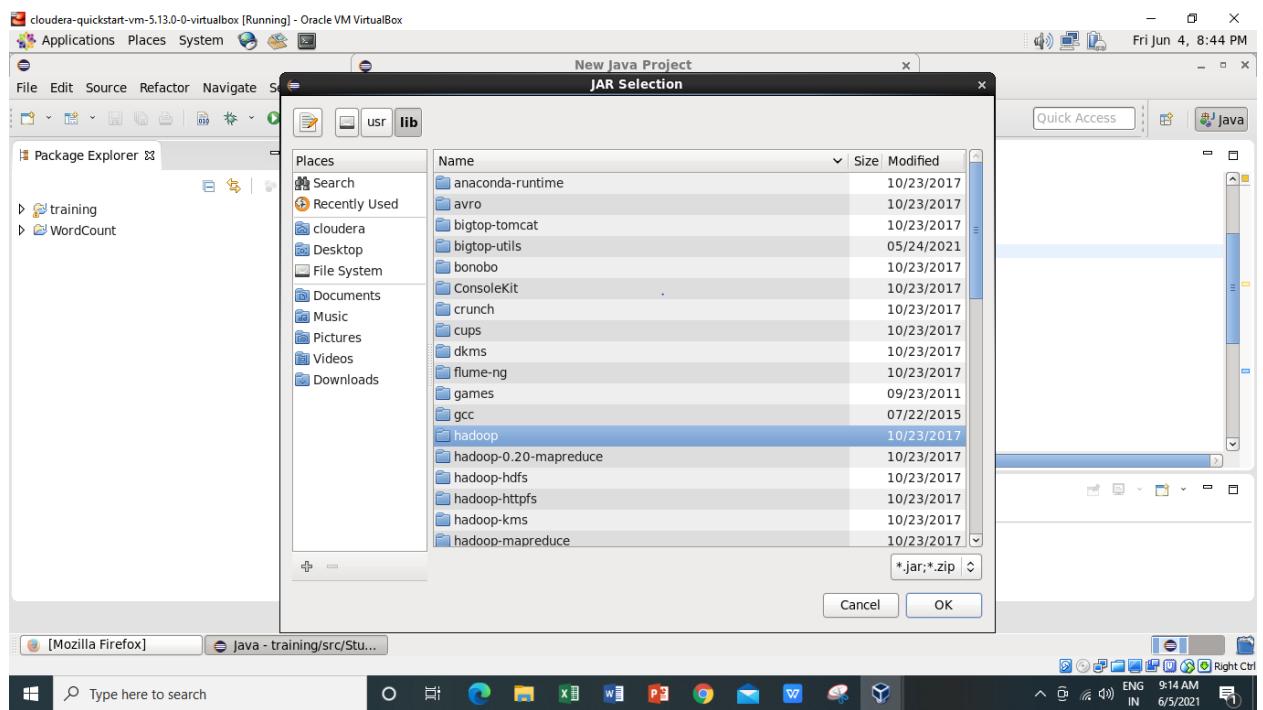
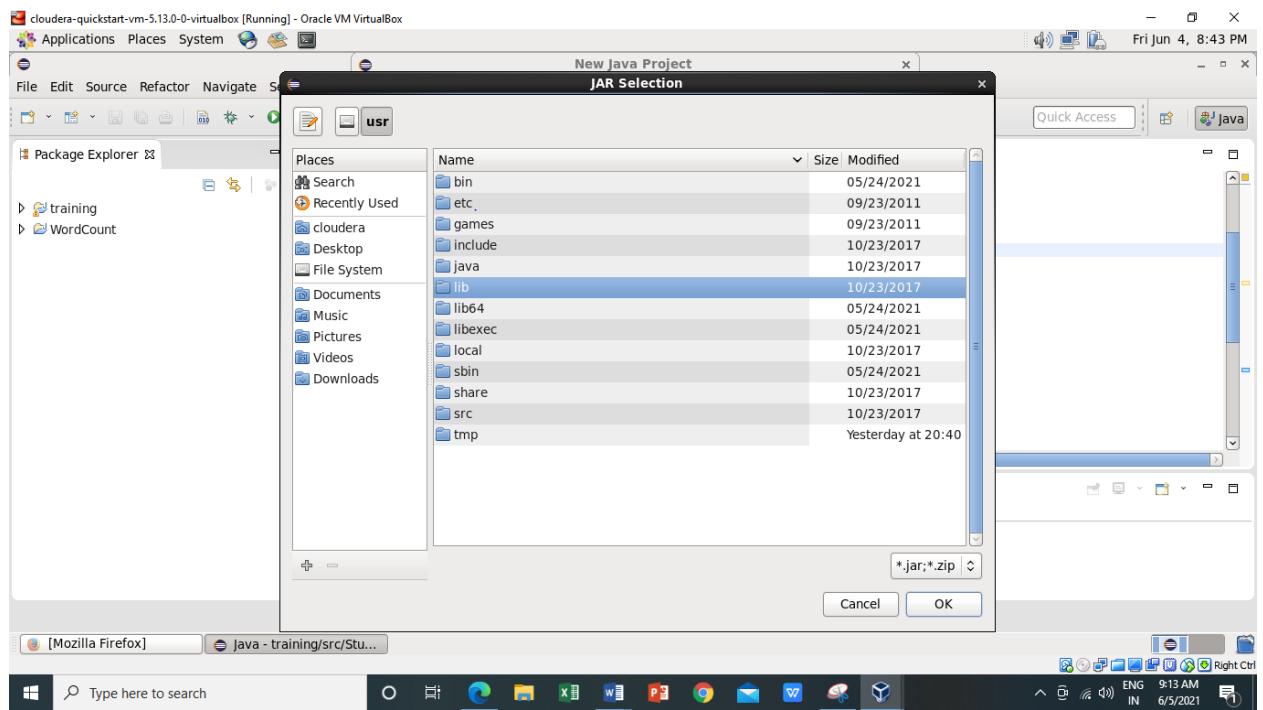


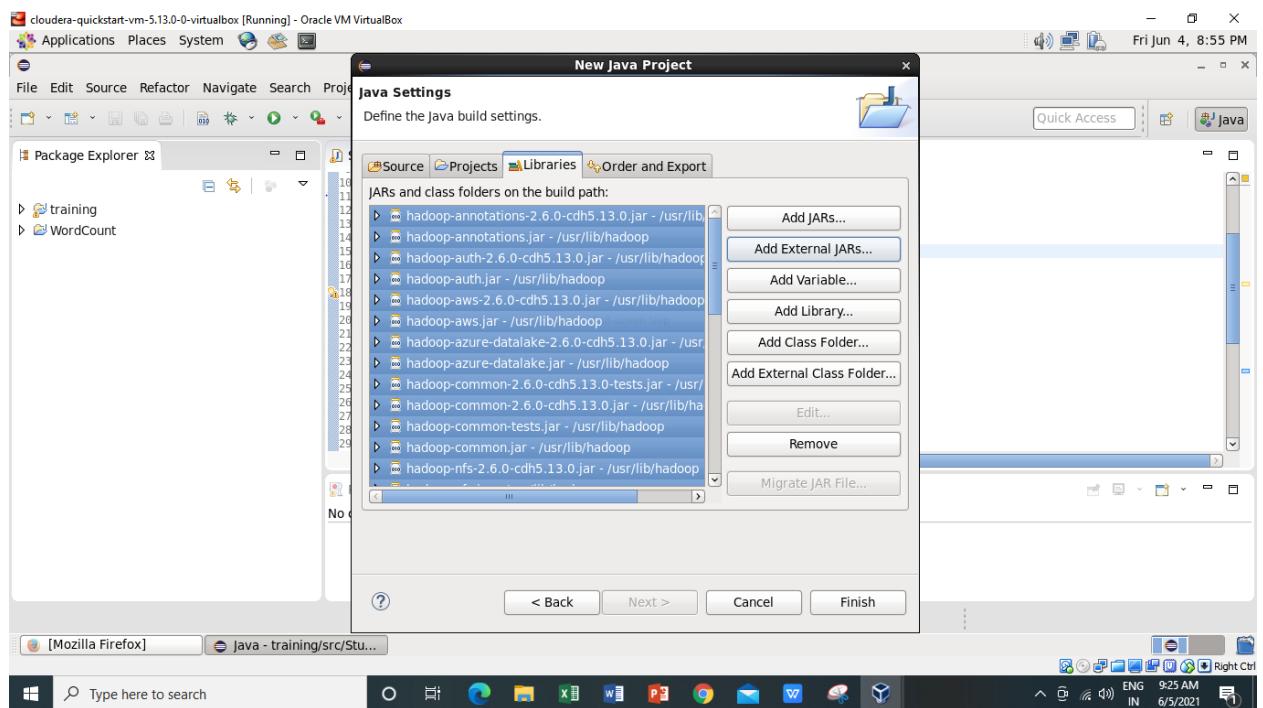
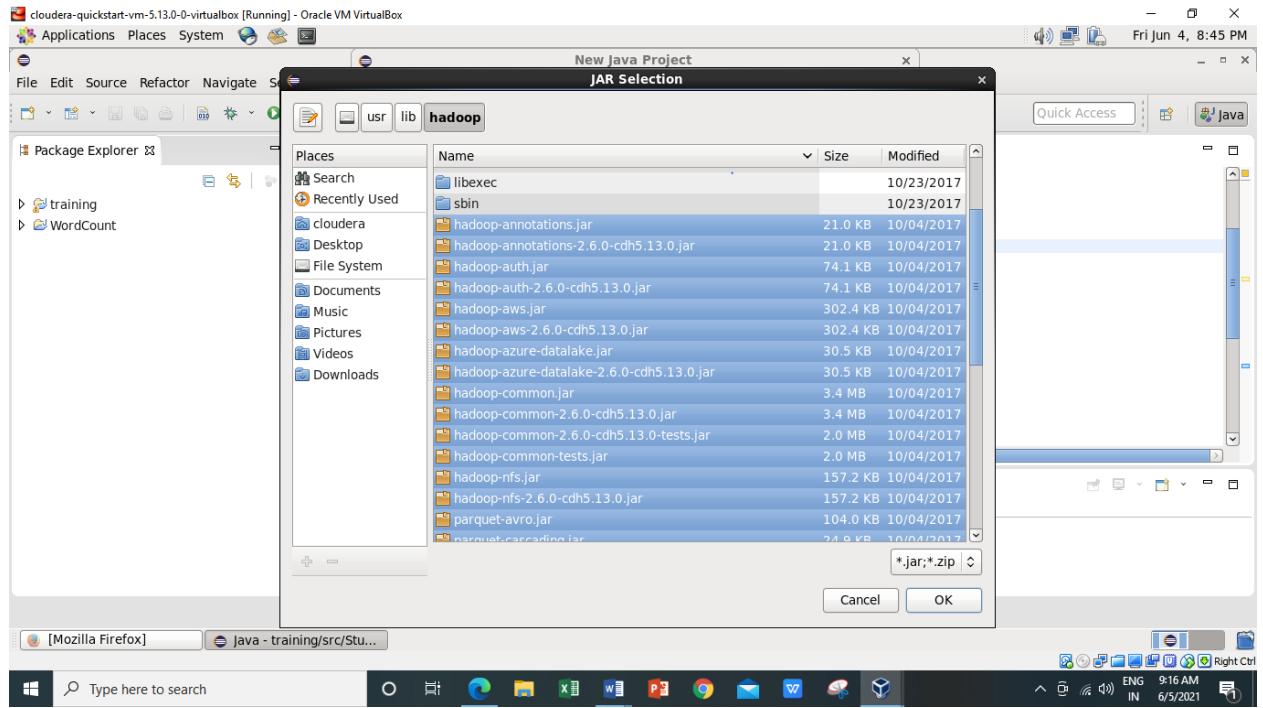
3) Create a new Java project clicking: File -> New -> Project -> Java Project -> Next ("WordCount" is the project name).

