



CSE488: Big Data Analytics [SPRING 2023]

Assignment 2

**How to install cloudera in virtual machine and
finally run a Hadoop-MapReduce program**

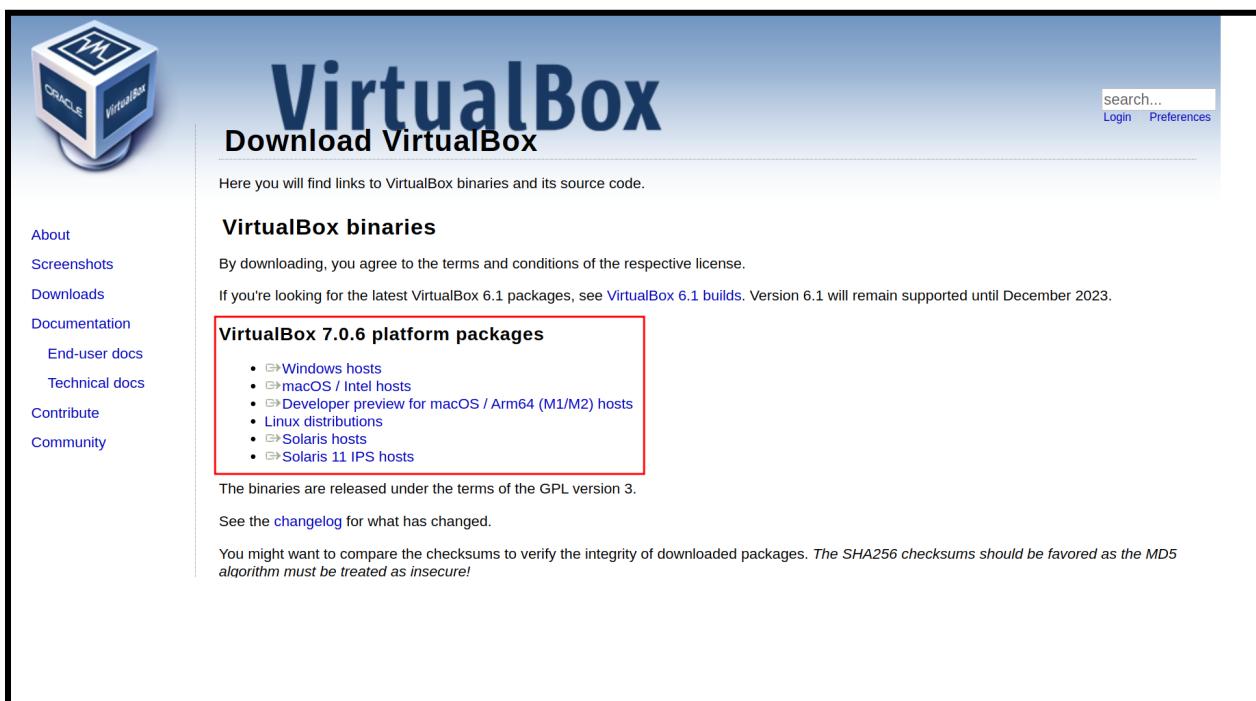
Submitted by:

Student ID: 2020-1-60-215
Student Name: Md. Abdul Ahad Rifat

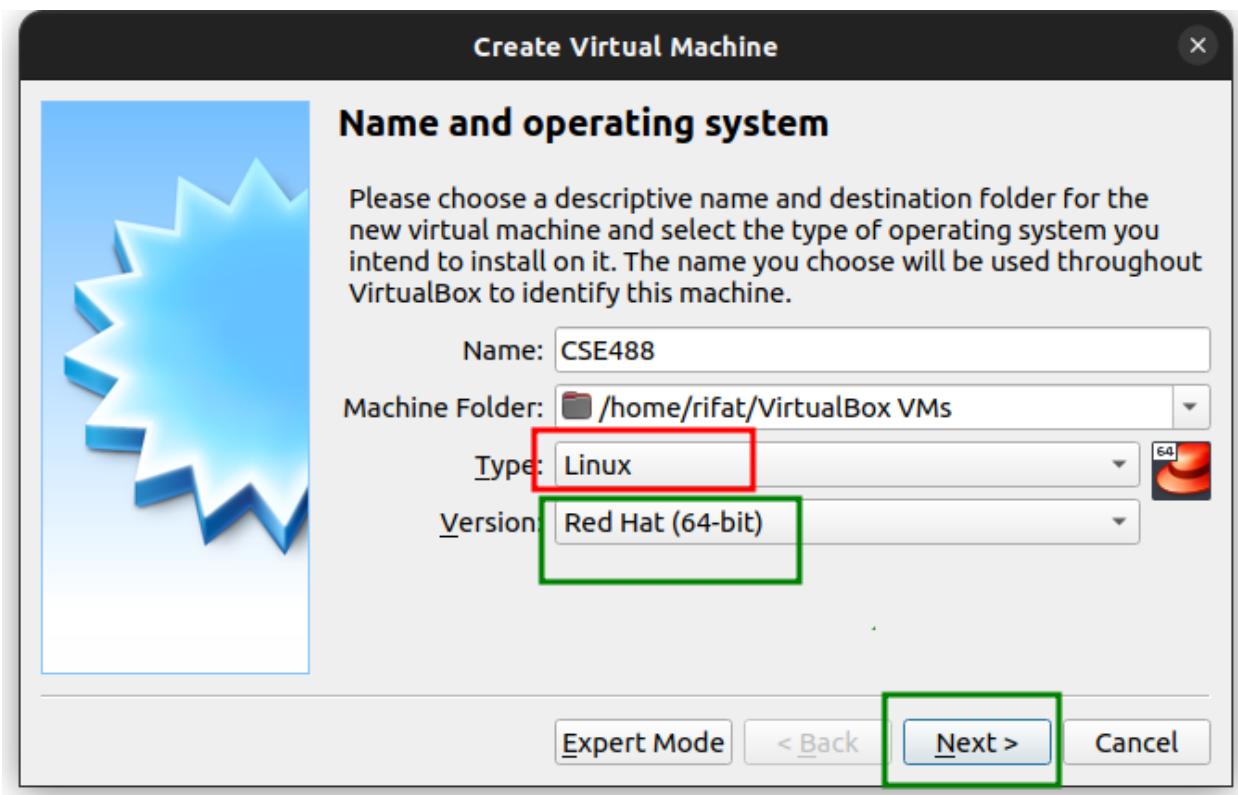
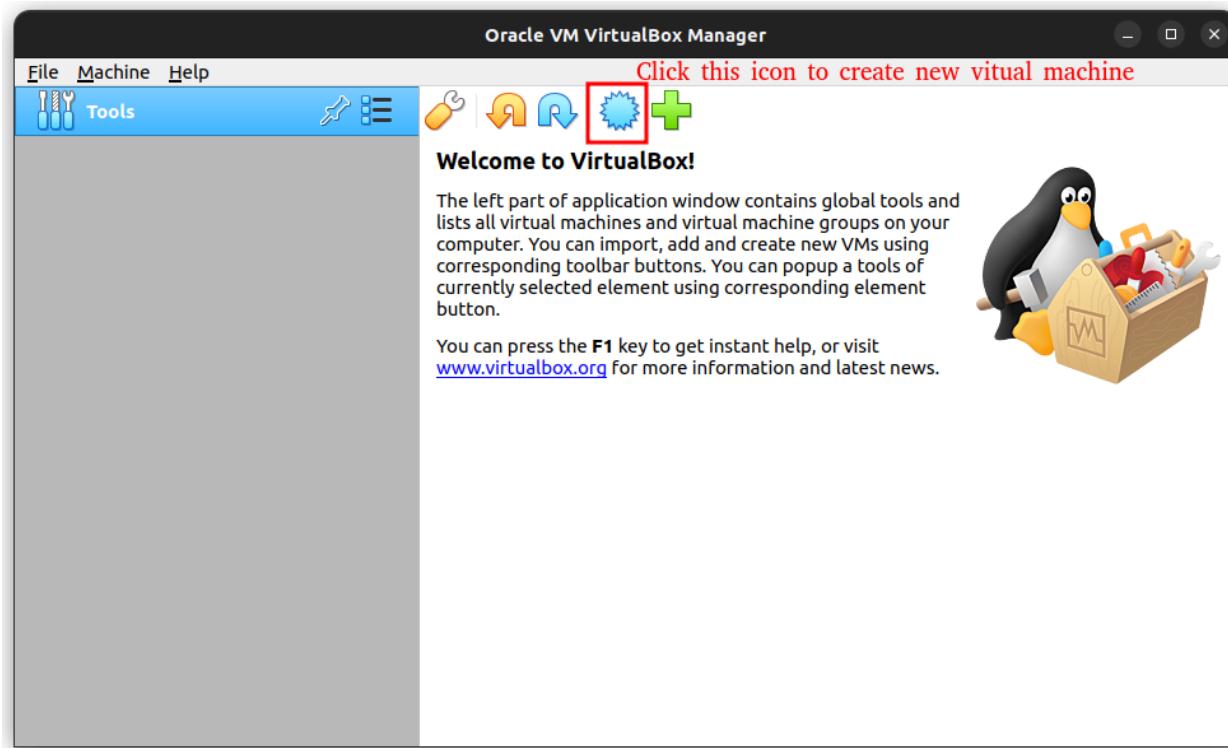
Step 1 :

Download and Installation Process of Virtual Machine and Cloudera

1. First we have to download a virtual machine just like Virtual Box on our local machine.



2. Virtual Box download link <https://www.virtualbox.org/wiki/Downloads>
3. Here we download which one our operating system recommends and after that install it.
4. Cloudera download link
https://downloads.cloudera.com/demo_vm/virtualbox/cloudera-quickstart-vm-5.13.0-0-virtualbox.zip
5. After finishing the download of cloudera we will extract the zip file so we can use it
6. Now we opening our virtual box and install our cloudera operating system as a guest operating system