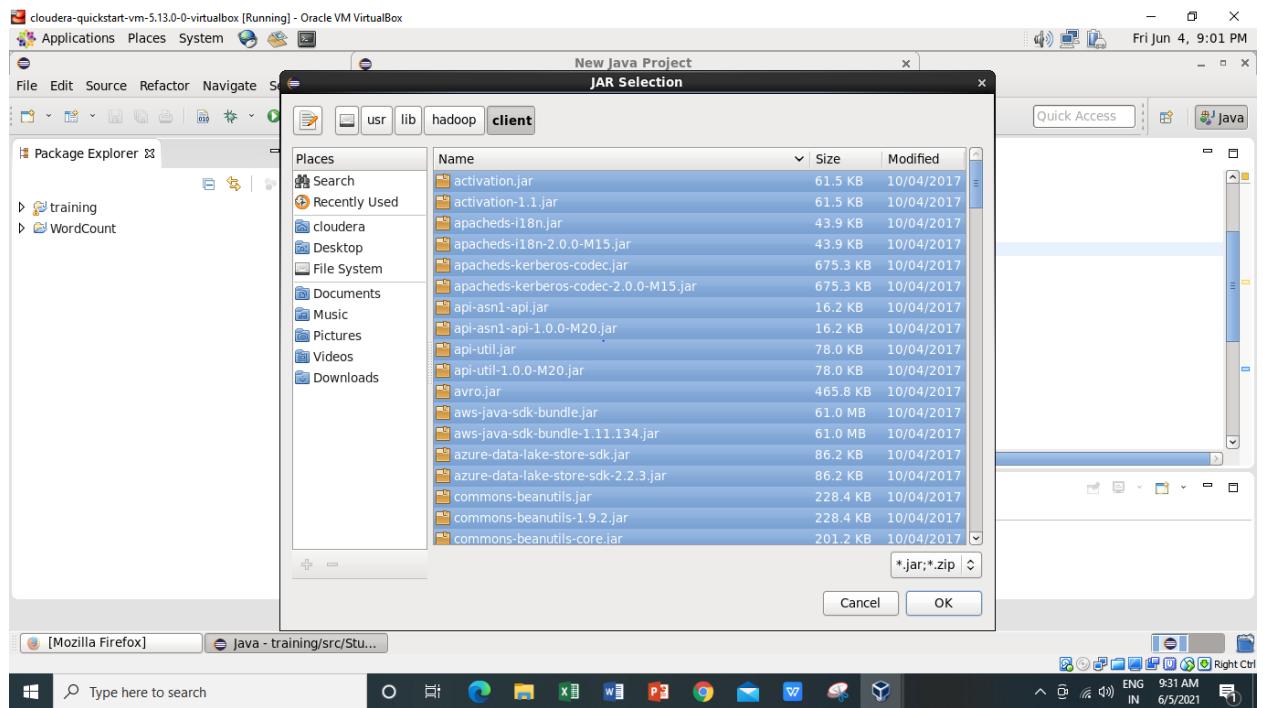
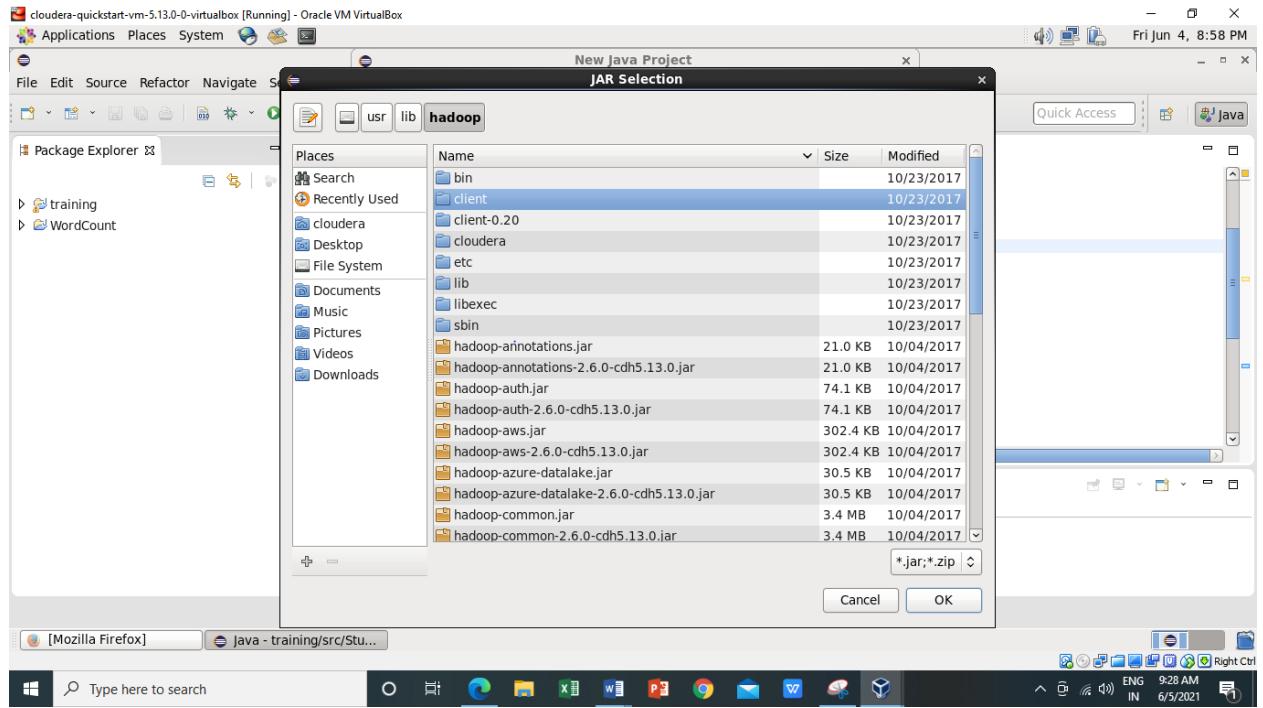


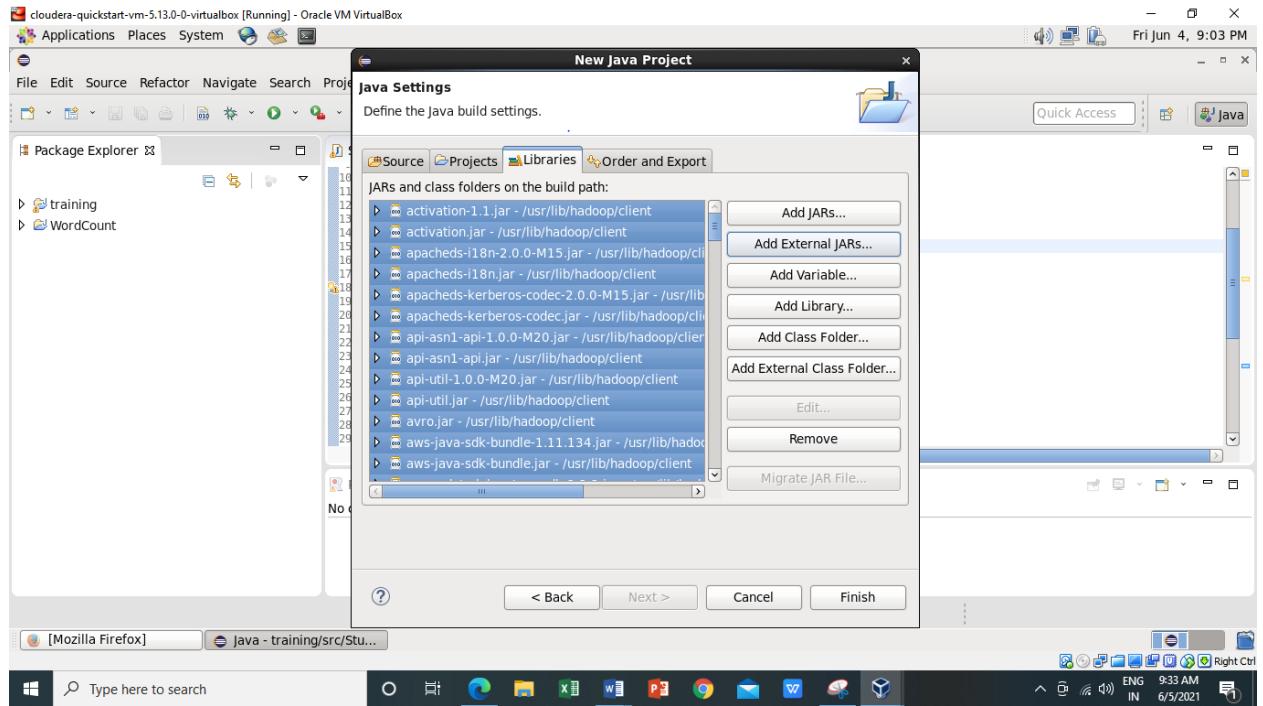
- 4) Adding the Hadoop libraries to the project Click on Libraries -> Add External JARs Click on File System -> usr -> lib -> hadoop Select all the libraries (JAR Files) -> click OK Click on Add External jars, -> client -> select all jar files -> ok -> Finish



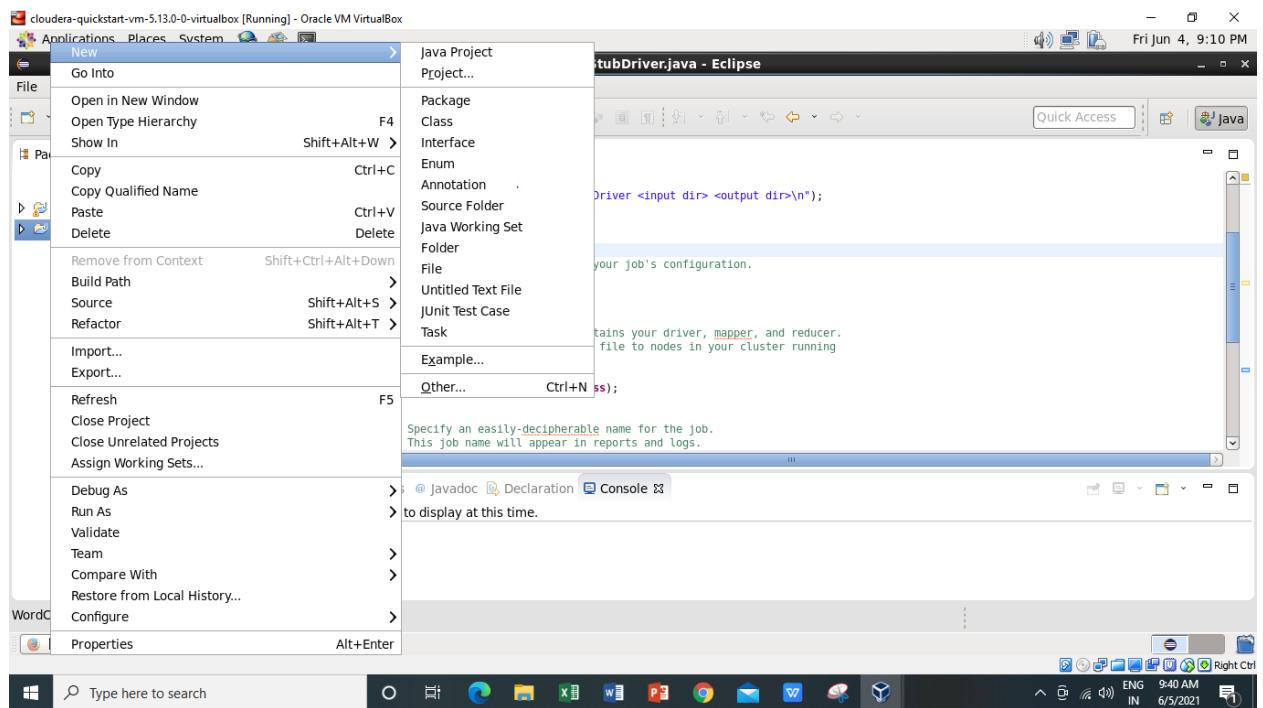


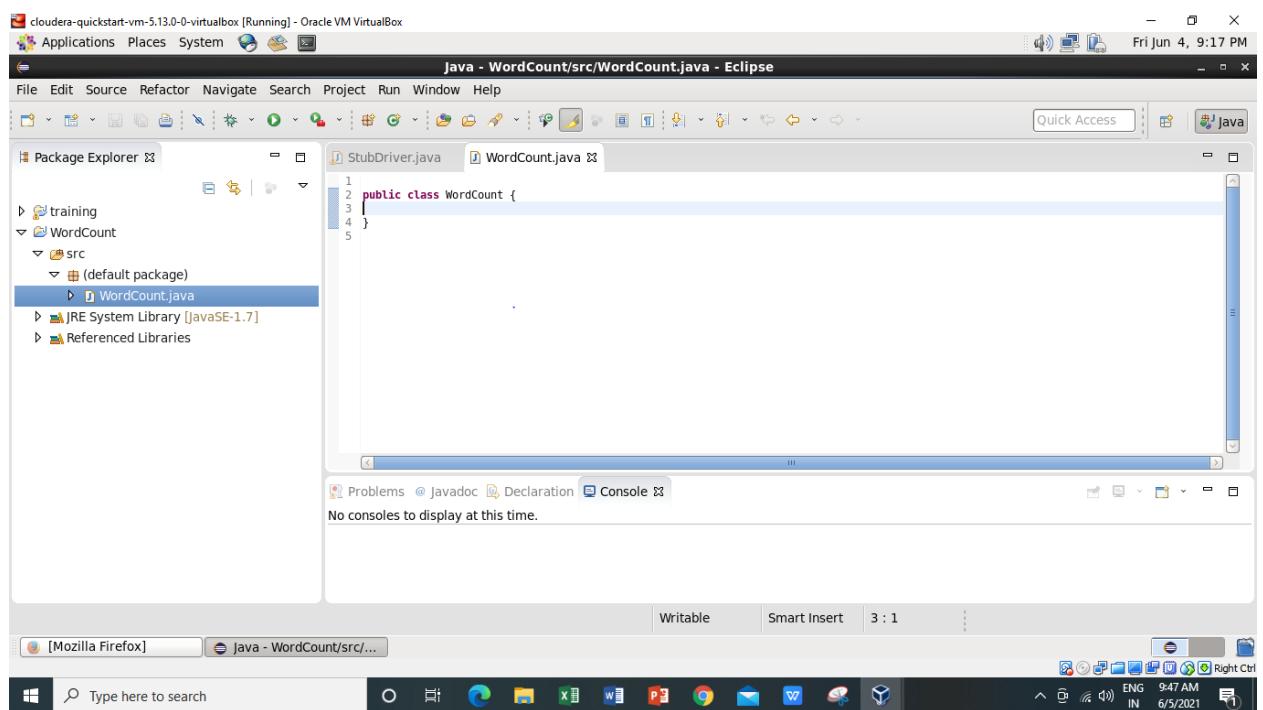
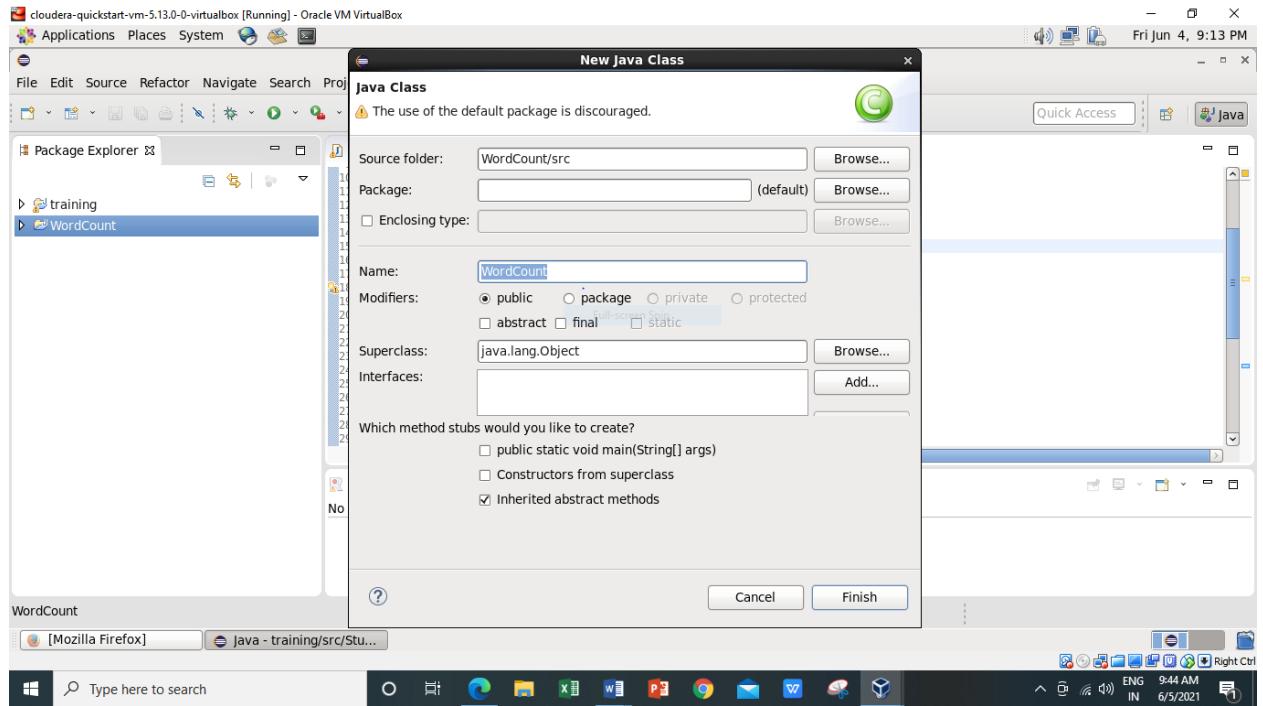