CSE488 - Hard Disk Selector

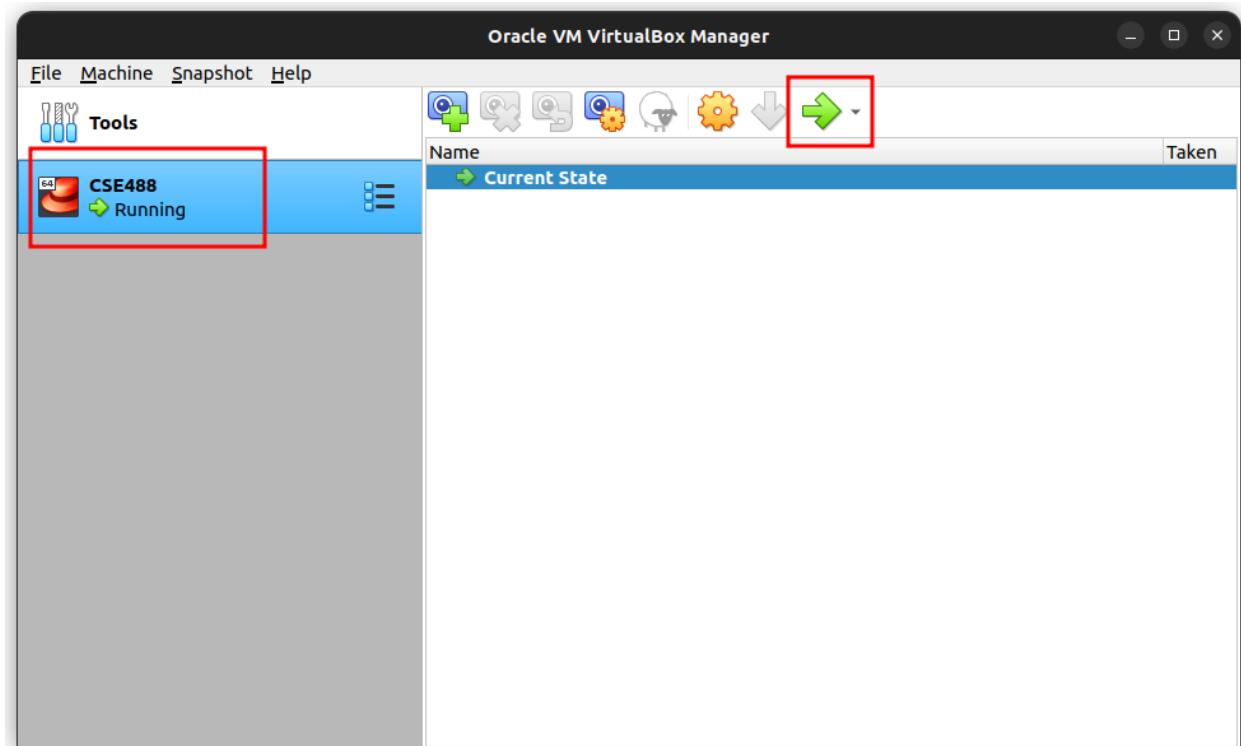
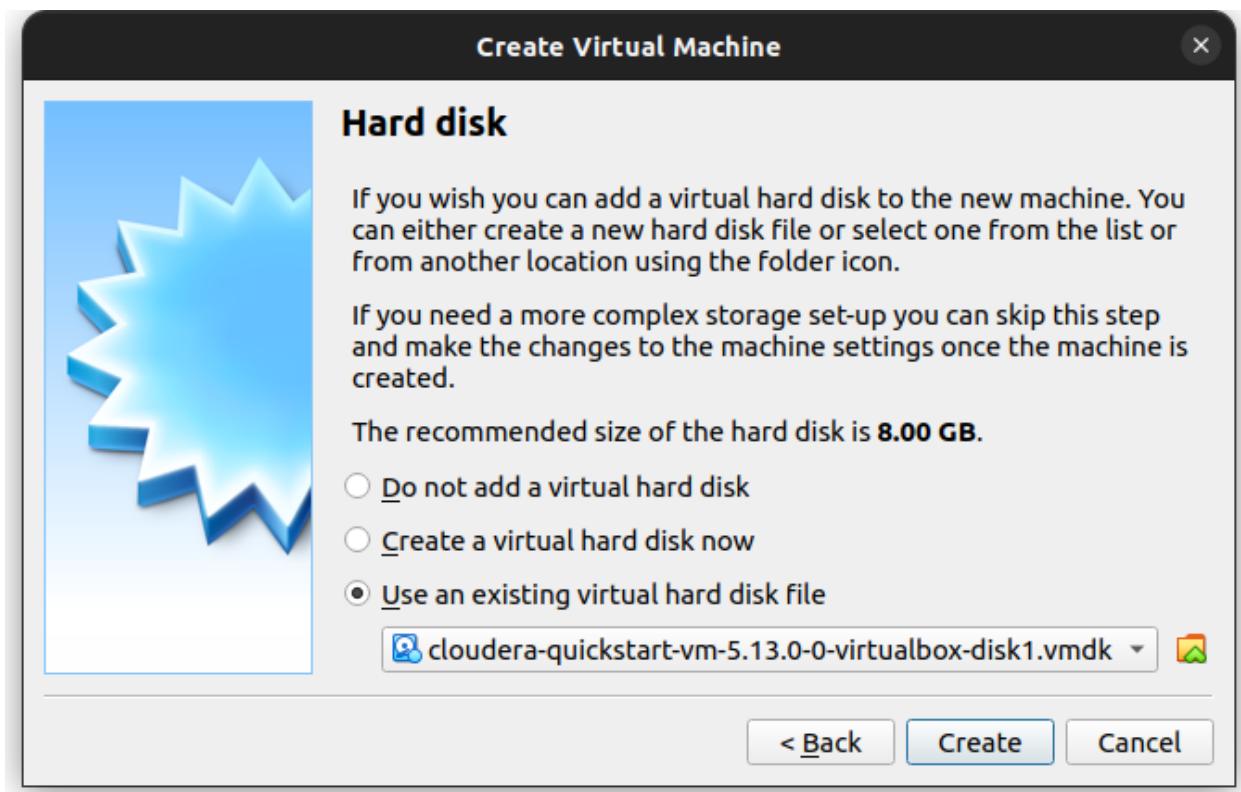
Medium

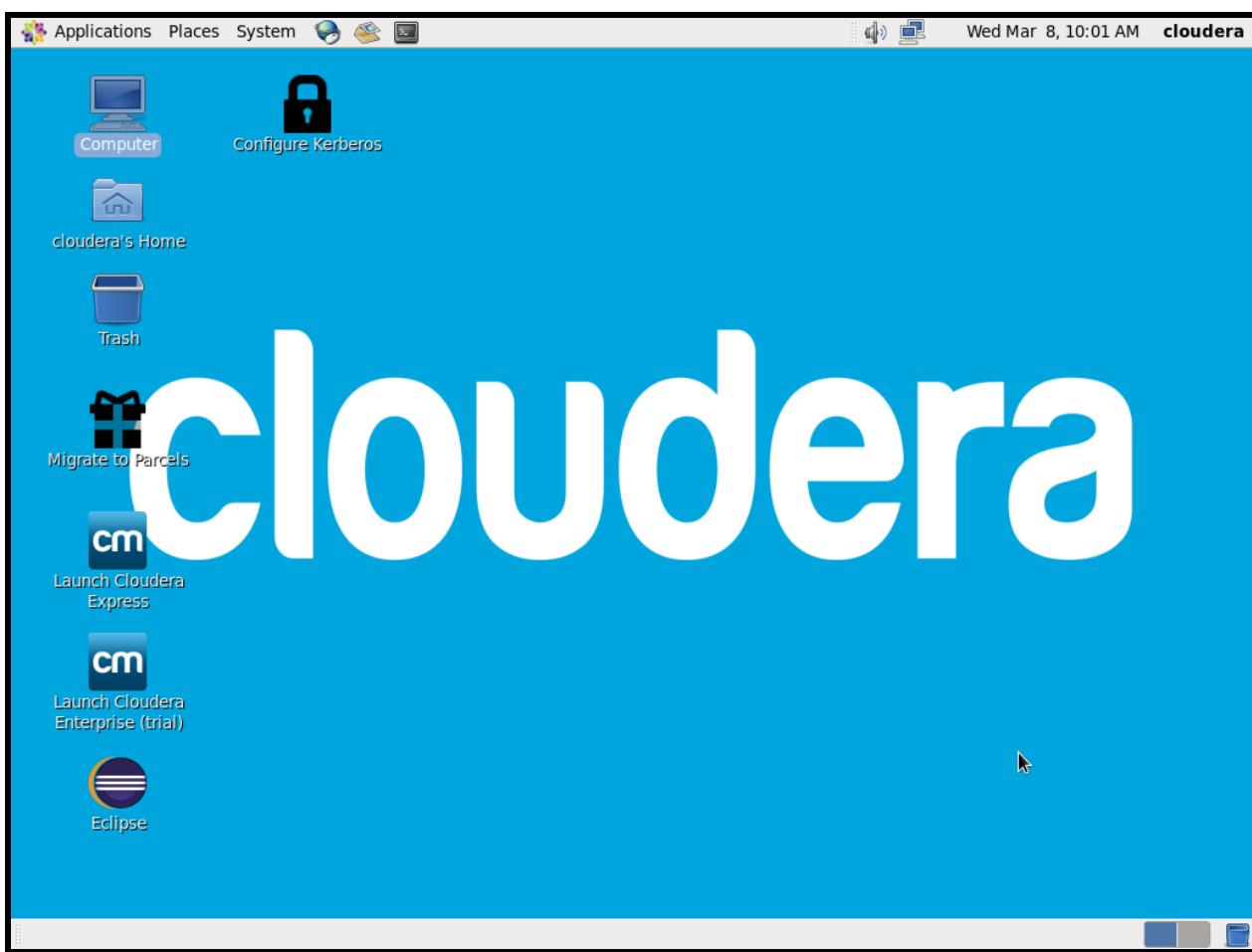
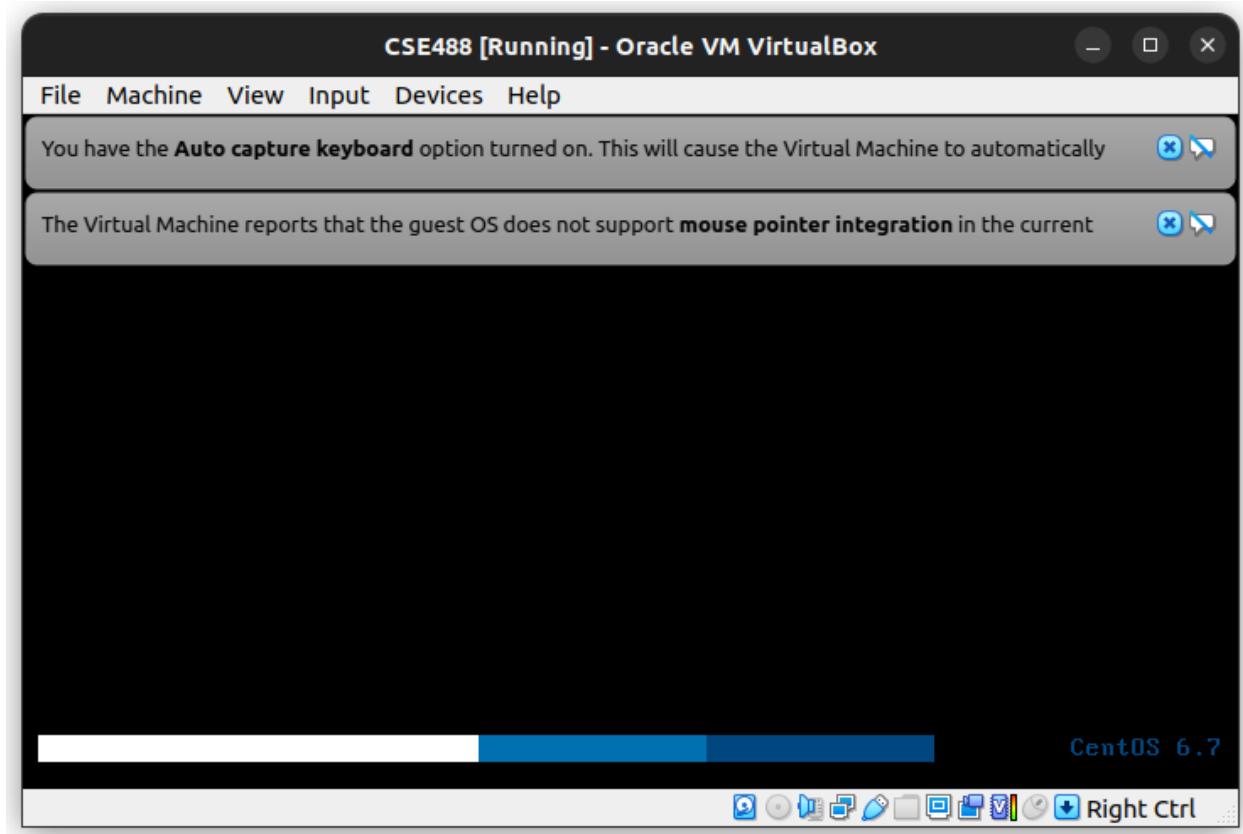
 Add  Refresh

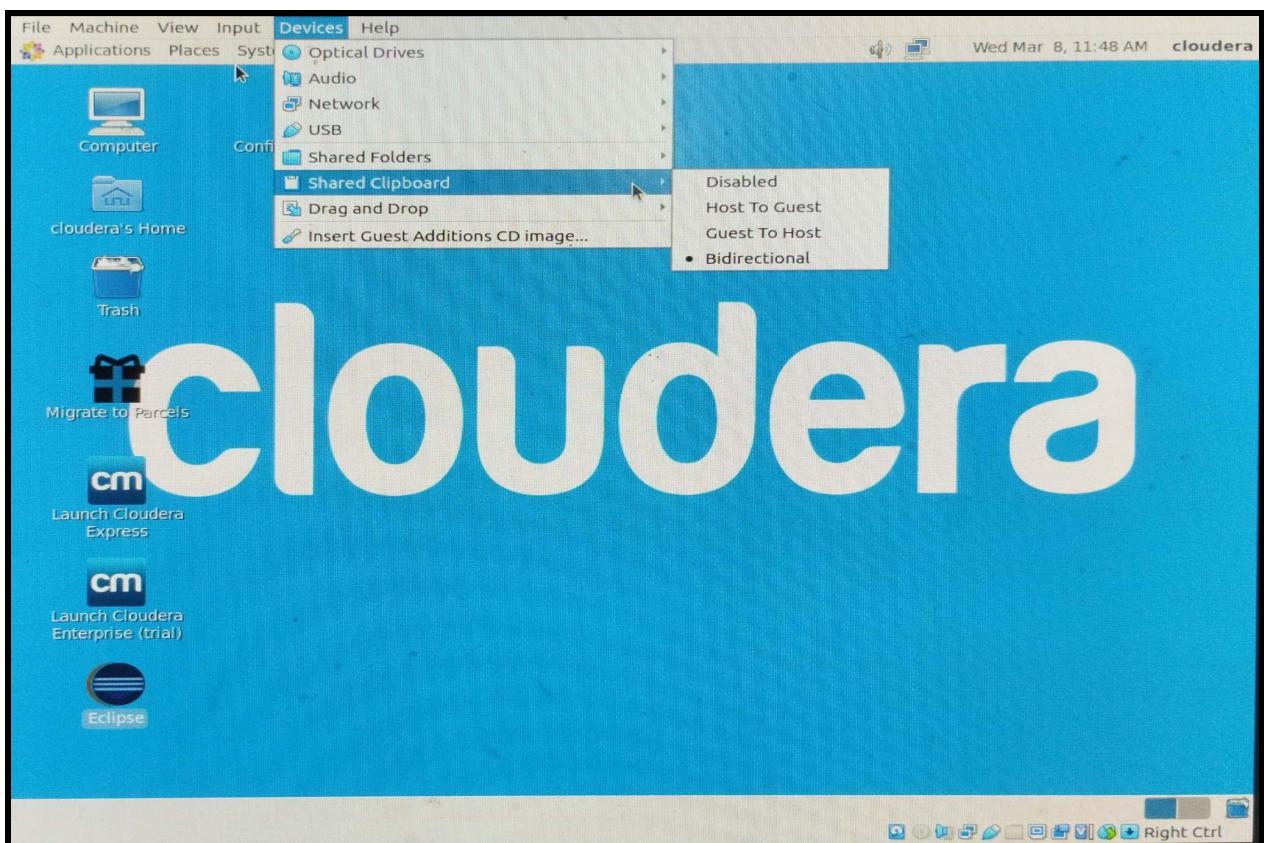
Name	Virtual Size	Actual Size
Not Attached cloudera-quickstart-vm-5.13.0-0-virtualbox-disk1.vmdk	64.00 GB	5.54 GB