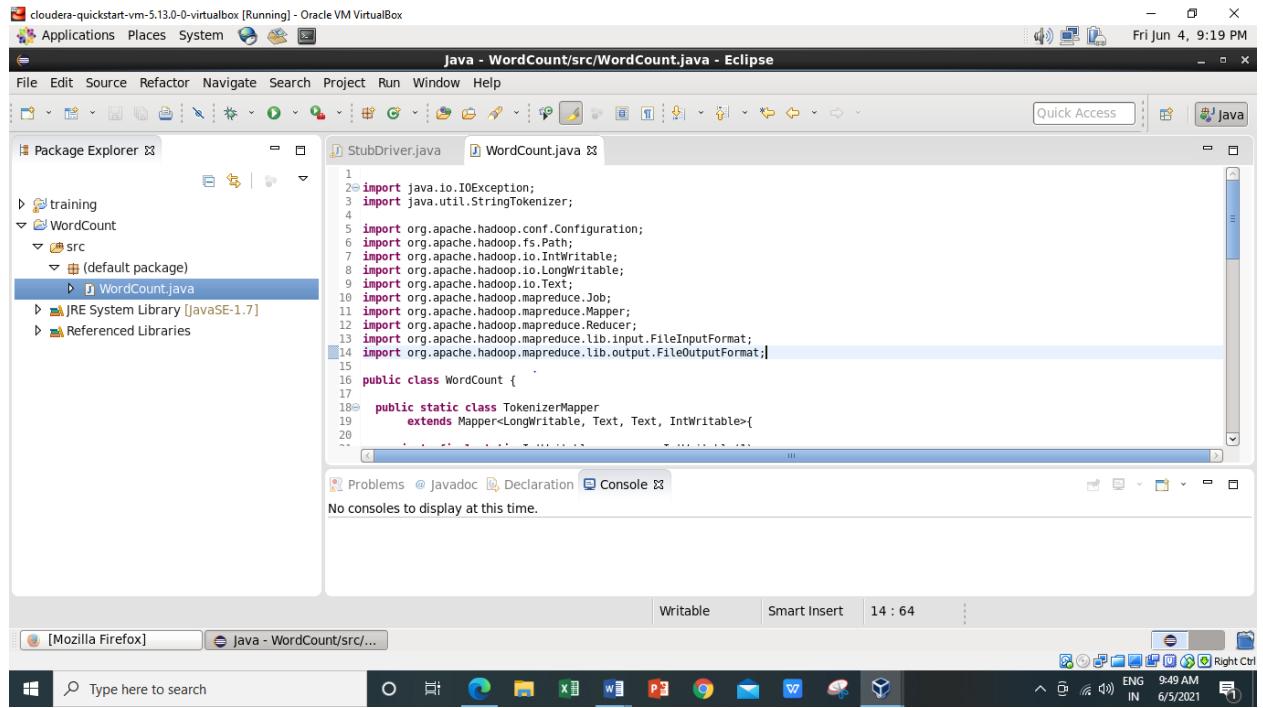




5) Right Click on the name of Project “WordCount” -> New -> class Don’t write anything for package Write Name Textbox write “WordCount” -> Finish Then WordCount.java window will pop up







Source code:

Packages

```
1 import java.io.IOException;
2 import java.util.StringTokenizer;
3
4 import org.apache.hadoop.conf.Configuration;
5 import org.apache.hadoop.fs.Path;
6 import org.apache.hadoop.io.IntWritable;
7 import org.apache.hadoop.io.LongWritable;
8 import org.apache.hadoop.io.Text;
9 import org.apache.hadoop.mapreduce.Job;
10 import org.apache.hadoop.mapreduce.Mapper;
11 import org.apache.hadoop.mapreduce.Reducer;
12 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
13 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
14
```

Mapper Logic

```

15
16 public class WordCount {
17
18     public static class TokenizerMapper
19         extends Mapper<LongWritable; Text, Text, IntWritable>{
20
21         private final static IntWritable one = new IntWritable(1);
22         private Text word = new Text();
23
24         public void map(LongWritable key, Text value, Context context
25                         ) throws IOException, InterruptedException {
26             StringTokenizer itr = new StringTokenizer(value.toString());
27             while (itr.hasMoreTokens()) {
28                 word.set(itr.nextToken());
29                 context.write(word, one);
30             }
31         }
32     }
33 }
```

Reducer logic

```

34     public static class IntSumReducer
35         extends Reducer<Text,IntWritable,Text,IntWritable> {
36         private IntWritable result = new IntWritable();
37
38         public void reduce(Text key, Iterable<IntWritable> values,
39                            Context context
40                            ) throws IOException, InterruptedException {
41             int sum = 0;
42             for (IntWritable val : values) {
43                 sum += val.get();
44             }
45             result.set(sum);
46             context.write(key, result);
47         }
48     }
```

Main function

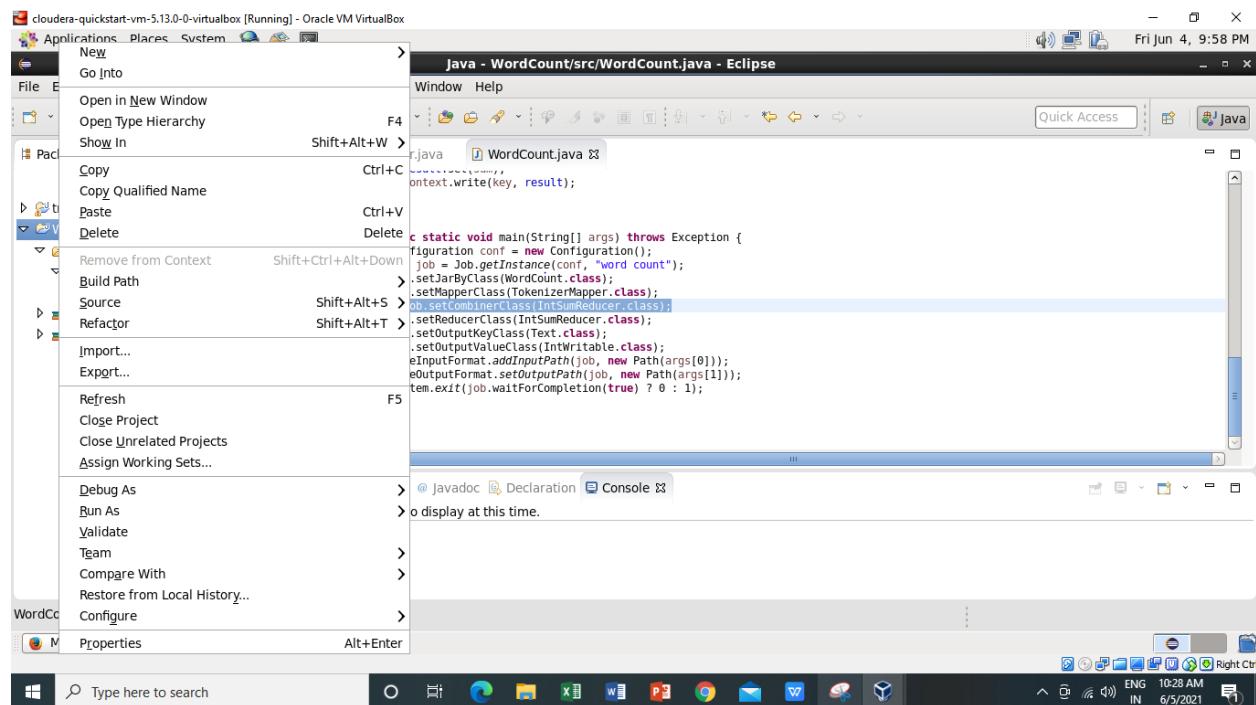
We are running the code without combiner. That is why we commented the combiner line in main function.