Search By Name  

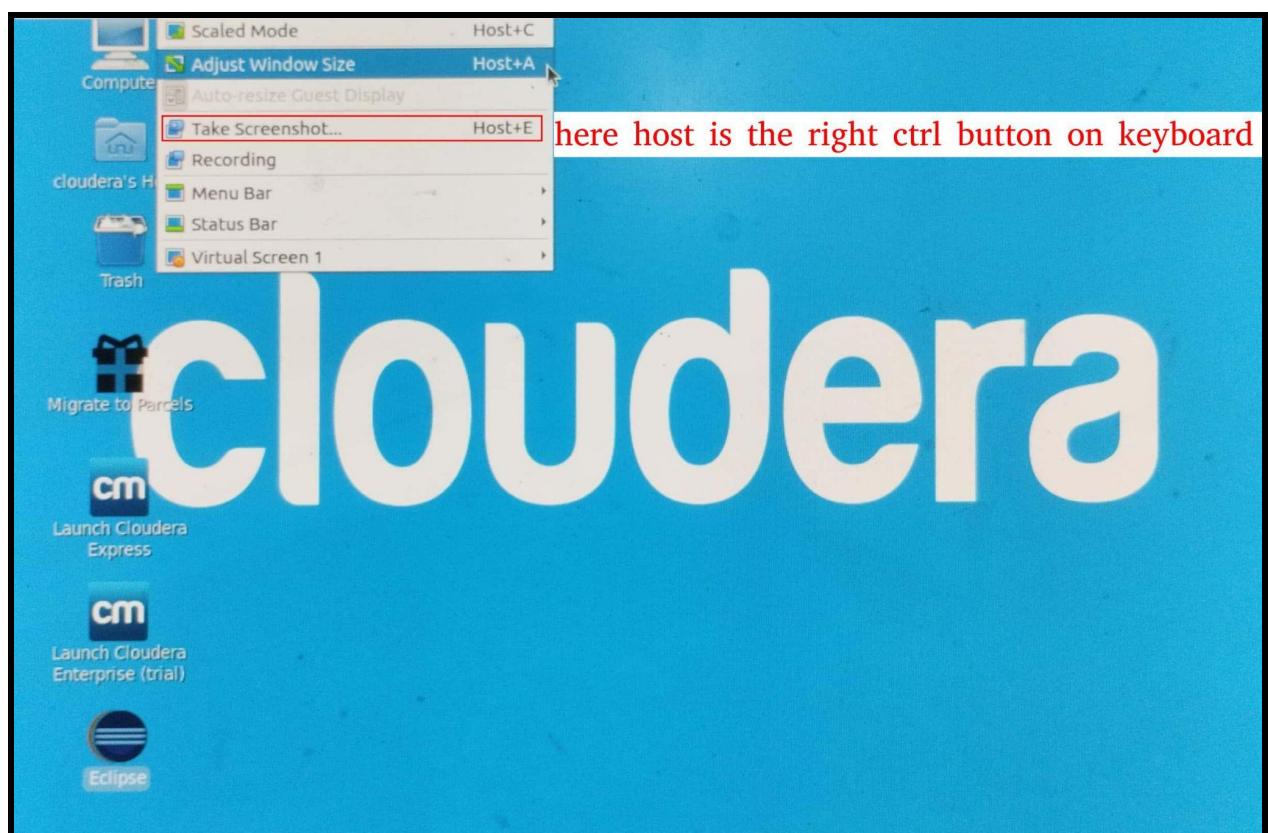
Cancel Choose







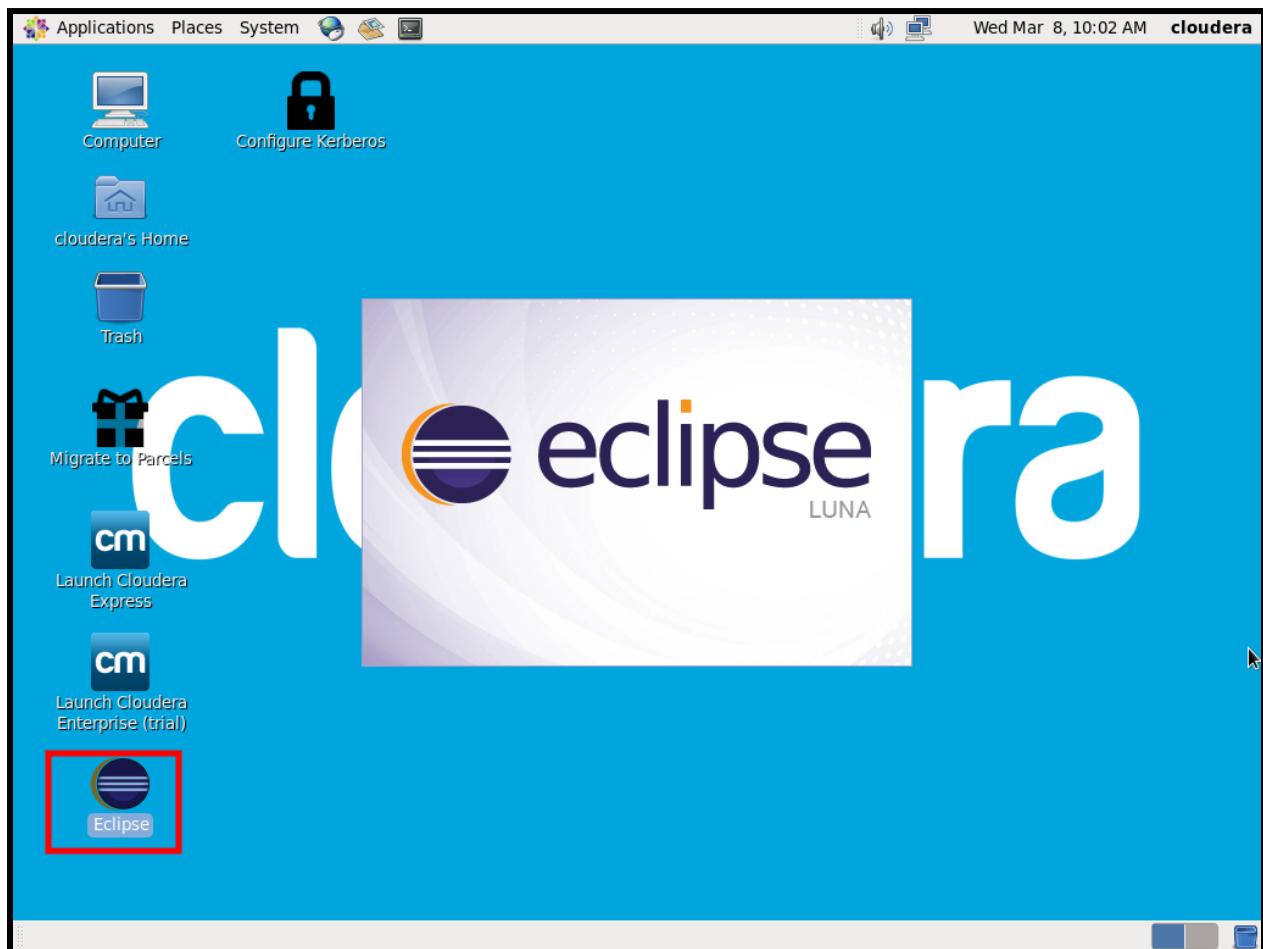
7. If we select bidirectional in a shared clipboard then we can copy and paste guest OS to host Os vice versa. Also drag and drop is the same feature.

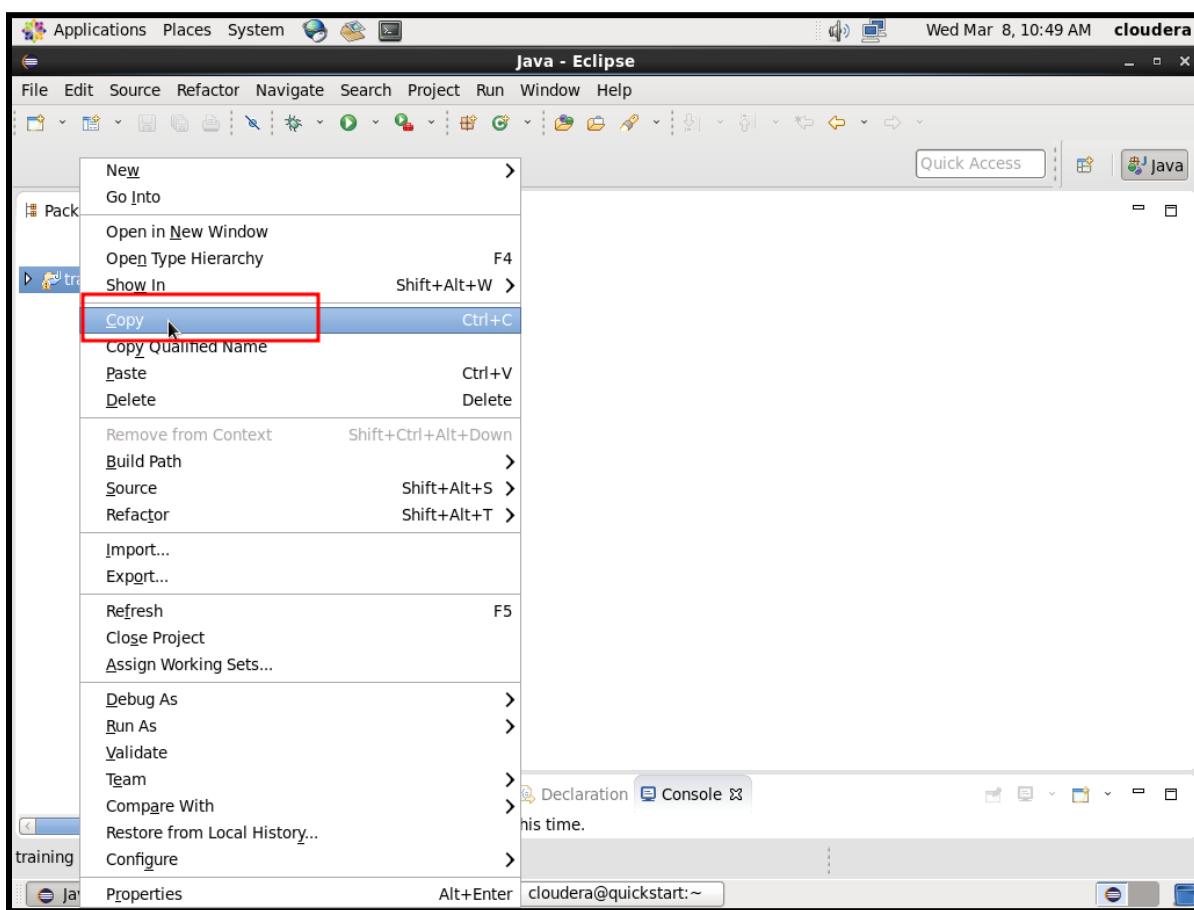
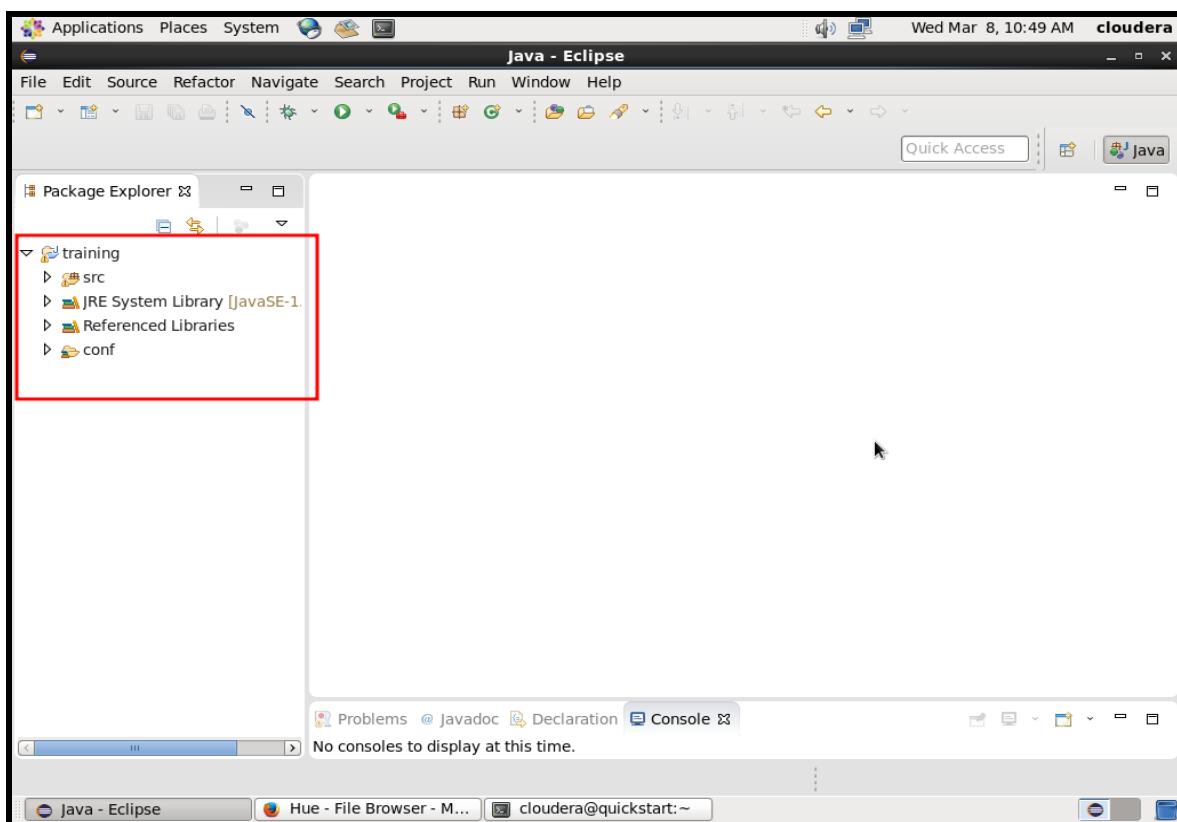