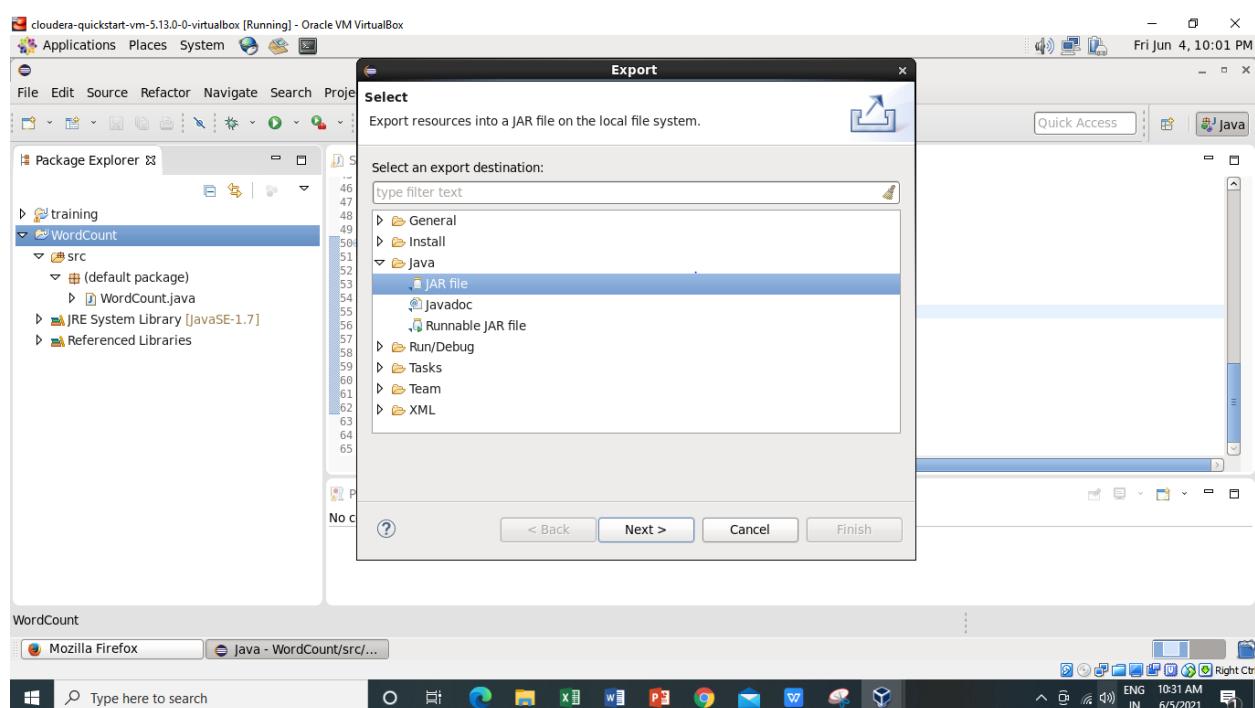
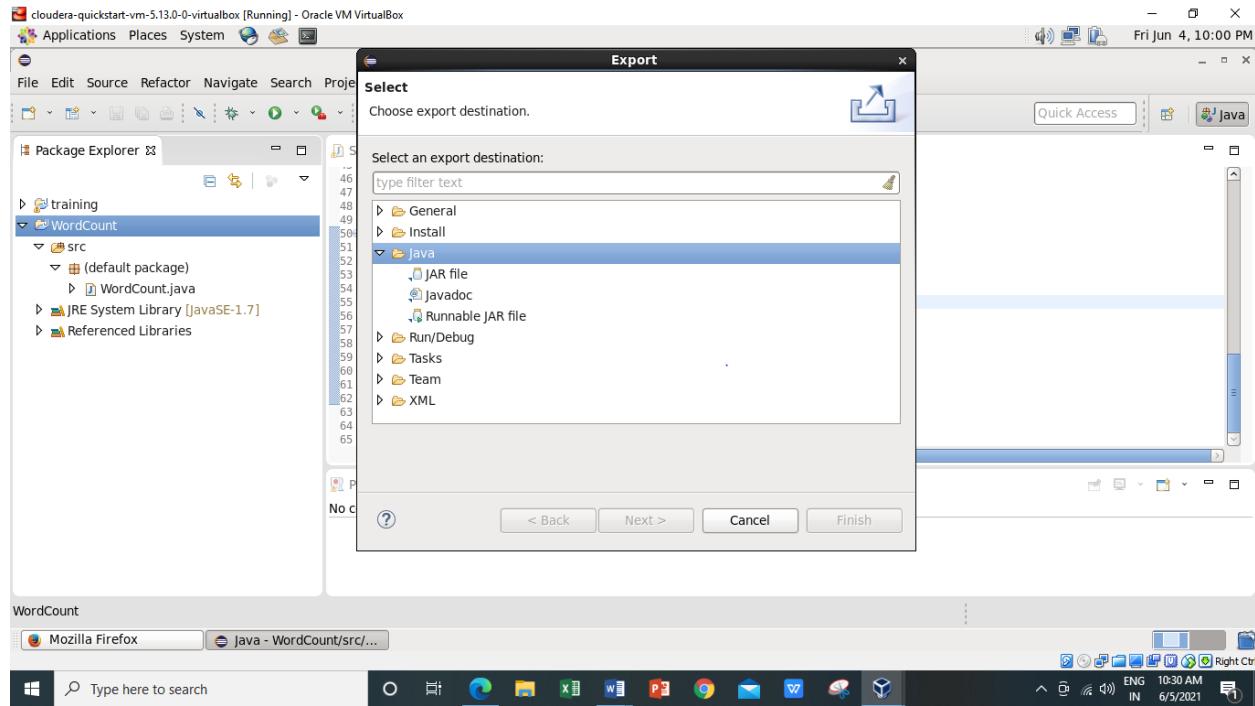
```

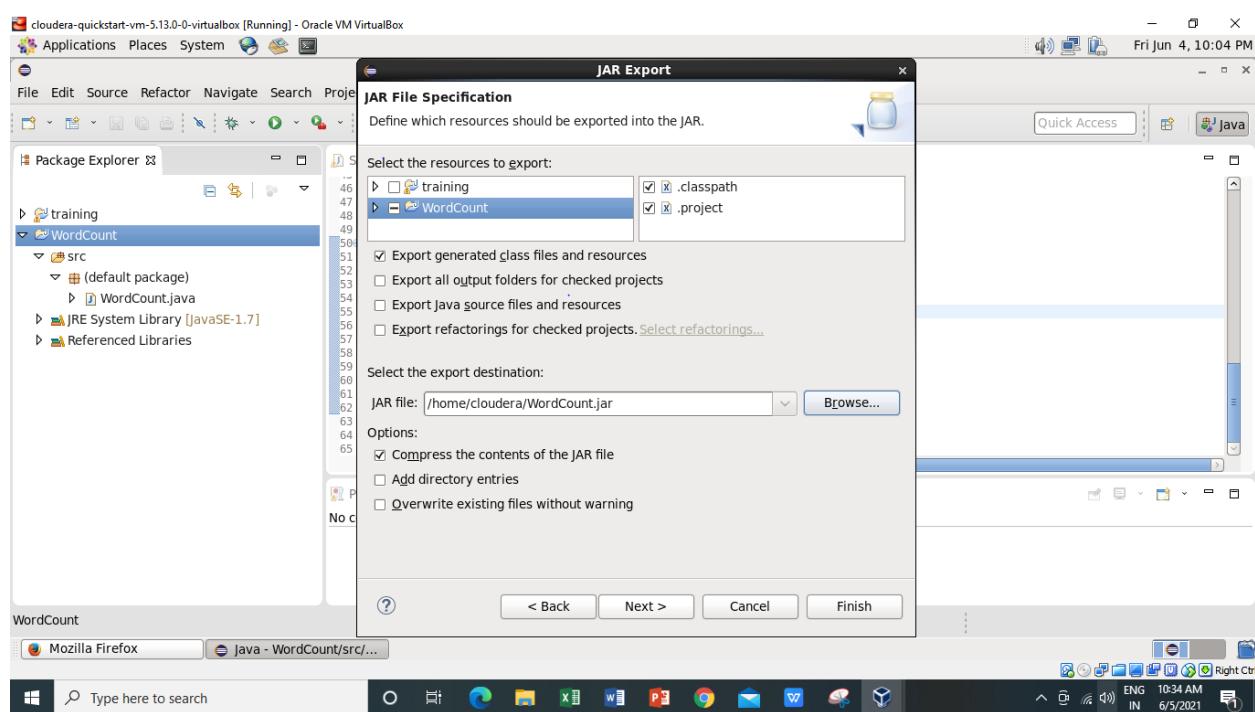
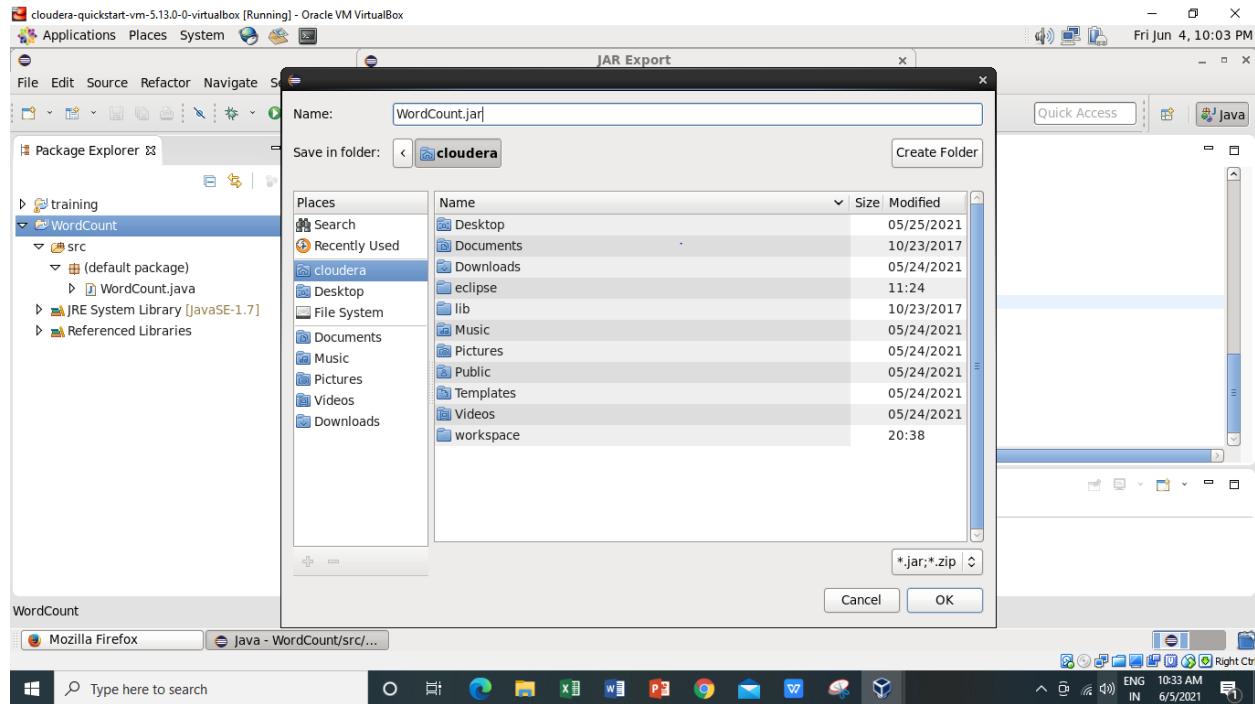
50 public static void main(String[] args) throws Exception {
51     Configuration conf = new Configuration();
52     Job job = Job.getInstance(conf, "word count");
53     job.setJarByClass(WordCount.class);
54     job.setMapperClass(TokenizerMapper.class);
55     //job.setCombinerClass(IntSumReducer.class);
56     job.setReducerClass(IntSumReducer.class);
57     job.setOutputKeyClass(Text.class);
58     job.setOutputValueClass(IntWritable.class);
59     FileInputFormat.addInputPath(job, new Path(args[0]));
60     FileOutputFormat.setOutputPath(job, new Path(args[1]));
61     System.exit(job.waitForCompletion(true) ? 0 : 1);
62 }
63 }
64
65

```

- 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

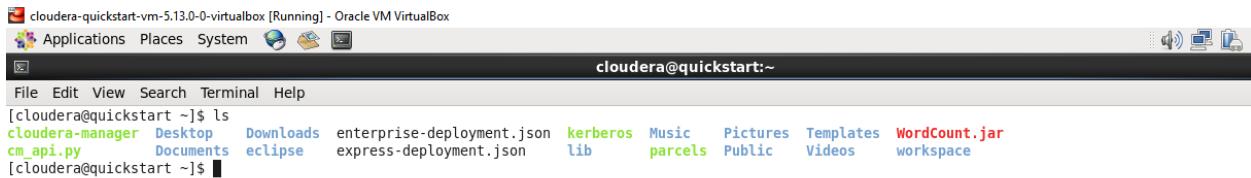






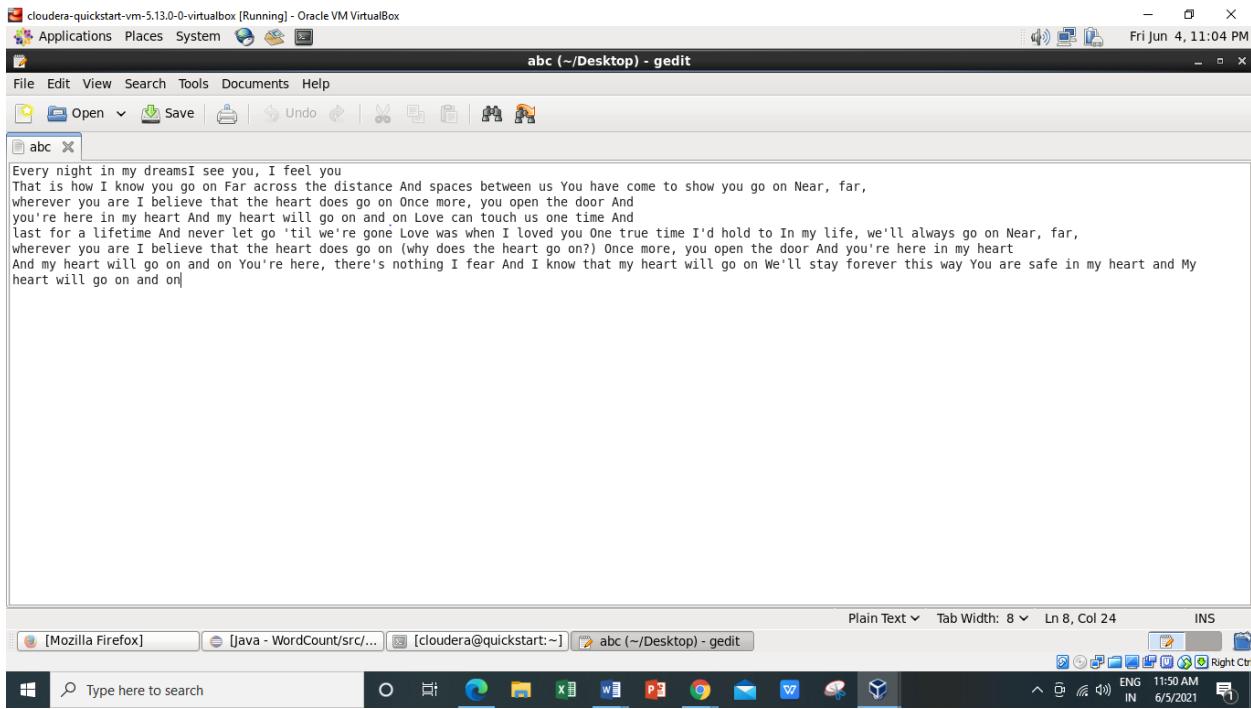
- 7) Verify jar file from terminal by using Open terminal & type " ls " There it will show WordCount.jar
Check current working directory
->pwd

```
[cloudera@quickstart ~]$ pwd  
/home/cloudera  
[cloudera@quickstart ~]$
```



```
cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox  
Applications Places System cloudera@quickstart:~  
File Edit View Search Terminal Help  
[cloudera@quickstart ~]$ ls  
cloudera-manager Desktop Downloads enterprise-deployment.json kerberos Music Pictures Templates WordCount.jar  
cm api.py Documents eclipse express-deployment.json lib parcels Public Videos workspace  
[cloudera@quickstart ~]$
```

- 8) We need to create an input file in local file system
Creating an input file named as “abc”.



Here listing all the directory present in hdfs using **hdfs dfs -ls** / command

```
[cloudera@quickstart ~]$ hdfs dfs -ls /
Found 10 items
-rw-r--r-- 1 cloudera supergroup          27 2021-05-24 12:04 /Sample_01
drwxrwxrwx  - hdfs    supergroup          0 2017-10-23 09:15 /benchmarks
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:58 /forcopy
drwxr-xr-x  - hbase   supergroup          0 2021-06-04 07:57 /hbase
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:20 /newdir
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:36 /rjc
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:55 /solr
drwxrwxrwt  - hdfs   supergroup          0 2021-05-24 10:39 /tmp
drwxr-xr-x  - hdfs   supergroup          0 2017-10-23 09:17 /user
drwxr-xr-x  - hdfs   supergroup          0 2017-10-23 09:17 /var
[cloudera@quickstart ~]$
```

- 9) Now we have to move this input file to hdfs. For this we create a directory on hdfs using command **hdfs dfs -mkdir /inputnew**.

```
[cloudera@quickstart ~]$ hdfs dfs -mkdir /inputdir
[cloudera@quickstart ~]$
```

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

```
[cloudera@quickstart ~]$ hdfs dfs -ls /
Found 11 items
-rw-r--r-- 1 cloudera supergroup          27 2021-05-24 12:04 /Sample_01
drwxrwxrwx  - hdfs    supergroup          0 2017-10-23 09:15 /benchmarks
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:58 /forcopy
drwxr-xr-x  - hbase   supergroup          0 2021-06-04 07:57 /hbase
drwxr-xr-x  - cloudera supergroup          0 2021-06-04 23:34 /inputdir
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:20 /newdir
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:36 /rjc
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:55 /solr
drwxrwxrwt  - hdfs   supergroup          0 2021-05-24 10:39 /tmp
drwxr-xr-x  - hdfs   supergroup          0 2017-10-23 09:17 /user
drwxr-xr-x  - hdfs   supergroup          0 2017-10-23 09:17 /var
[cloudera@quickstart ~]$
```

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

```
[cloudera@quickstart ~]$ hdfs dfs -put /home/cloudera/Desktop/abc /inputdir/
[cloudera@quickstart ~]$
```

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

```
[cloudera@quickstart ~]$ hdfs dfs -ls /inputdir
Found 1 items
-rw-r--r-- 1 cloudera supergroup      813 2021-06-05 00:06 /inputdir/abc
[cloudera@quickstart ~]$
```

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

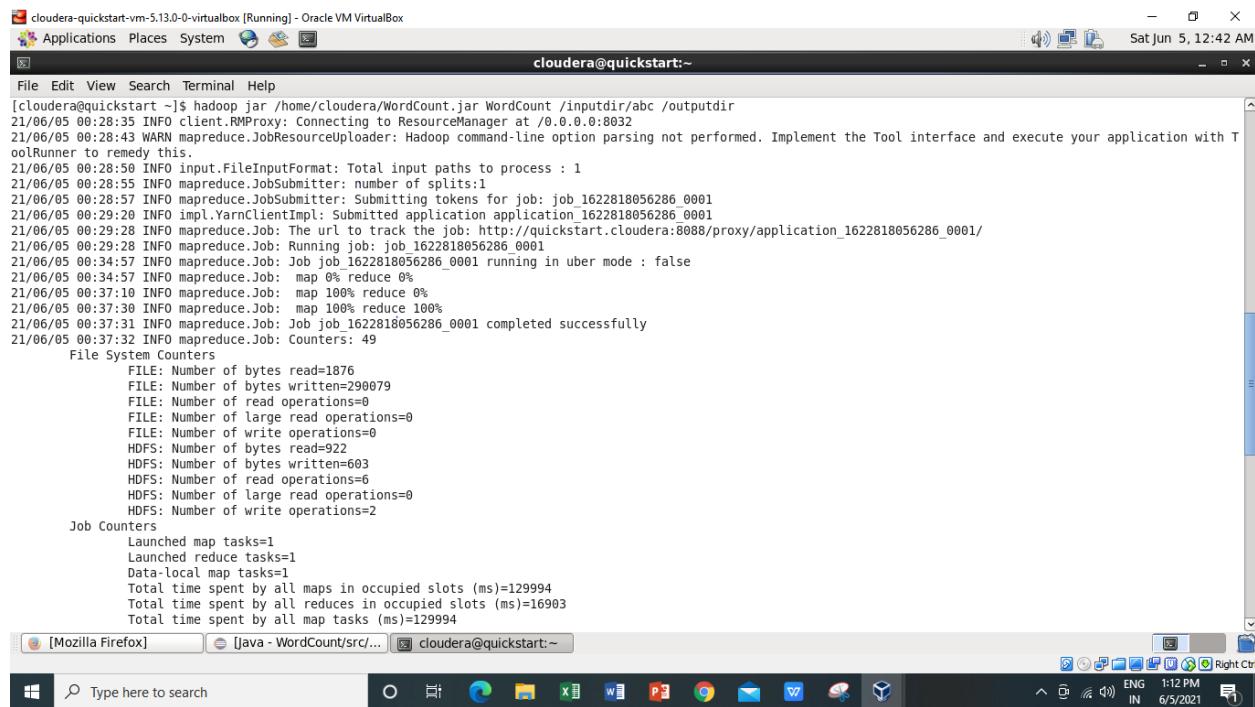
```

-rw-r--r-- 1 cloudera supergroup 0 15 Jun 2021 00:00 00.00 /inputdir/abc
[cloudera@quickstart ~]$ hdfs dfs -cat /inputdir/abc
Every night in my dreams I see you, I feel you
That is how I know you go on Far across the distance And spaces between us You have come to show you go on Near, far,
wherever you are I believe that the heart does go on Once more, you open the door And
you're here in my heart And my heart will go on and on Love can touch us one time And
last for a lifetime And never let go 'til we're gone Love was when I loved you One true time I'd hold to In my life, we'll always go on Near, far,
wherever you are I believe that the heart does go on (why does the heart go on?) Once more, you open the door And you're here in my heart
And my heart will go on and on You're here, there's nothing I fear And I know that my heart will go on We'll stay forever this way You are safe in my heart and My
heart will go on and on
[cloudera@quickstart ~]$ 

```

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**



The screenshot shows a Windows desktop environment with a terminal window open. The terminal window title is "cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox". The terminal session is as follows:

```

cloudera@quickstart ~]$ hadoop jar /home/cloudera/WordCount.jar WordCount /inputdir/abc /outputdir
21/06/05 00:28:35 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
21/06/05 00:28:43 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
21/06/05 00:28:50 INFO input.FileInputFormat: Total input paths to process : 1
21/06/05 00:28:57 INFO mapreduce.JobSubmitter: number of splits:1
21/06/05 00:29:20 INFO impl.YarnClientImpl: Submitted application application_1622818056286_0001
21/06/05 00:29:28 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1622818056286_0001/
21/06/05 00:34:57 INFO mapreduce.Job: Job job_1622818056286_0001 running in uber mode : false
21/06/05 00:34:57 INFO mapreduce.Job: map 0% reduce 0%
21/06/05 00:37:10 INFO mapreduce.Job: map 100% reduce 0%
21/06/05 00:37:30 INFO mapreduce.Job: map 100% reduce 100%
21/06/05 00:37:31 INFO mapreduce.Job: Job job_1622818056286_0001 completed successfully
21/06/05 00:37:32 INFO mapreduce.Job: Counters: 49
  File System Counters
    FILE: Number of bytes read=1876
    FILE: Number of bytes written=290079
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=922
    HDFS: Number of bytes written=603
    HDFS: Number of read operations=0
    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)=129994
    Total time spent by all reduces in occupied slots (ms)=16903
    Total time spent by all map tasks (ms)=129994

```

Below the terminal window, the taskbar shows icons for Mozilla Firefox, Java - WordCount/src/..., and cloudera@quickstart:~. The system tray indicates the date as 6/5/2021 and the time as 1:12 PM.

Map-Reduce Framework

```
Total megabyte-milliseconds taken by all reduce tasks=17308672
Map-Reduce Framework
  Map input records=8
  Map output records=177
  Map output bytes=1516
  Map output materialized bytes=1876
  Input split bytes=109
  Combine input records=0
  Combine output records=0
  Reduce input groups=84
  Reduce shuffle bytes=1876
  Reduce input records=177
  Reduce output records=84
  Spilled Records=354
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=3570
  CPU time spent (ms)=3570
  Physical memory (bytes) snapshot=309850112
  Virtual memory (bytes) snapshot=3015163904
  Total committed heap usage (bytes)=152965120
Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=813
File Output Format Counters
  Bytes Written=603
[cloudera@quickstart ~]$
```

As we can see in the above output,

Combine input records=0

Combine output records=0

We are getting this because we have commented the Combiner line in main function.

And Reduce shuffle bytes coming as,

Reduce shuffle bytes=1876

So when we are not using combiner 1876 bytes acting as an input for the reducer.

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 /

```
[cloudera@quickstart ~]$ hdfs dfs -ls /
Found 12 items
-rw-r--r-- 1 cloudera supergroup          27 2021-05-24 12:04 /Sample_01
drwxrwxrwx  - hdfs    supergroup          0 2017-10-23 09:15 /benchmarks
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:58 /forcopy
drwxr-xr-x  - hbase   supergroup          0 2021-06-04 07:57 /hbase
drwxr-xr-x  - cloudera supergroup          0 2021-06-05 00:06 /inputdir
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:20 /newdir
drwxr-xr-x  - cloudera supergroup          0 2021-06-05 00:37 /outputdir
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:36 /rjc
drwxr-xr-x  - cloudera supergroup          0 2021-05-24 13:55 /solr
drwxrwxrwt  - hdfs   supergroup          0 2021-05-24 10:39 /tmp
drwxr-xr-x  - hdfs   supergroup          0 2017-10-23 09:17 /user
drwxr-xr-x  - hdfs   supergroup          0 2017-10-23 09:17 /var
[cloudera@quickstart ~]$ █
```

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

```
[cloudera@quickstart ~]$ hdfs dfs -ls /outputdir
Found 2 items
-rw-r--r-- 1 cloudera supergroup      0 2021-06-05 00:37 /outputdir/_SUCCESS
-rw-r--r-- 1 cloudera supergroup  603 2021-06-05 00:37 /outputdir/part-r-00000
[cloudera@quickstart ~]$ █
```

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**

```
[cloudera@quickstart ~]$ hdfs dfs -cat /outputdir/part-r-00000
'til      1
(why      1
And      8
Every    1
Far      1
I        7
I'd      1
In       1
Love     2
My       1
Near,    2
Once     2
One      1
That     1
We'll   1
You     2
You're  1
a        1
across   1
always   1
and      4
are      3
believe  2
between  1
can      1
come     1
distance    1
does     3
door     2
dreamsI  1
```

```
more,    2
my      8
never   1
night   1
nothing 1
on      12
on?)    1
one     1
open    2
safe    1
see     1
show    1
spaces  1
stay    1
that    3
the     6
there's 1
this    1
time    2
to      2
touch   1
true    1
us      2
was     1
way     1
we'll   1
we're   1
when   1
wherever      2
will    4
you     8
you're  2
you,    1
[cloudera@quickstart ~]$ █
```

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->Utilities ->Browse the file system

Browsing HDFS - Mozilla Firefox

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	cloudera	supergroup	27 B	Mon May 24 12:04:47 -0700 2021	1	128 MB	Sample_01
drwxrwxrwx	hdfs	supergroup	0 B	Mon Oct 23 09:15:43 -0700 2017	0	0 B	benchmarks
drwxr-xr-x	cloudera	supergroup	0 B	Mon May 24 13:58:34 -0700 2021	0	0 B	forcopy
drwxr-xr-x	hbase	supergroup	0 B	Fri Jun 04 07:57:07 -0700 2021	0	0 B	hbase
drwxr-xr-x	cloudera	supergroup	0 B	Sat Jun 05 00:06:22 -0700 2021	0	0 B	inputdir
drwxr-xr-x	cloudera	supergroup	0 B	Mon May 24 13:20:10 -0700 2021	0	0 B	newdir
drwxr-xr-x	cloudera	supergroup	0 B	Sat Jun 05 00:37:29 -0700 2021	0	0 B	outputdir
drwxr-xr-x	cloudera	supergroup	0 B	Mon May 24 13:36:01 -0700 2021	0	0 B	rjc
drwxr-xr-x	cloudera	supergroup	0 B	Mon May 24 13:55:18 -0700 2021	0	0 B	solr
drwxrwxrwt	hdfs	supergroup	0 B	Mon May 24 10:39:39 -0700 2021	0	0 B	tmp
drwxr-xr-x	hdfs	supergroup	0 B	Mon Oct 23 09:17:33 -0700 2017	0	0 B	user
drwxr-xr-x	hdfs	supergroup	0 B	Mon Oct 23 09:17:24 -0700 2017	0	0 B	var

Browsing HDFS - Mozilla Firefox

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	cloudera	supergroup	0 B	Sat Jun 05 00:37:29 -0700 2021	1	128 MB	_SUCCESS
-rw-r--r--	cloudera	supergroup	603 B	Sat Jun 05 00:37:28 -0700 2021	1	128 MB	part-r-00000