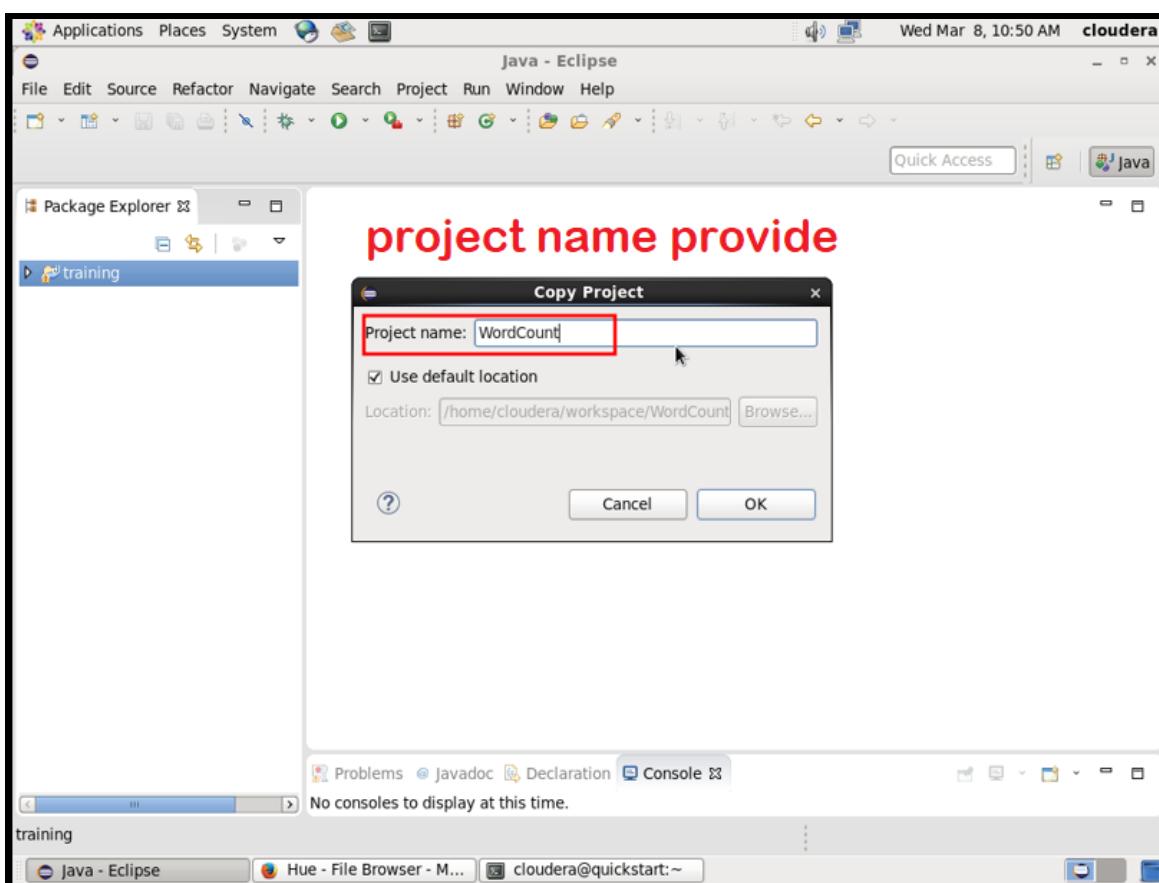
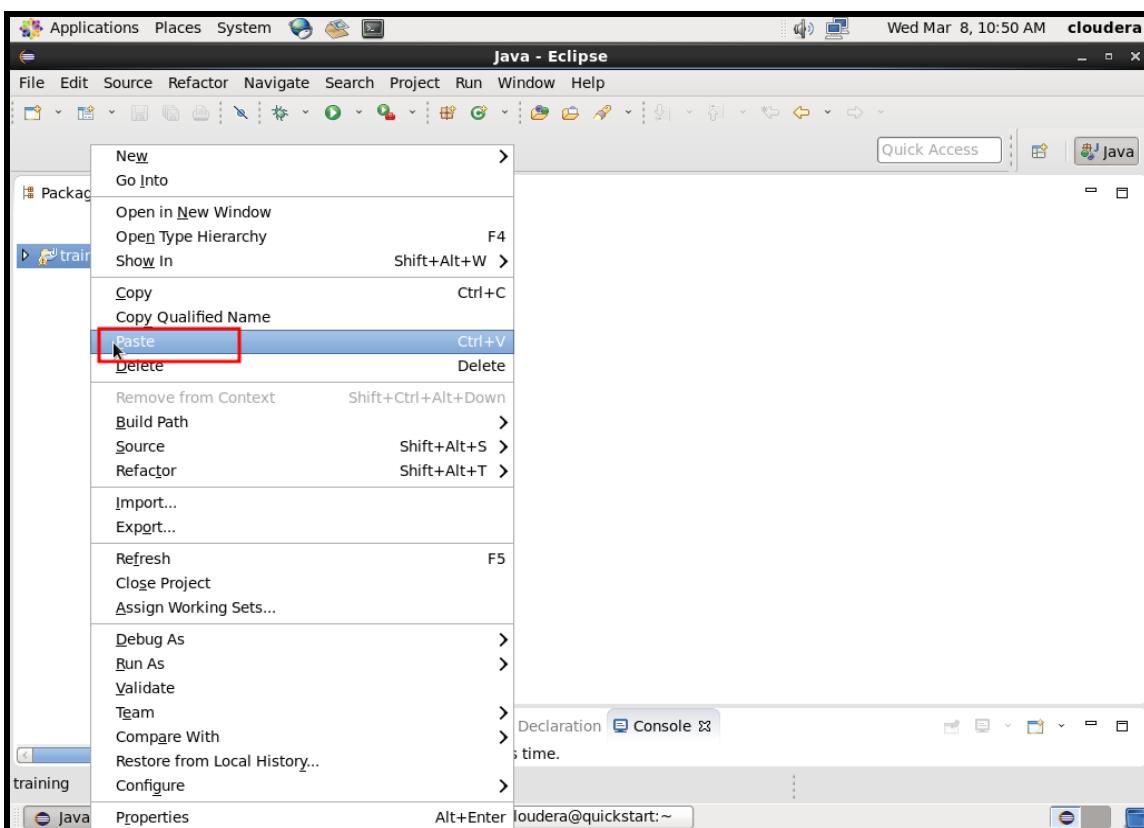
Step 2 :

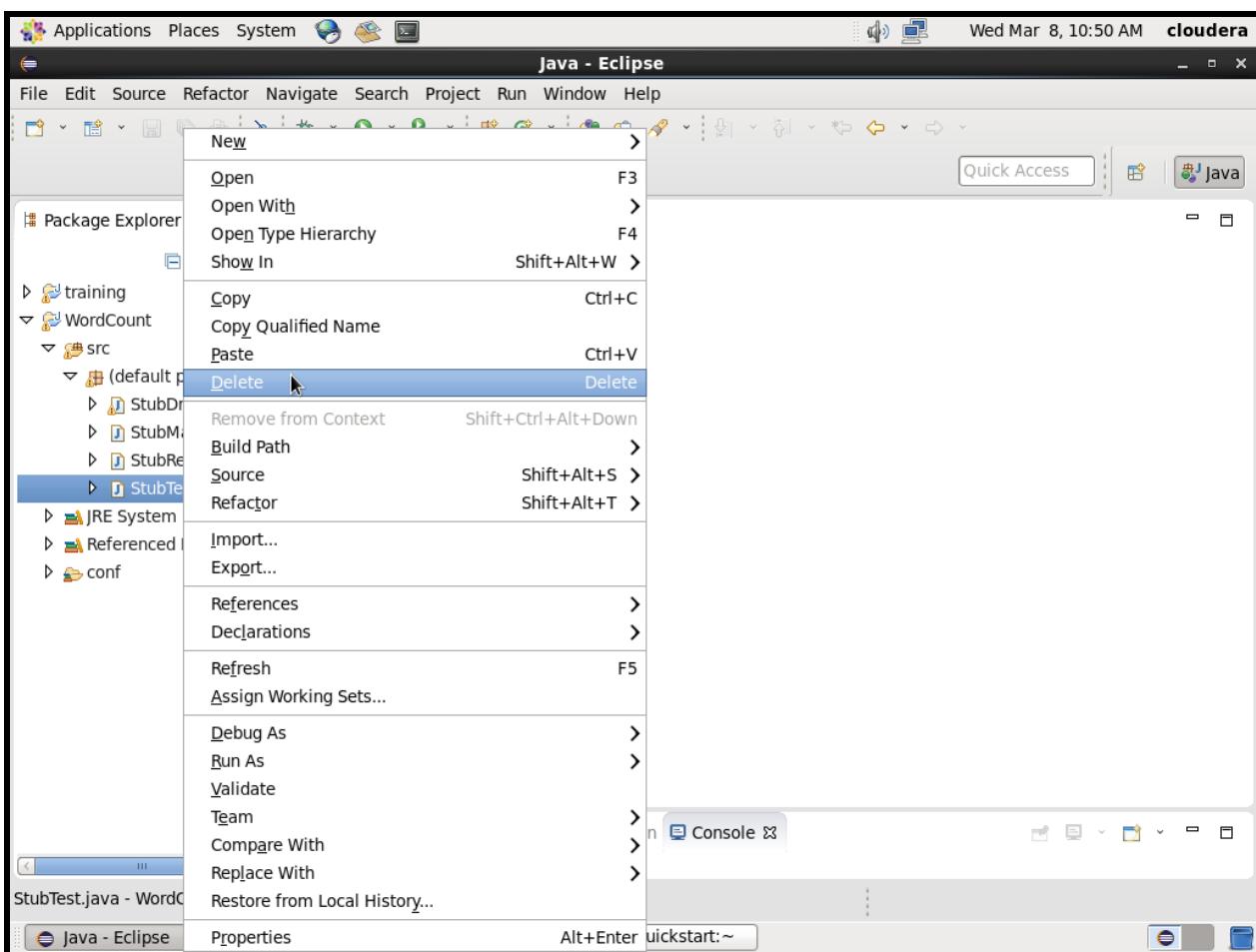
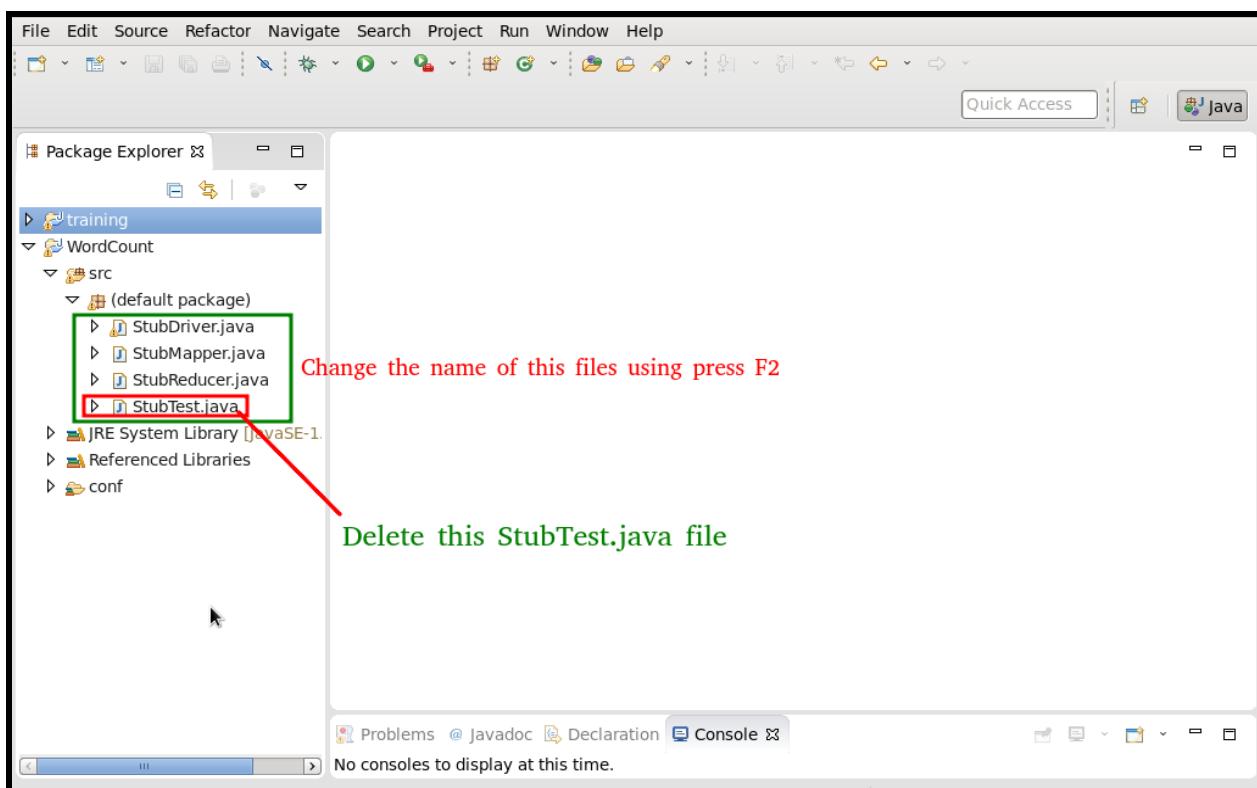
Hadoop Map-Reduce Program