Browse Directory

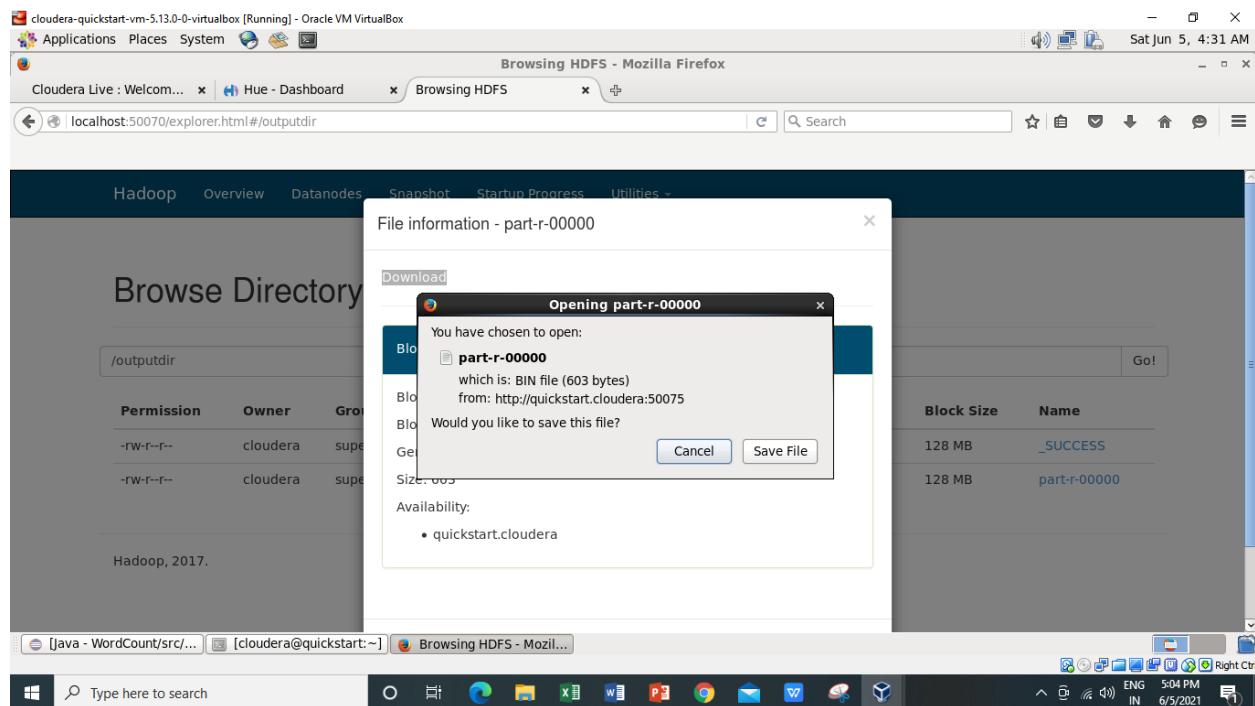
/outputdir

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	cloudera	supergroup	0 B	Sat Jun 05 00:37:29 -0700 2021	1	128 MB	_SUCCESS
-rw-r--r--	cloudera	supergroup	603 B	Sat Jun 05 00:37:28 -0700 2021	1	128 MB	part-r-00000

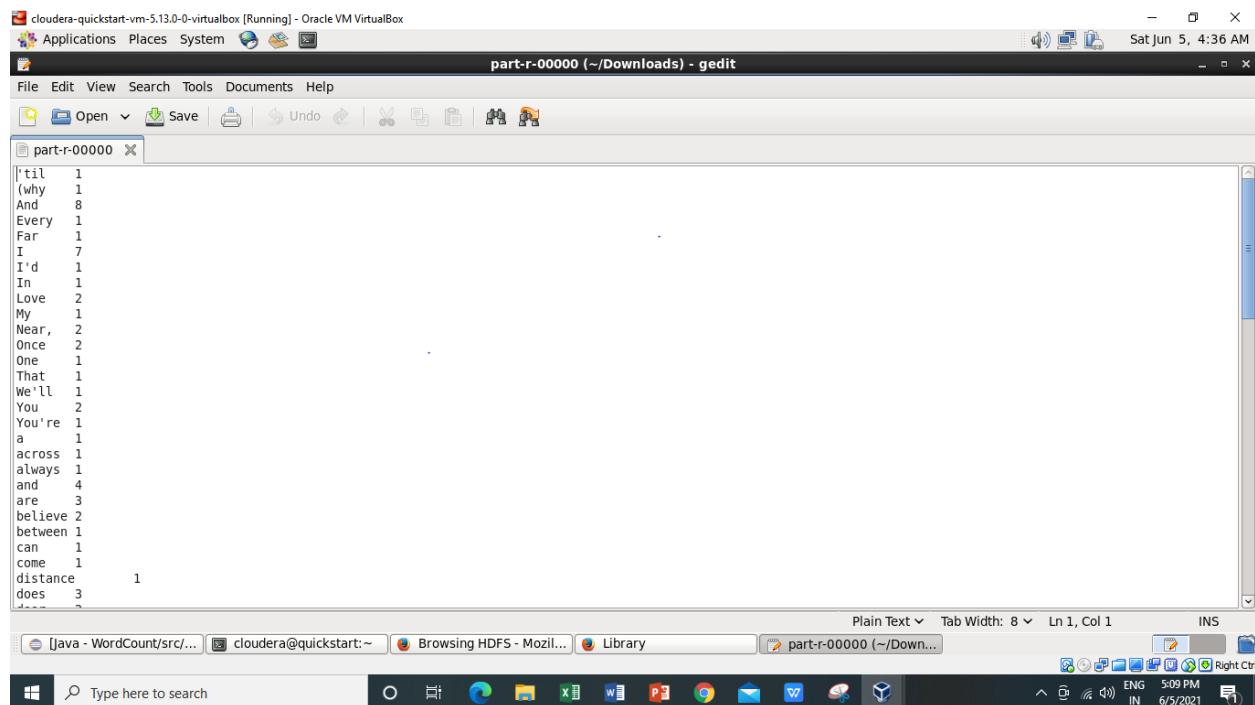
Hadoop, 2017.

Browsing HDFS - Mozilla Firefox

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**



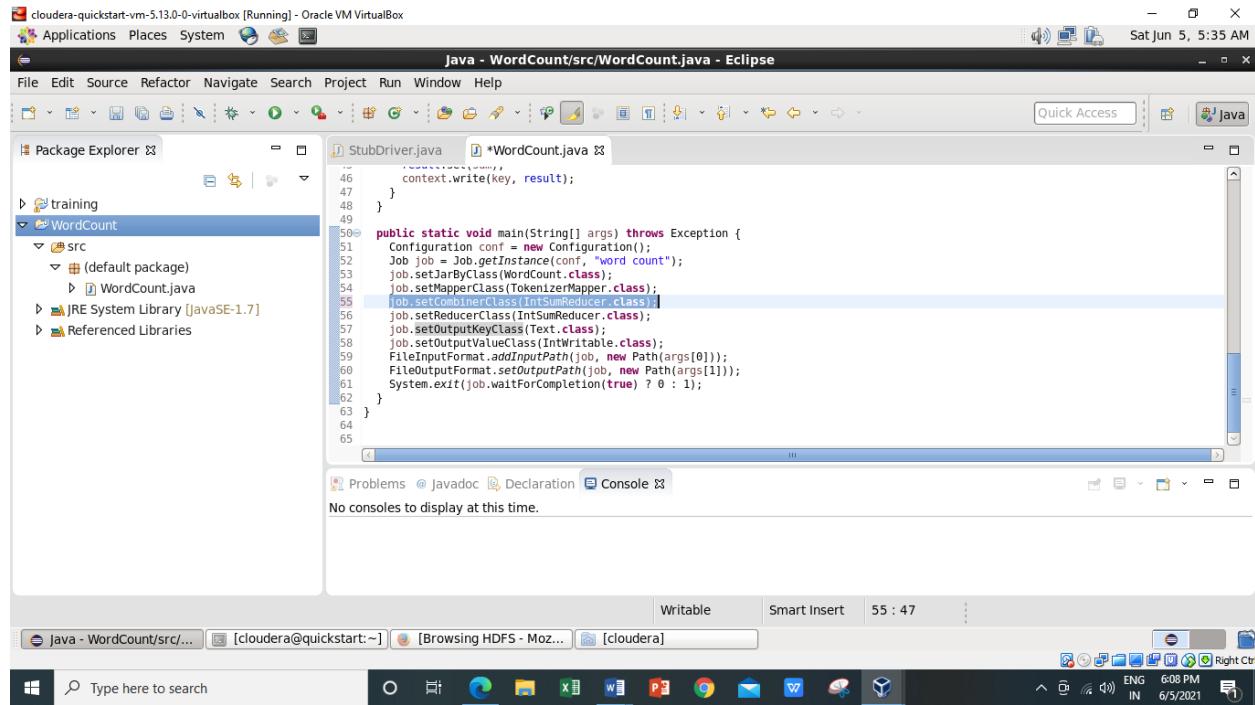
The screenshot shows a Linux desktop environment with several windows open. In the top-left corner, there's a terminal window titled "cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox" displaying the command "hadoop jar wordcount.jar WordCount /user/cloudera/words /user/cloudera/output". Below it is a file viewer window titled "part-r-00000 (~/Downloads) - gedit" showing the word count output as a list of words and their counts. The bottom of the screen features a dock with icons for various applications like Java, a browser, and file management tools. The system tray in the bottom-right corner shows network status, battery level, and system information.

```
word 12
on? 1
one 1
open 2
safe 1
see 1
show 1
spaces 1
stay 1
that 3
the 6
there's 1
this 1
time 2
to 2
touch 1
true 1
us 2
was 1
way 1
we'll 1
we're 1
when 1
wherever 2
will 4
you 8
you're 2
you, 1
```

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 (With Combiner) in Eclipse:

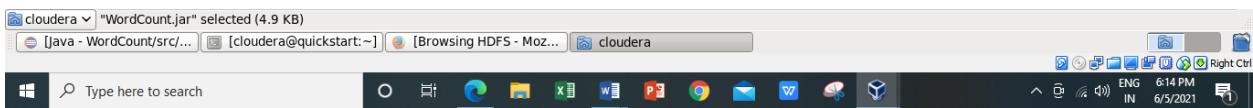
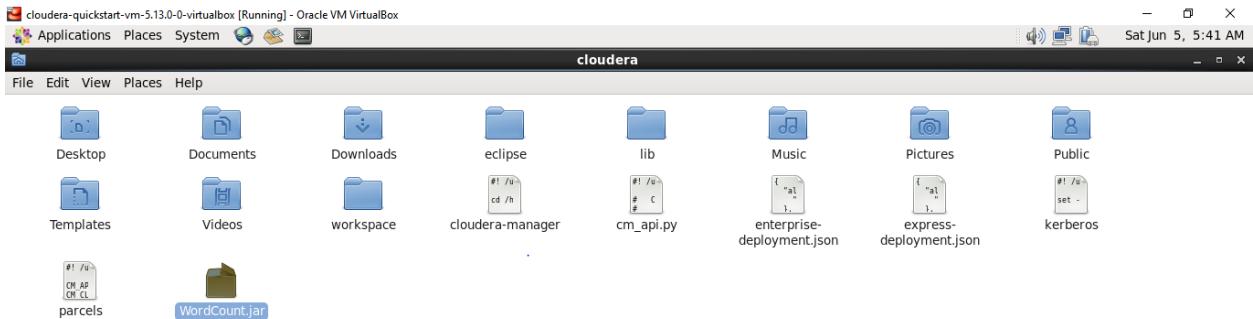
- 1) We will perform the same steps as we have done above for WordCount (without using combiner) in that we just uncommenting the combiner line in main function.



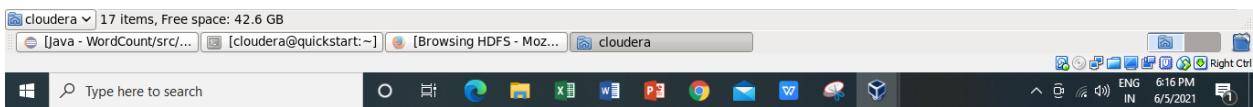
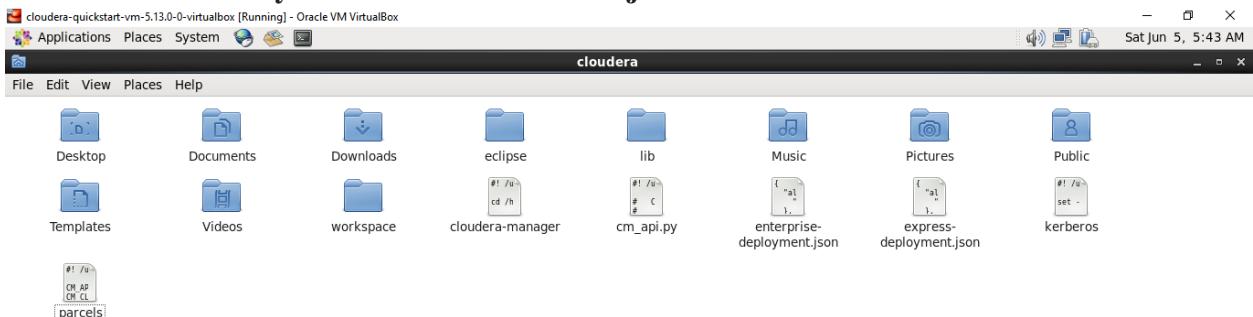
The screenshot shows the Eclipse IDE interface with the following details:

- Title Bar:** cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox
- Menu Bar:** File Edit Source Refactor Navigate Search Project Run Window Help
- Toolbar:** Standard Eclipse toolbar with various icons for file operations.
- Package Explorer:** Shows the project structure with a package named "WordCount" containing a file "WordCount.java".
- Editor:** Displays the Java code for "WordCount.java". The line `job.setCombinerClass(IntSumReducer.class);` is highlighted in blue, indicating it has been uncommented.
- Bottom Status Bar:** Shows the current time as Sat Jun 5, 5:35 AM.
- Taskbar:** Shows multiple open tasks including "Java - WordCount/src/...", "[cloudera@quickstart:~]", "[Browsing HDFS - Mozilla...]", and "[cloudera]".
- System Tray:** Shows system status icons like battery, signal strength, and date/time (6/5/2021, 6:08 PM).