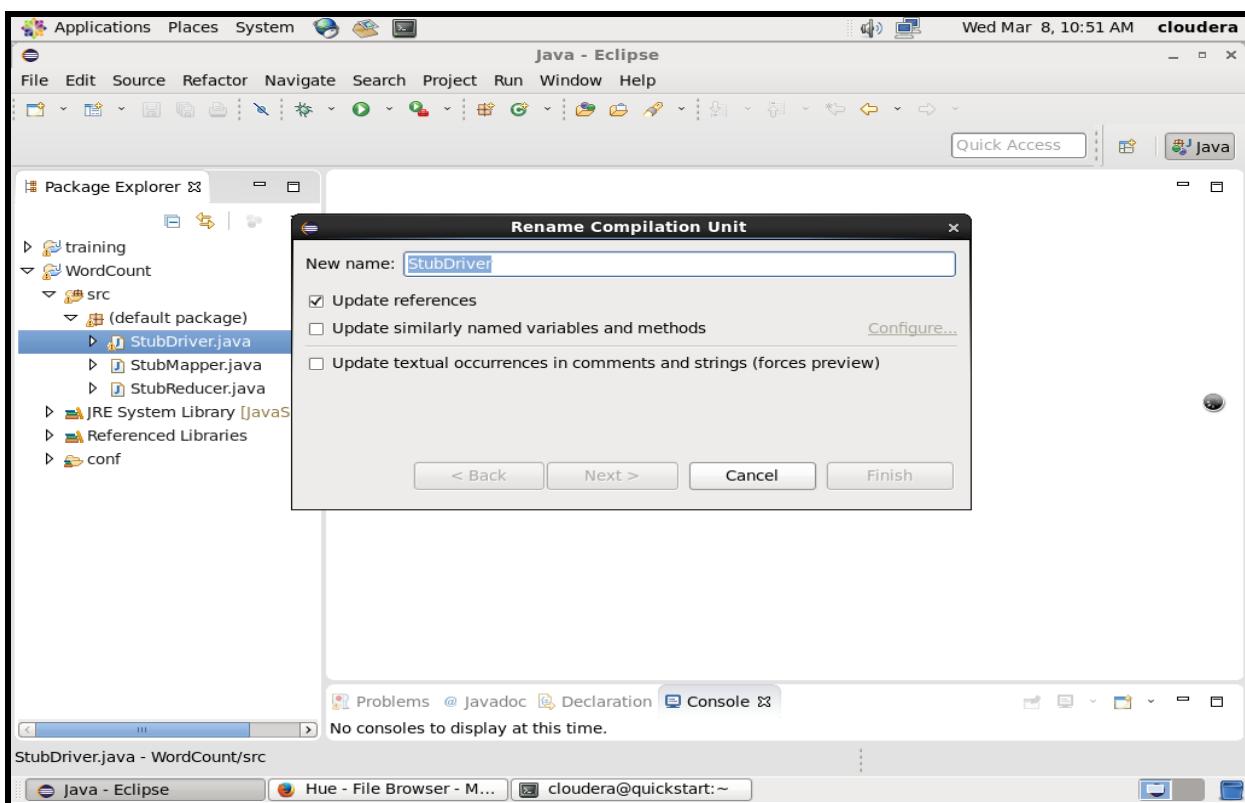
1. Opening **eclipse ide** for written java code for our Hadoop Map-reduce program and follow the instructions which are provide into the images



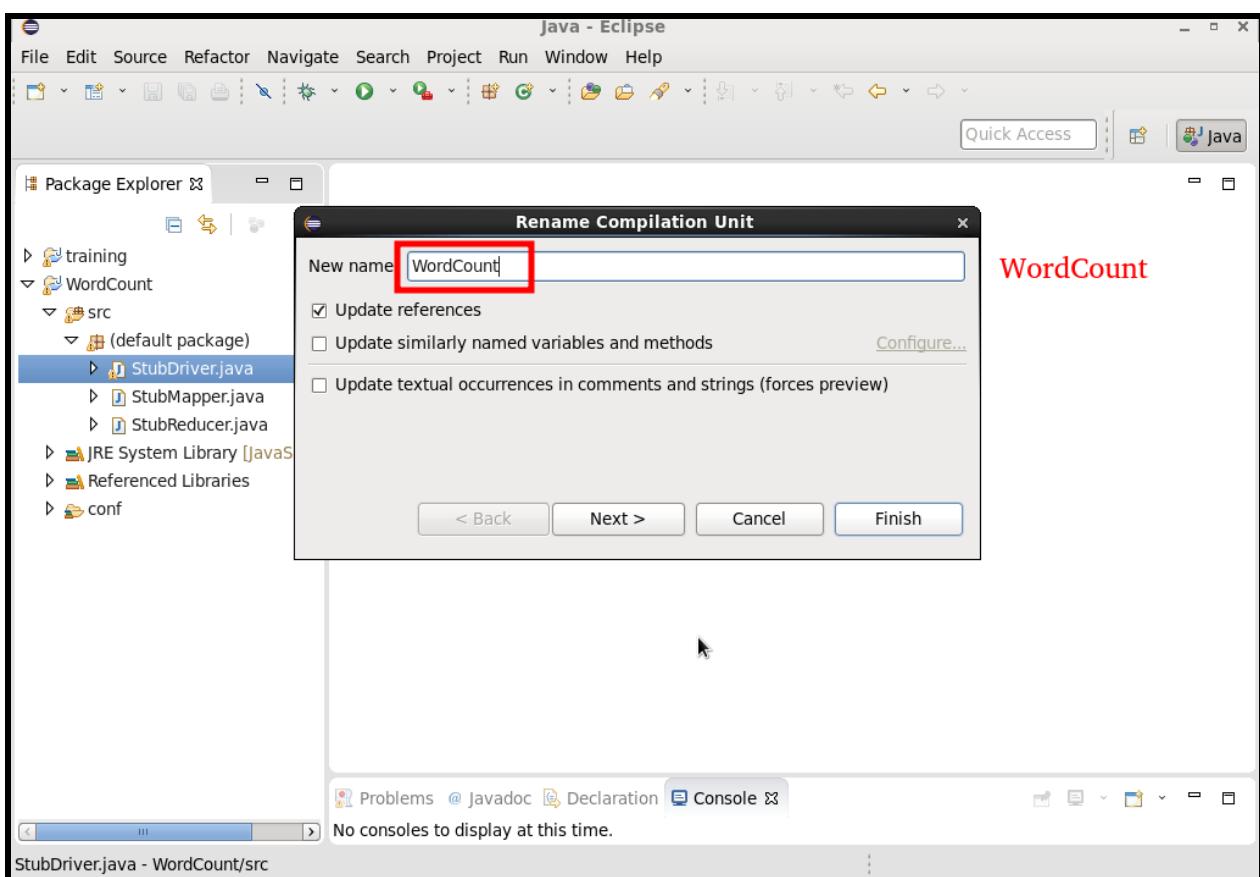


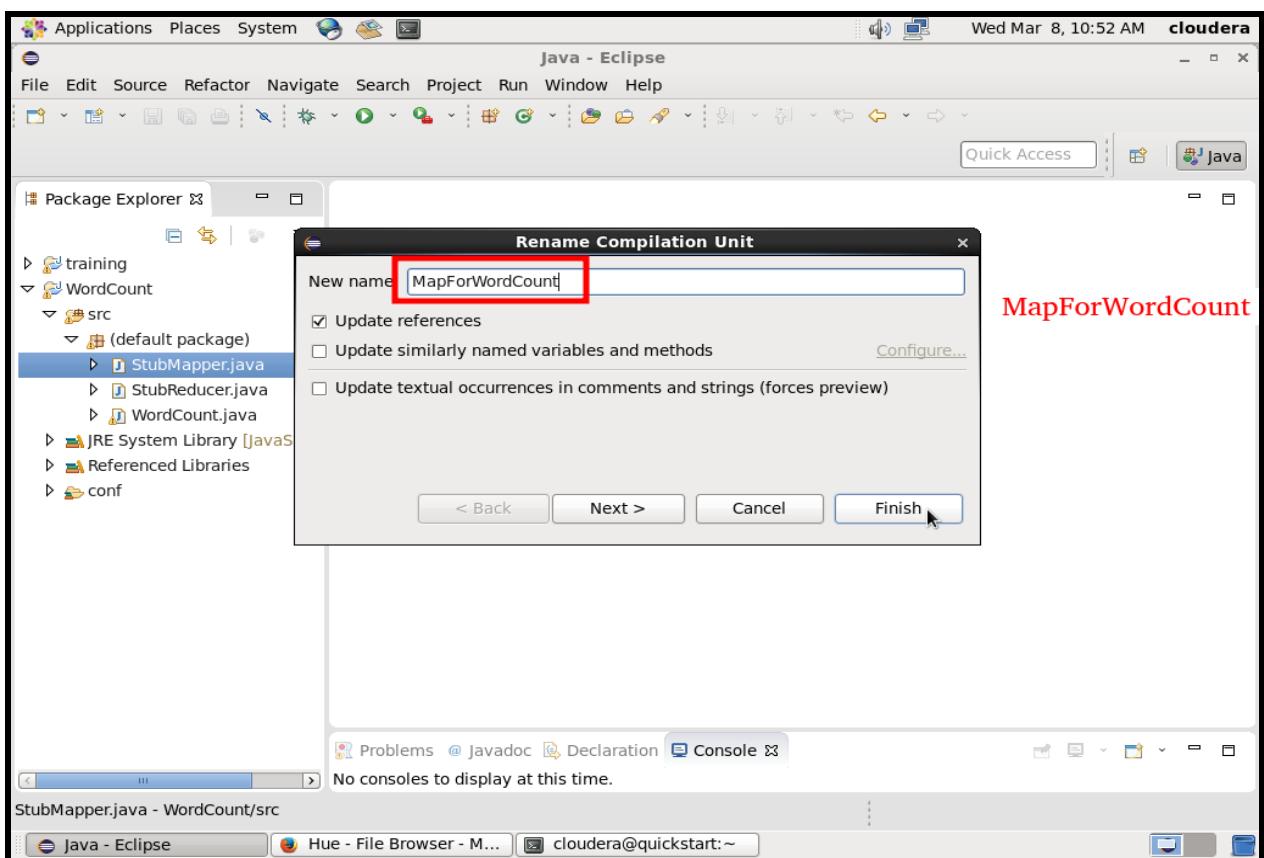




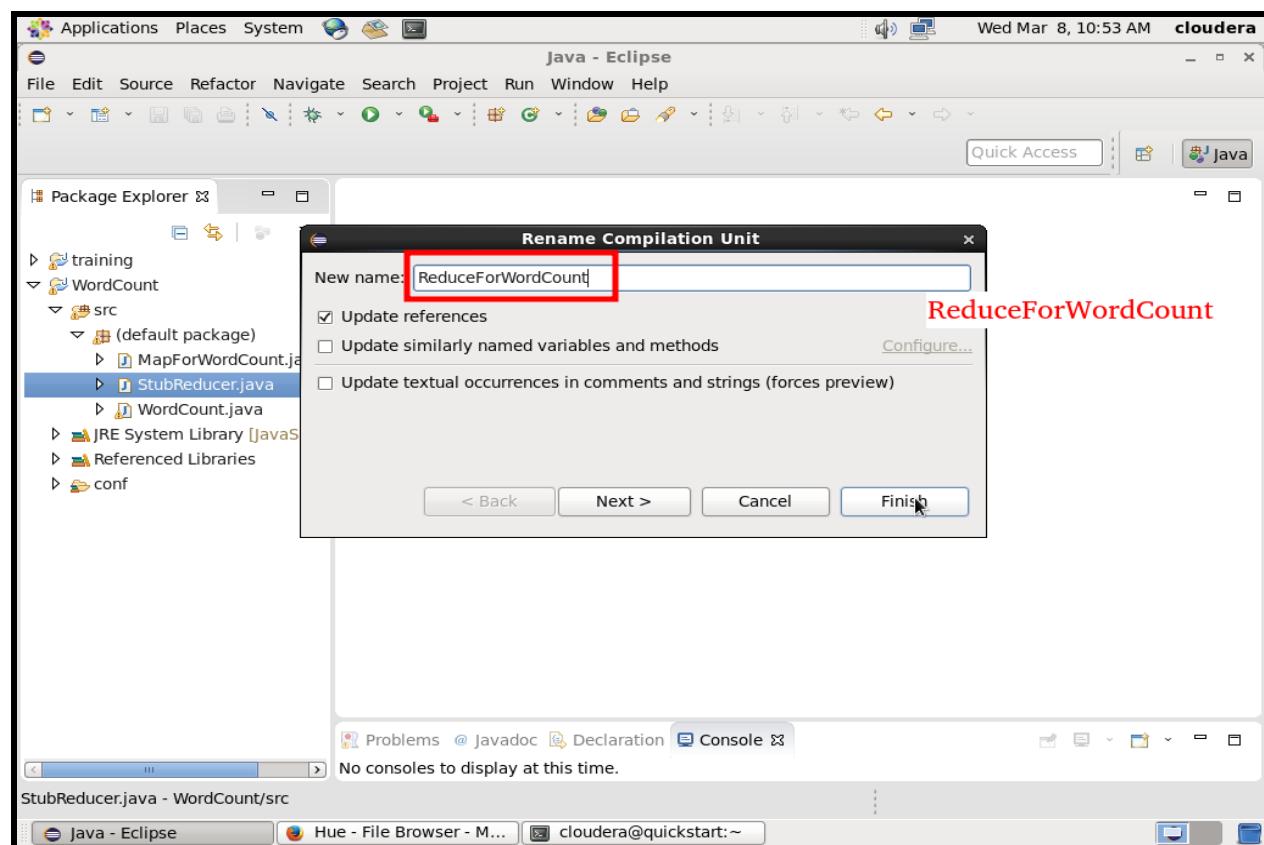


2. Rename The Files





MapForWordCount



ReduceForWordCount

3. The StubDriver.java, StubMapper.java and StubReducer.java files must be renamed WordCount.java, MapForWordCount.java and ReduceForWordCount.java respectively and the required codes must be written or copy

The screenshot shows the Eclipse IDE interface with the title bar "Java - WordCount/src/WordCount.java - Eclipse". The menu bar includes File, Edit, Source, Refactor, Navigate, Search, Project, Run, Window, Help. The toolbar has various icons for file operations. The Package Explorer view on the left shows a project structure with a "WordCount" package containing "src", "MapForWordCount.java", "ReduceForWordCount.java", and "WordCount.java". The src folder also contains "JRE System Library [javaSE-1.8]", "Referenced Libraries", and "conf". The central editor area displays the code for WordCount.java:

```
1 import org.apache.hadoop.mapreduce.Job;
2
3 public class WordCount {
4
5     public static void main(String[] args) throws Exception {
6
7         /*
8          * Validate that two arguments were passed from the command line.
9          */
10        if (args.length != 2) {
11            System.out.printf("Usage: StubDriver <input dir> <output dir>\n");
12            System.exit(-1);
13        }
14
15        /*
16         * Instantiate a Job object for your job's configuration.
17         */
18        Job job = new Job();
19
20        /*
21         * Specify the jar file that contains your driver, mapper, and reducer.
22         * Hadoop will transfer this jar file to nodes in your cluster running
23         * mapper and reducer tasks.
24         */
25        job.setJarByClass(WordCount.class);
26
27        /*
28         * Specify an easily-decipherable name for the job.
29         * This job name will appear in reports and logs.
30         */
31    }
32}
```

The bottom status bar shows "No consoles to display at this time.", "Writable", "Smart Insert", and "1 : 1". The bottom navigation bar includes tabs for Java - WordCount/src/..., Hue - File Browser - M..., and cloudera@quickstart:~.

WordCount.java

The screenshot shows the Eclipse IDE interface with the title bar "Java - WordCount/src/WordCount.java - Eclipse". The menu bar includes File, Edit, Source, Refactor, Navigate, Search, Project, Run, Window, Help. The toolbar has various icons for file operations. The Package Explorer view on the left shows a project structure with "training", "WordCount", and "src" folders containing "MapForWordCount.java", "ReduceForWordCount.java", and "WordCount.java". The editor view on the right displays the Java code for WordCount.java:

```
1+ import org.apache.hadoop.conf.Configuration;□
2  public class WordCount {
3      public static void main(String[] args) throws Exception {
4          /*
5             * Validate that two arguments were passed from the command line.
6             */
7          if (args.length != 2) {
8              System.out.printf("Usage: WordCount <input dir> <output dir>\n");
9              System.exit(-1);
10         Configuration config = new Configuration();
11         Path input = new Path(args[0]);
12         Path output = new Path(args[1]);
13
14         //Instantiate a Job object for your job's configuration.
15         @SuppressWarnings("deprecation")
16         Job job = new Job(config, "WordCount");
17         /*
18             * Specify the jar file that contains your driver, mapper, and reducer.
19             * Hadoop will transfer this jar file to nodes in your cluster running
20             * mapper and reducer tasks.
21         */
22         job.setJarByClass(WordCount.class);
23         job.setMapperClass(MapForWordCount.class);
24         job.setReducerClass(ReduceForWordCount.class);
25         job.setOutputKeyClass(Text.class);job.setOutputValueClass(IntWritable.class);
26         FileInputFormat.addInputPath(job, input);
27         FileOutputFormat.setOutputPath(job, output);
28
29         /*
30             * Start the MapReduce job and wait for it to finish.
31             * If it finishes successfully, return 0. If not, return 1.
32             */
33         job.waitForCompletion(true);
34     }
35
36     /*
37         * This function is called by the framework for every key-value pair.
38         */
39     @Override
40     public void map(LongWritable key, Text value, Context context)
41         throws IOException, InterruptedException {
42
43         String line = value.toString();
44         String[] words = line.split(",");
45         for (String word: words){
46             Text outputKey = new Text(word.toUpperCase().trim());IntWritable outputValue = new IntWritable(1);
47             context.write(outputKey, outputValue);
48         }
49     }
50 }
51
```

The "WordCount.java" tab is selected in the editor tabs. Below the editor are tabs for Problems, Javadoc, Declaration, and Console, all of which show "No consoles to display at this time." The status bar at the bottom shows "Java - WordCount/src/...".

MapForWordCount.java

The screenshot shows the Eclipse IDE interface with the title bar "Java - WordCount/src/MapForWordCount.java - Eclipse". The menu bar and toolbar are identical to the previous screenshot. The Package Explorer view on the left shows the same project structure. The editor view on the right displays the Java code for MapForWordCount.java:

```
1+ import java.io.IOException;
2  public class MapForWordCount extends Mapper<LongWritable, Text, Text, IntWritable> {
3      @Override
4      public void map(LongWritable key, Text value, Context context)
5          throws IOException, InterruptedException {
6
7         String line = value.toString();
8         String[] words = line.split(",");
9         for (String word: words){
10             Text outputKey = new Text(word.toUpperCase().trim());IntWritable outputValue = new IntWritable(1);
11             context.write(outputKey, outputValue);
12         }
13     }
14 }
```

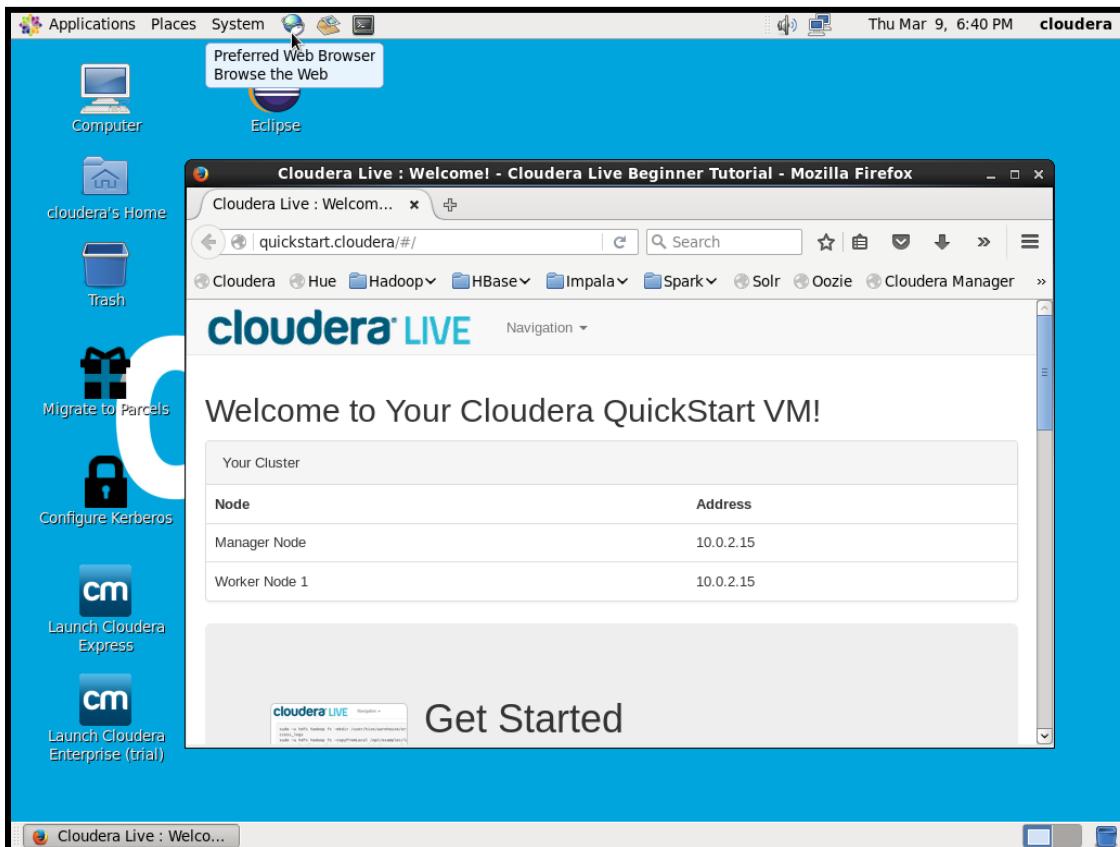
The "MapForWordCount.java" tab is selected in the editor tabs. Below the editor are tabs for Problems, Javadoc, Declaration, and Console, all of which show "No consoles to display at this time." The status bar at the bottom shows "Java - WordCount/src/...".