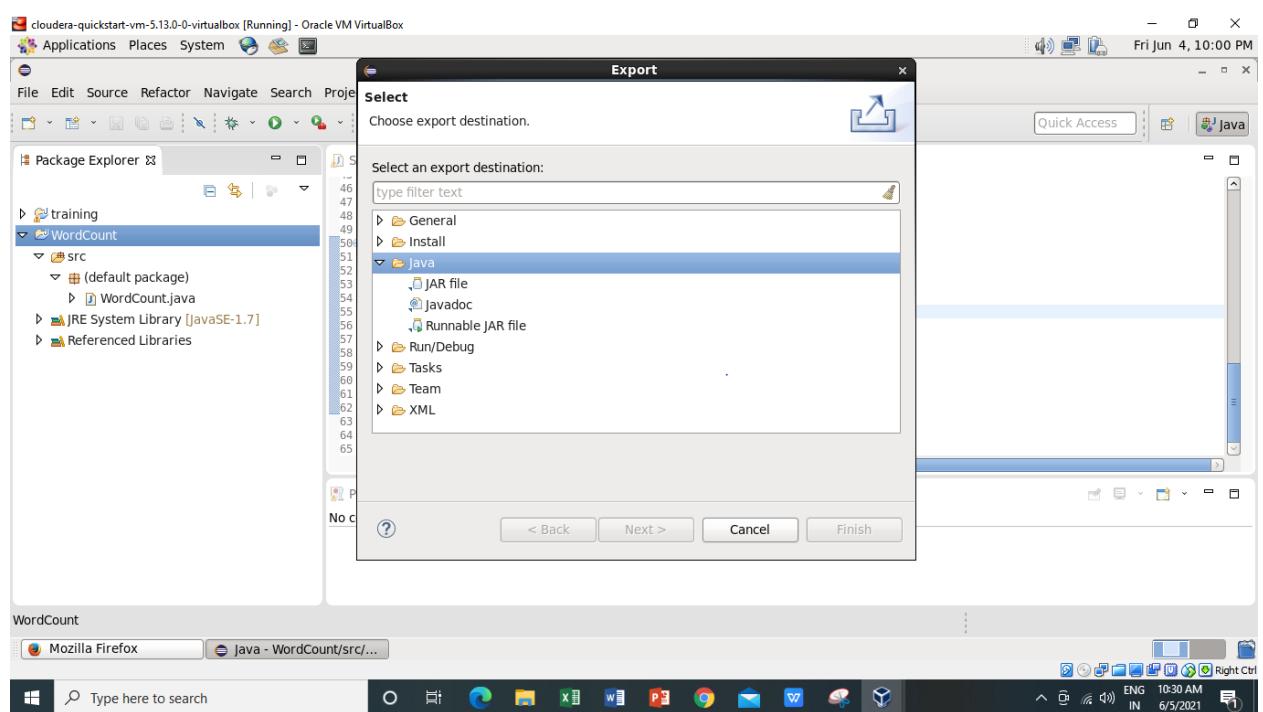
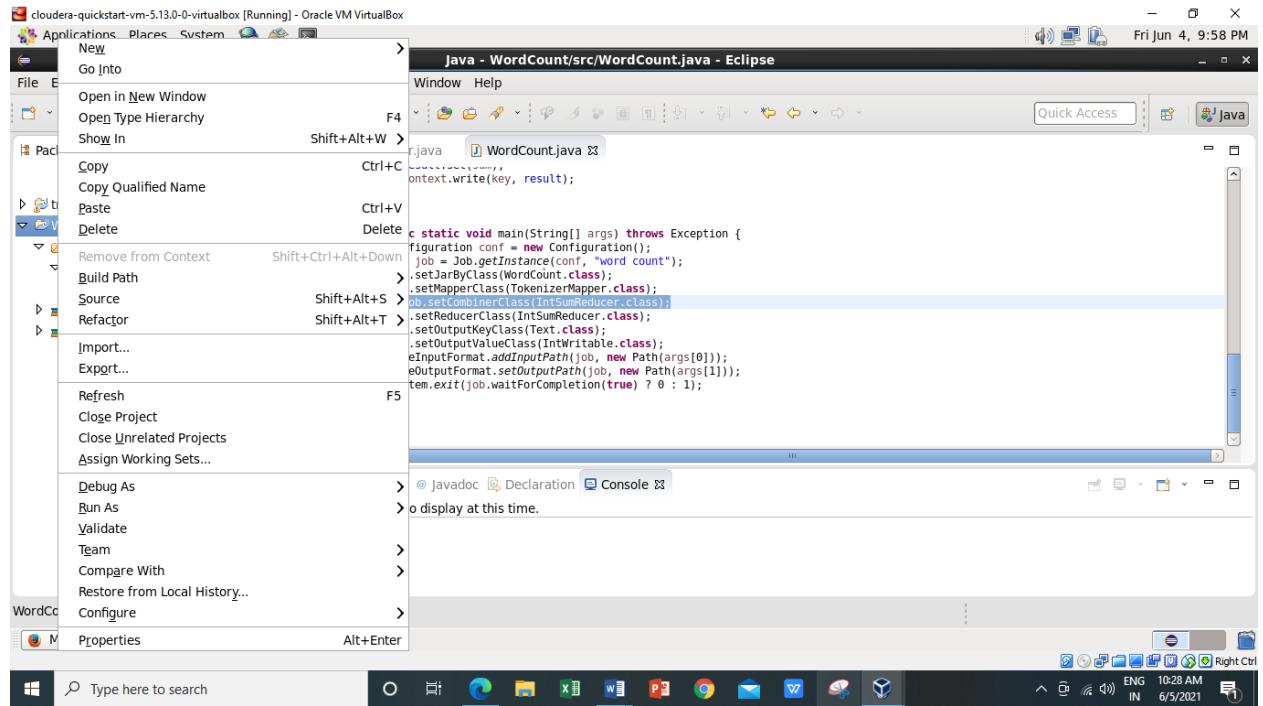
- 2) And will delete the WordCount.jar file in which all jar files are present from **/home/cloudera**.

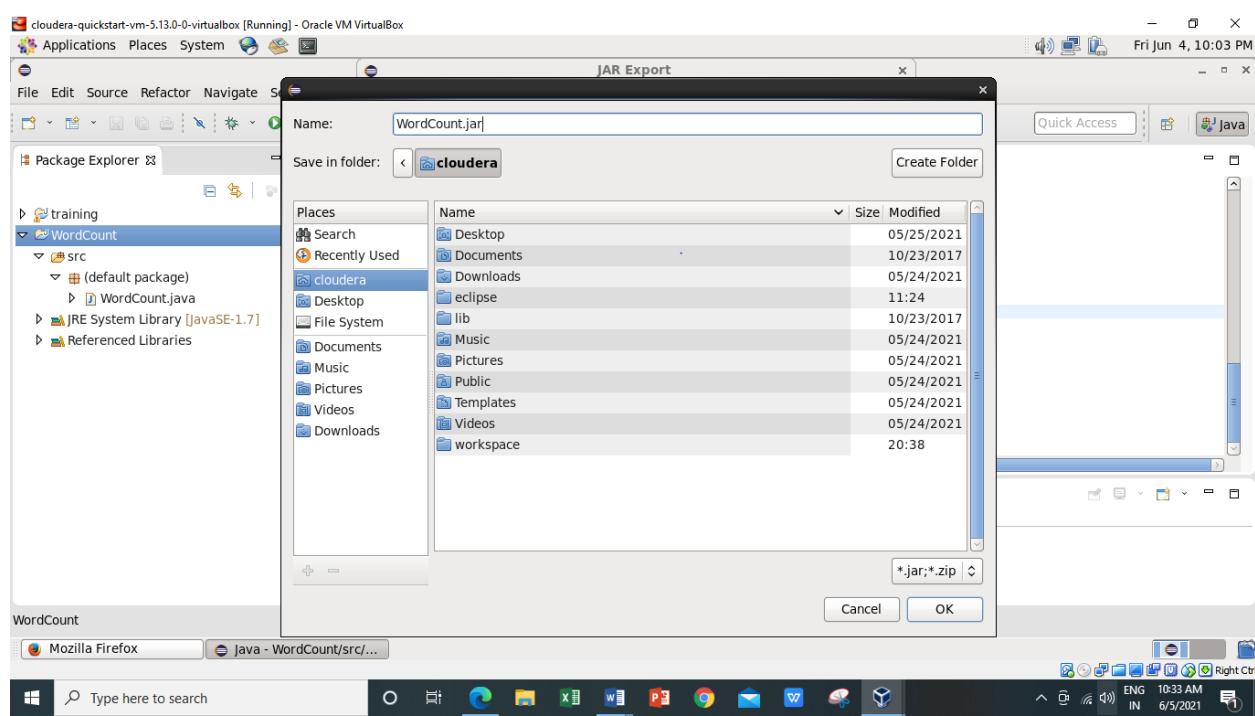
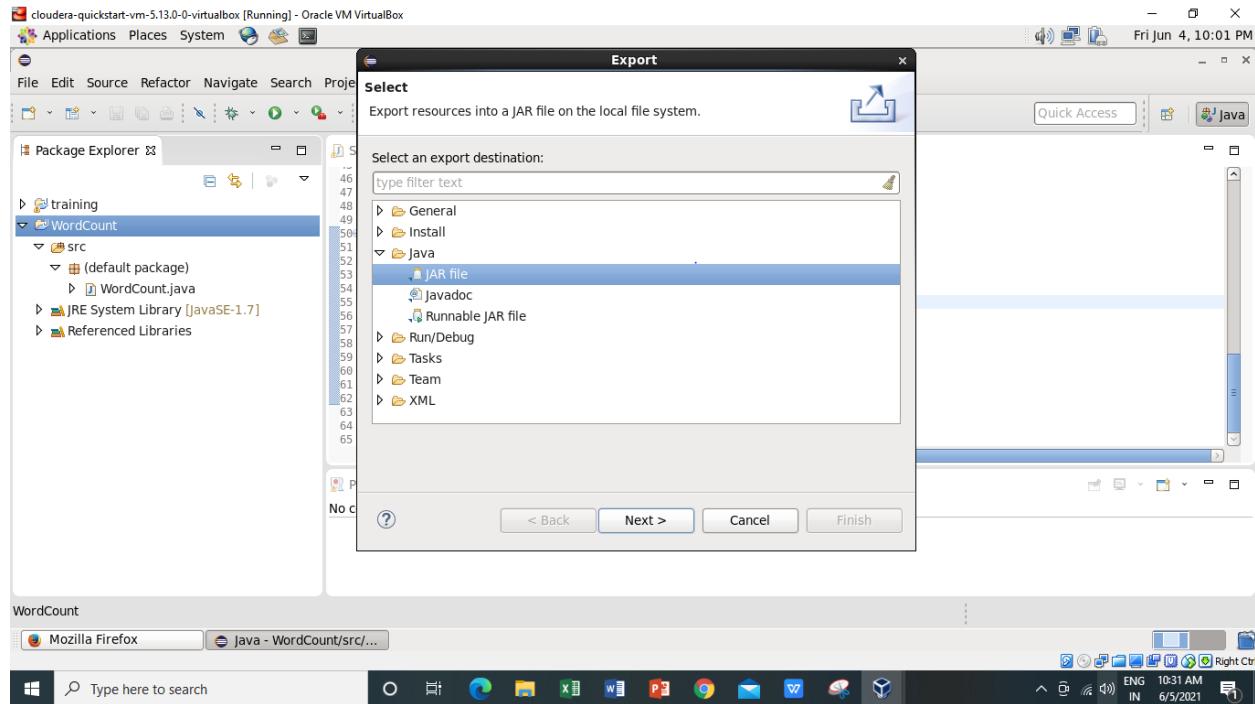