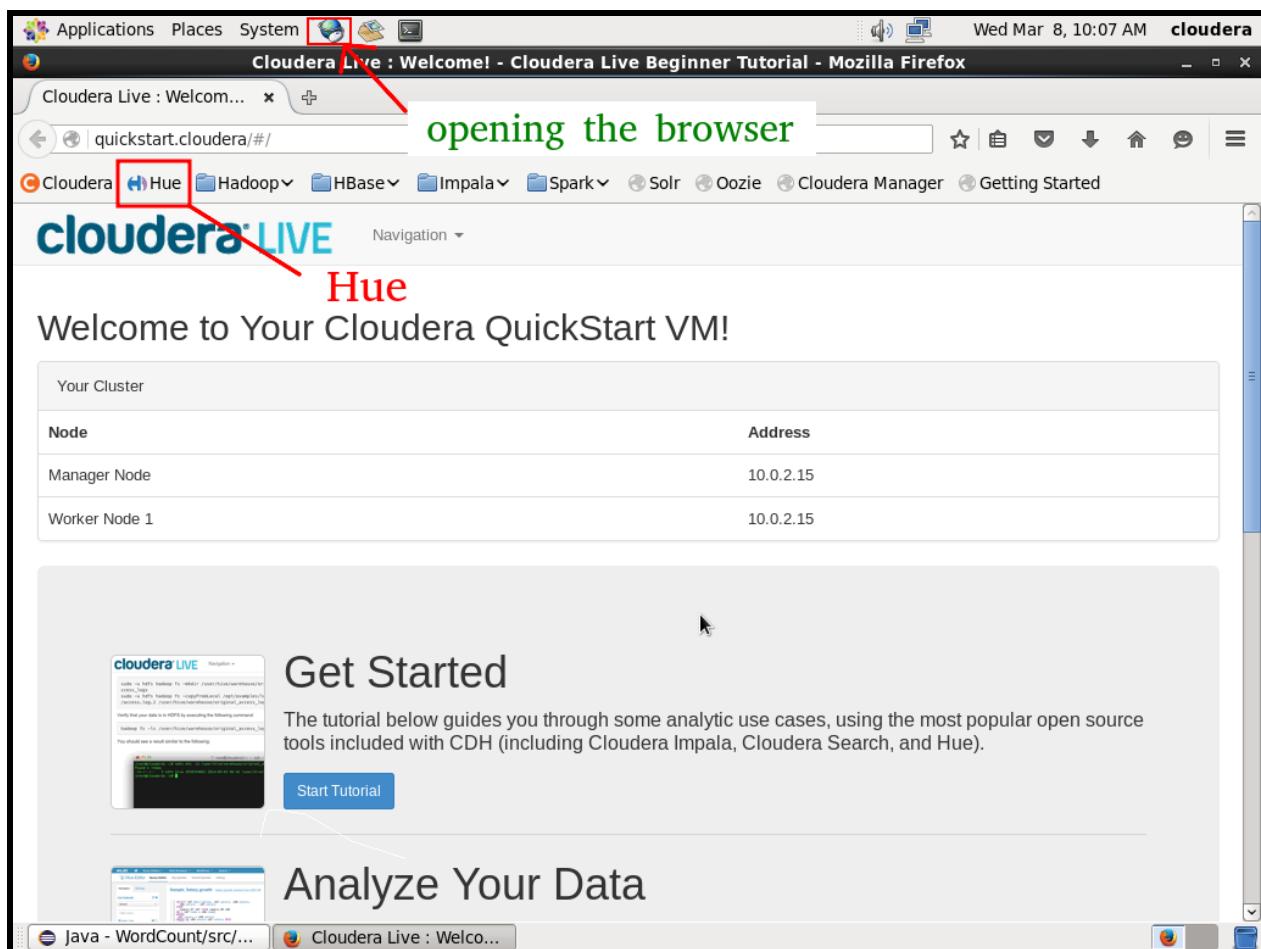
ReduceForWordCount.java

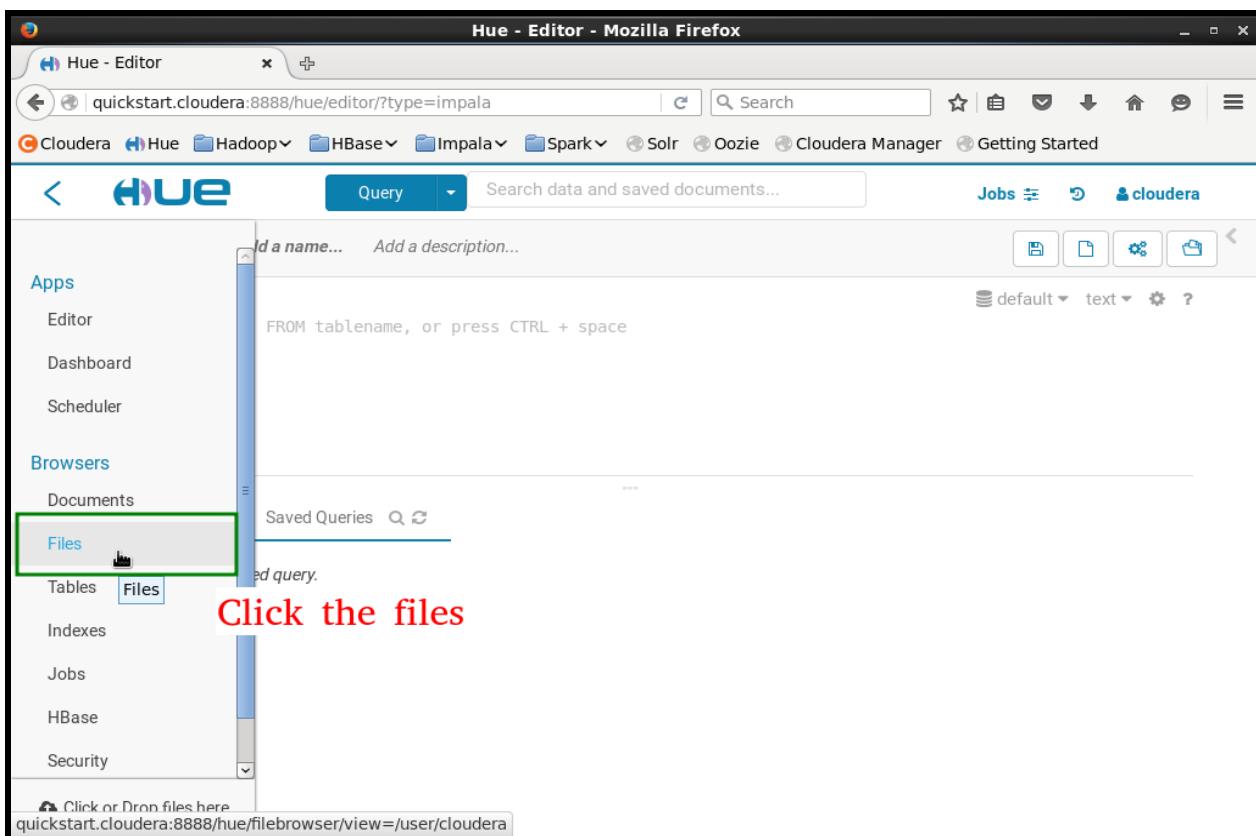
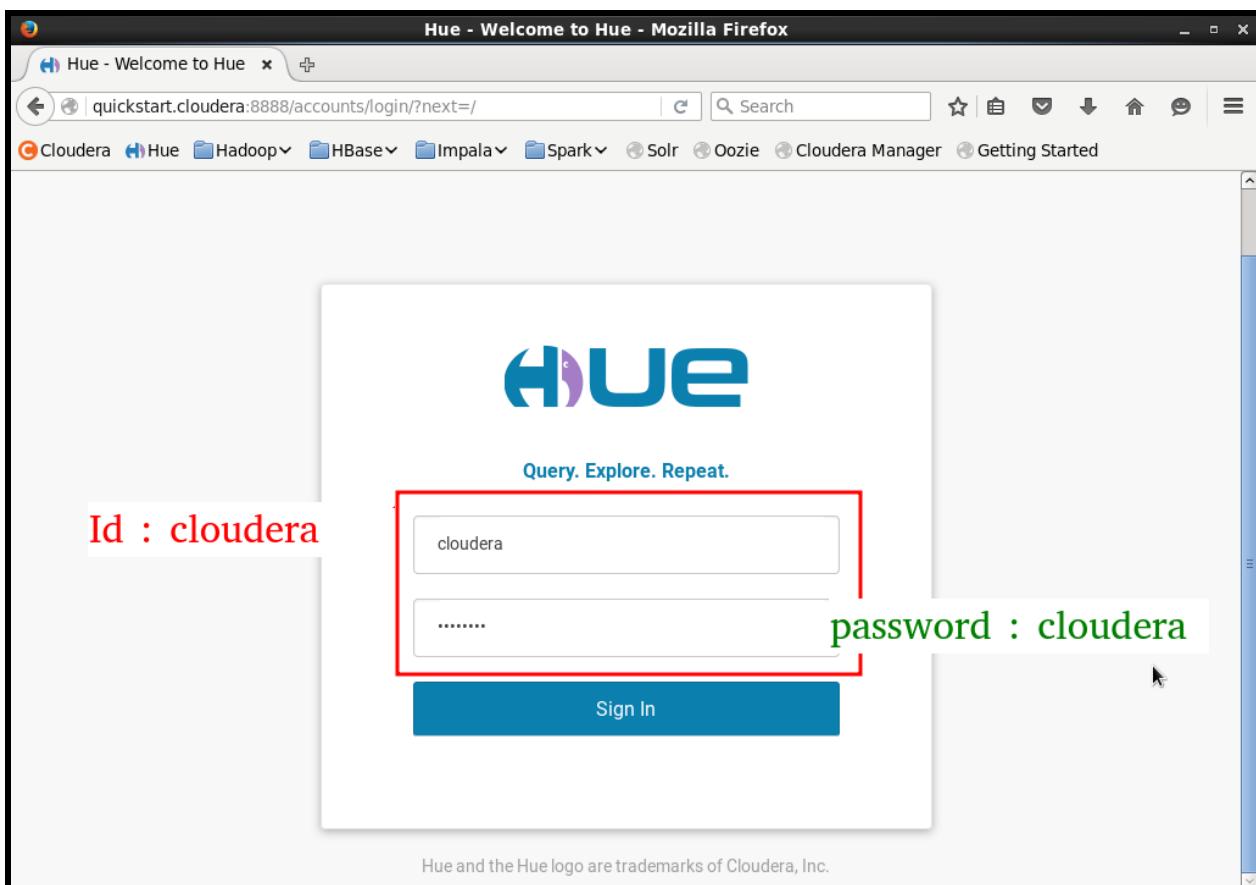
The screenshot shows the Eclipse IDE interface. The title bar reads "java - WordCount/src/ReduceForWordCount.java - Eclipse". The menu bar includes File, Edit, Source, Refactor, Navigate, Search, Project, Run, Window, Help. The toolbar has various icons for file operations. The Package Explorer view on the left shows a project structure with "WordCount" containing "src" which has "MapForWordCount.java" and "ReduceForWordCount.java" (highlighted with a blue box). Other items include "JRE System Library [JavaSE-1.8]", "Referenced Libraries", and "conf". The central editor area displays the Java code for "ReduceForWordCount.java". The code implements a Reducer that sums up values for each key. The bottom status bar shows "Writable", "Smart Insert", and the time "18 : 2".

```
import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;
public class ReduceForWordCount extends Reducer<Text, IntWritable, Text, IntWritable> {
    @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
        throws IOException, InterruptedException {
        int sum = 0;
        for(IntWritable value : values){
            sum += value.get();
        }
        context.write(key, new IntWritable(sum));
    }
}
```

4. Getting to the Cloudera distributed systems and do the necessary steps







A screenshot of the Hue File Browser interface. At the top, there's a navigation bar with links to Cloudera, Hue, Hadoop, HBase, Impala, Spark, Solr, Oozie, Cloudera Manager, and Getting Started. Below the navigation bar is the Hue logo and a search bar labeled "Search data and saved documents...". On the right side of the header, there are buttons for "Jobs", "cloudera", and user information. The main area is titled "File Browser" and shows a directory listing for "/user/cloudera". The listing includes two items: a folder named "WordCountProblem" and a file named ".". The "Actions" dropdown menu is open, and a context menu is displayed with options "Upload", "New", "File", and "Directory". The "Directory" option is highlighted with a red box. A large red box also highlights the "New" button in the main toolbar.

New Directory or Folder create

A screenshot of the Hue File Browser interface, similar to the one above but with a modal dialog box in the foreground. The dialog is titled "Create Directory" and contains a single input field labeled "Directory Name" with the value "WordCountProblem" entered. Below the input field are "Cancel" and "Create" buttons. The background shows the same directory listing for "/user/cloudera" with the same two items: "WordCountProblem" and ".". A large red box highlights the text "Directory name provide any valid name" at the bottom of the screen.