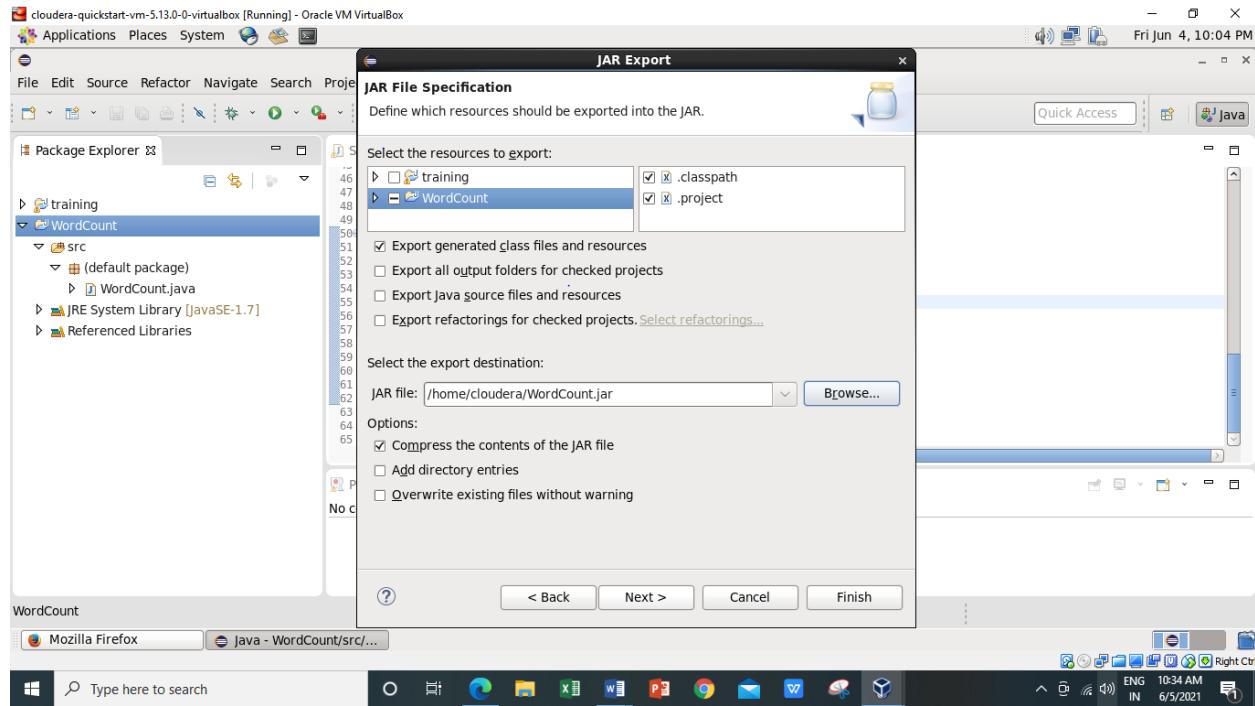
We have successfully deleted the WordCount.jar file



- 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

```
[cloudera@quickstart ~]$ ls
cloudera-manager Desktop Downloads enterprise-deployment.json kerberos Music Pictures Templates WordCount.jar
cm_api.py Documents eclipse express-deployment.json lib parcels Public Videos workspace
[cloudera@quickstart ~]$
```

- 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’.

```

cloudera@quickstart:~$ hadoop jar /home/cloudera/WordCount.jar WordCount /inputdir/abc /op1
21/06/05 06:31:13 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
21/06/05 06:31:22 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
21/06/05 06:31:28 INFO input.FileInputFormat: Total input paths to process : 1
21/06/05 06:31:29 INFO mapreduce.JobSubmitter: number of splits:1
21/06/05 06:31:32 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1622818056286_0002
21/06/05 06:31:44 INFO impl.YarnClientImpl: Submitted application application_1622818056286_0002
21/06/05 06:31:50 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1622818056286_0002/
21/06/05 06:31:50 INFO mapreduce.Job: Job: Running job: job_1622818056286_0002
21/06/05 06:33:29 INFO mapreduce.Job: Job job_1622818056286_0002 running in uber mode : false
21/06/05 06:33:29 INFO mapreduce.Job: map 0% reduce 0%
21/06/05 06:34:12 INFO mapreduce.Job: map 100% reduce 0%
21/06/05 06:34:49 INFO mapreduce.Job: map 100% reduce 100%
21/06/05 06:35:38 INFO mapreduce.Job: Job job_1622818056286_0002 completed successfully
21/06/05 06:35:41 INFO mapreduce.Job: Counters: 49
File System Counters
FILE: Number of bytes read=942
FILE: Number of bytes written=288551
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
HDFS: Number of bytes read=922
HDFS: Number of bytes written=603
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)=37060
Total time spent by all reduces in occupied slots (ms)=34834
Total time spent by all map tasks (ms)=37060

```



```

Total megabyte-milliseconds taken by all reduce tasks=35670016
Map-Reduce Framework
Map input records=8
Map output records=177
Map output bytes=1516
Map output materialized bytes=942
Input split bytes=109
Combine input records=177
Combine output records=84
Reduce input groups=84
Reduce shuffle bytes=942
Reduce input records=84
Reduce output records=84
Spilled Records=168
Shuffled Maps =1
Failed Shuffles=0
Merged Map outputs=1
GC time elapsed (ms)=957
CPU time spent (ms)=3680
Physical memory (bytes) snapshot=324685824
Virtual memory (bytes) snapshot=3015766016
Total committed heap usage (bytes)=152965120
Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
File Input Format Counters
Bytes Read=813
File Output Format Counters
Bytes Written=603

```

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

Combine input records=177

Combine output records=84

Earlier it was coming out as “zero” while executing WordCount (without combiner).

Combine input records=0

Combine output records=0

- 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->Utilities ->Browse the file system

cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox

Applications Places System

Browsing HDFS - Mozilla Firefox

Cloudera Live : Welcome | Browsing HDFS

localhost:50070/explorer.html#/

Search

Browse Directory

/

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	cloudera	supergroup	27 B	Mon May 24 12:04:47 -0700 2021	1	128 MB	Sample_01
drwxrwxrwx	hdfs	supergroup	0 B	Mon Oct 23 09:15:43 -0700 2017	0	0 B	benchmarks
drwxr-xr-x	cloudera	supergroup	0 B	Mon May 24 13:58:34 -0700 2021	0	0 B	forcopy
drwxr-xr-x	hbase	supergroup	0 B	Fri Jun 04 07:57:07 -0700 2021	0	0 B	hbase
drwxr-xr-x	cloudera	supergroup	0 B	Sat Jun 05 00:06:22 -0700 2021	0	0 B	inputdir
drwxr-xr-x	cloudera	supergroup	0 B	Mon May 24 13:20:10 -0700 2021	0	0 B	newdir
drwxr-xr-x	cloudera	supergroup	0 B	Sat Jun 05 06:34:49 -0700 2021	0	0 B	op1
drwxr-xr-x	cloudera	supergroup	0 B	Sat Jun 05 00:37:29 -0700 2021	0	0 B	outputdir
drwxr-xr-x	cloudera	supergroup	0 B	Mon May 24 13:36:01 -0700 2021	0	0 B	ric

Java - WordCount/src... cloudera@quickstart:~ Browsing HDFS - Mozilla Firefox cloudera

Type here to search

7:40 PM Sat Jun 5, 2021

cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox

Applications Places System

Browsing HDFS - Mozilla Firefox

Cloudera Live : Welcome | Browsing HDFS

localhost:50070/explorer.html#/op1

Search

Hadoop Overview Datanodes Snapshot Startup Progress Utilities

Browse Directory

/op1

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	cloudera	supergroup	0 B	Sat Jun 05 06:34:49 -0700 2021	1	128 MB	_SUCCESS
-rw-r--r--	cloudera	supergroup	603 B	Sat Jun 05 06:34:47 -0700 2021	1	128 MB	part-r-00000

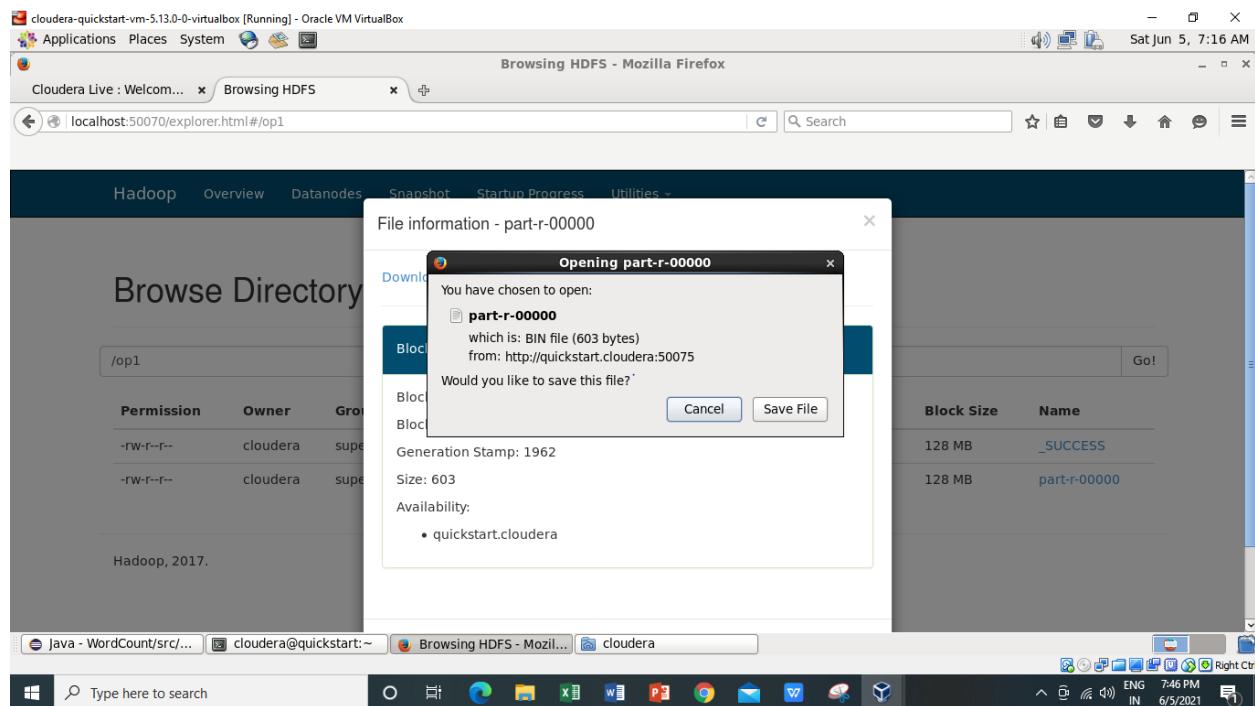
Hadoop, 2017.

Java - WordCount/src... cloudera@quickstart:~ Browsing HDFS - Mozilla Firefox cloudera

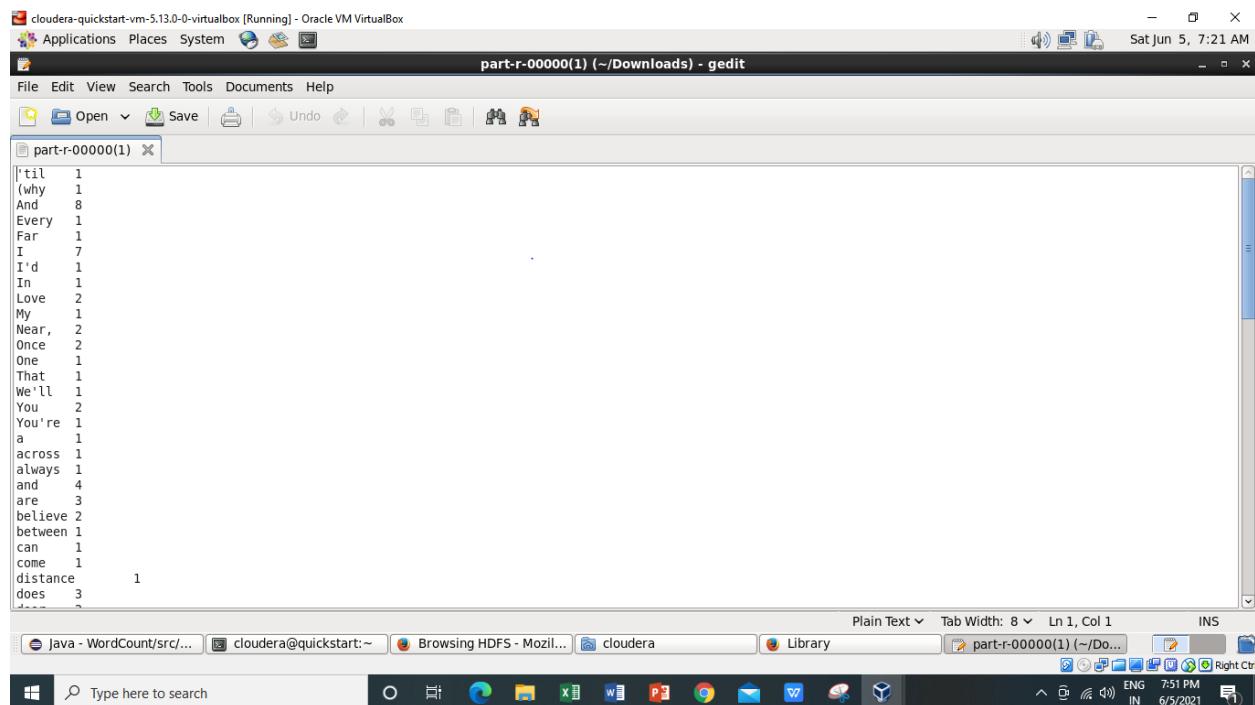
Type here to search

7:44 PM Sat Jun 5, 2021

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**



CloudBees Quickstart VM - Oracle VM VirtualBox [Running] - Oracle VM VirtualBox

File Edit View Search Tools Documents Help

part-r-00000(1) (~/Downloads) - gedit

part-r-00000(1) x

```
on      12
on?)   1
one    1
open   2
safe   1
see    1
show   1
spaces 1
stay   1
that   3
the    6
there's 1
this   1
time   2
to     2
touch  1
true   1
us     2
was    1
way    1
we'll  1
we're  1
when   1
wherever 2
will   4
you    8
you're 2
you,   1
```

Plain Text ▾ Tab Width: 8 ▾ Ln 1, Col 1 INS

Java - WordCount/src/... cloudera@quickstart:~ Browsing HDFS - Mozilla... cloudera Library part-r-00000(1) (~/Do... Right Ctrl