Directory name provide any valid name

A screenshot of the Hue File Browser interface. At the top, there's a navigation bar with links to Cloudera, Hue, Hadoop, HBase, Impala, Spark, Solr, Oozie, Cloudera Manager, and Getting Started. Below the navigation bar is the Hue logo and a search bar labeled "Search data and saved documents...". On the left, there's a sidebar with a "File Browser" section. The main area shows a list of files under the path "/user/cloudera/WordCountProblem". The list includes two entries: a folder named 't' and a folder named '.'. The columns in the list are Name, Size, User, Group, Permissions, and Date. A red box highlights the "New" button in the top right corner of the toolbar, which has options for "Upload" and "New". Below the list, there's a pagination control showing "Page 1 of 1" and navigation icons.

New File Create

A screenshot of the Hue File Browser interface in Mozilla Firefox. A "Create File" dialog box is open in the center. It has a "File Name" input field containing "Input.txt". Below the input field are "Cancel" and "Create" buttons. The background shows the Hue File Browser interface with a list of files under the path "/user/cloudera/WordCountProblem". The list includes a folder named 't' and a folder named '.'. The columns in the list are Name, Size, User, Group, Permissions, and Date. A red box highlights the "File Name" input field in the "Create File" dialog. Below the dialog, there's a pagination control showing "Page 1 of 1" and navigation icons.

Any Valid File Name

Hue - File Browser - Mozilla Firefox

Wed Mar 8, 10:11 AM cloudera

Hue - File Browser

quickstart.cloudera:8888/hue/filebrowser/view=/user/cloudera/WordCountProblem/Input.txt

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

HUE

File Browser

View as binary

Edit file

Download

View file location

Refresh

Last modified
03/08/2023 6:11 PM

User
cloudera

Group
cloudera

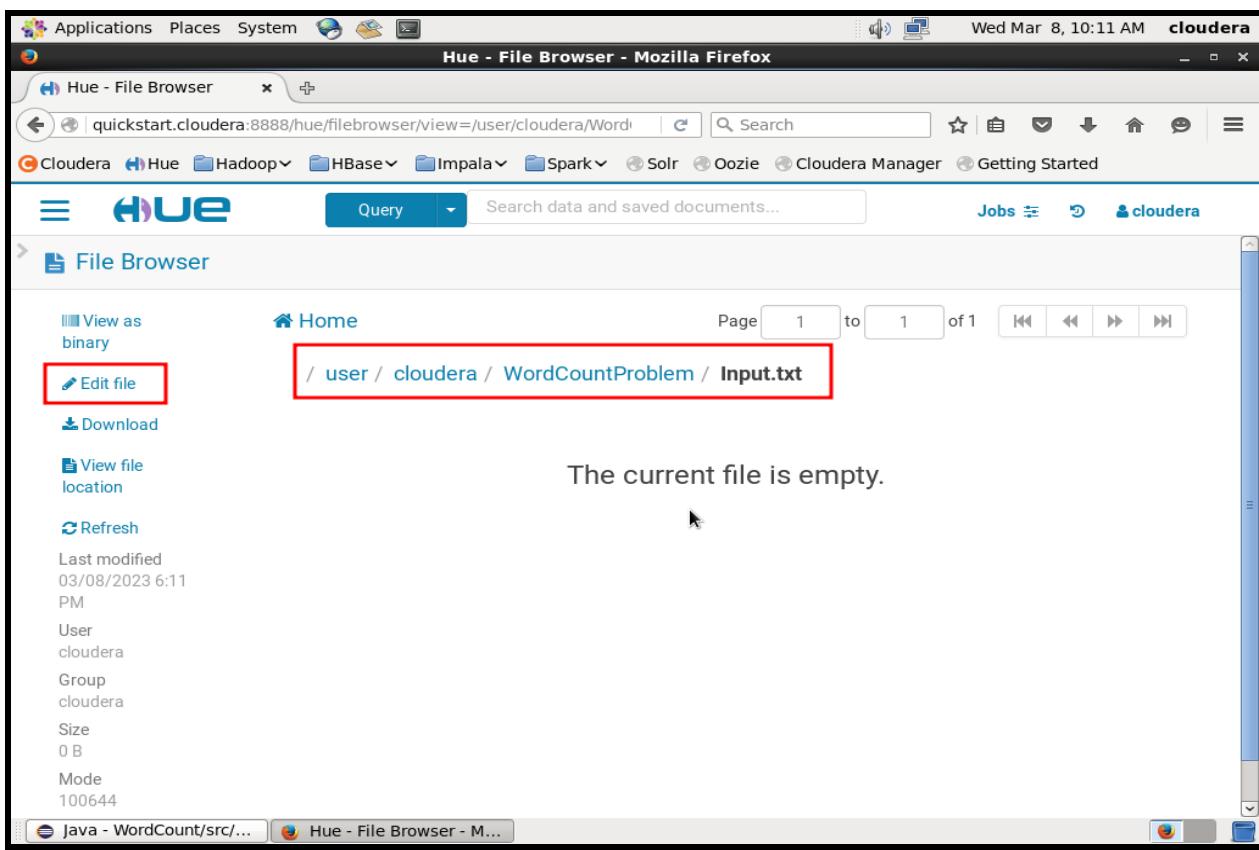
Size
0 B

Mode
100644

Home / user / cloudera / WordCountProblem / Input.txt

The current file is empty.

Java - WordCount/src/... Hue - File Browser - M...



Hue - File Browser - Mozilla Firefox

Wed Mar 8, 10:11 AM cloudera

Hue - File Browser

quickstart.cloudera:8888/hue/filebrowser/view=/user/cloudera/WordCountProblem/Input.txt

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

HUE

File Browser

View file

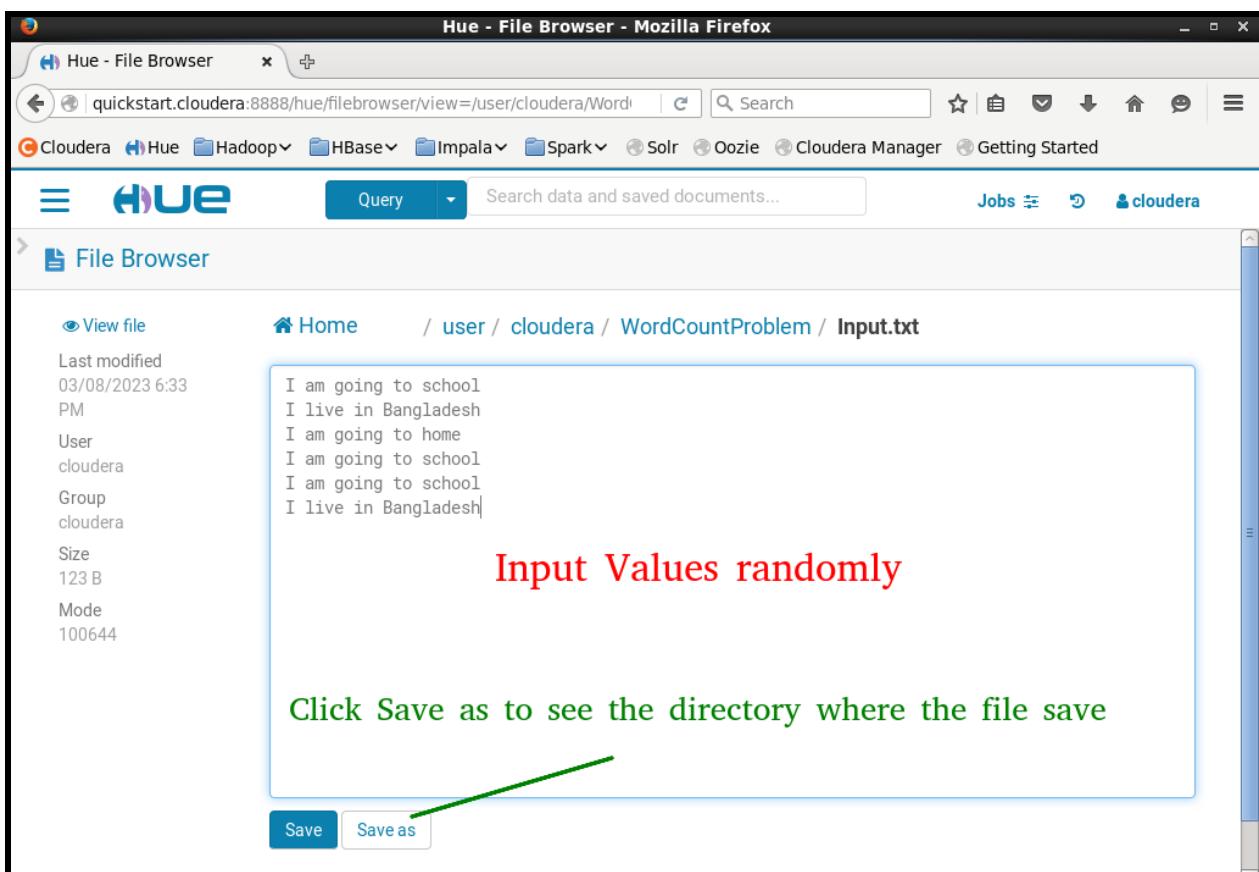
Home / user / cloudera / WordCountProblem / Input.txt

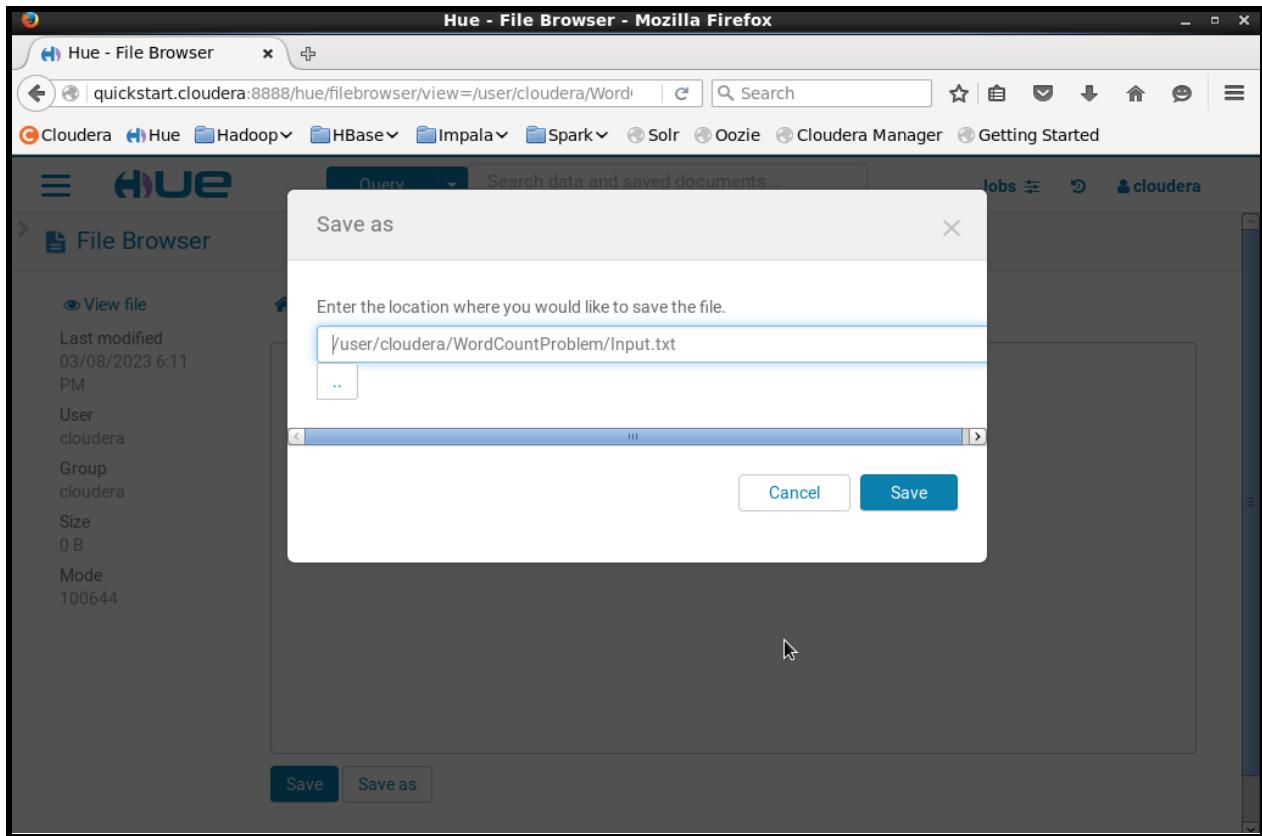
I am going to school
I live in Bangladesh
I am going to home
I am going to school
I am going to school
I live in Bangladesh

Input Values randomly

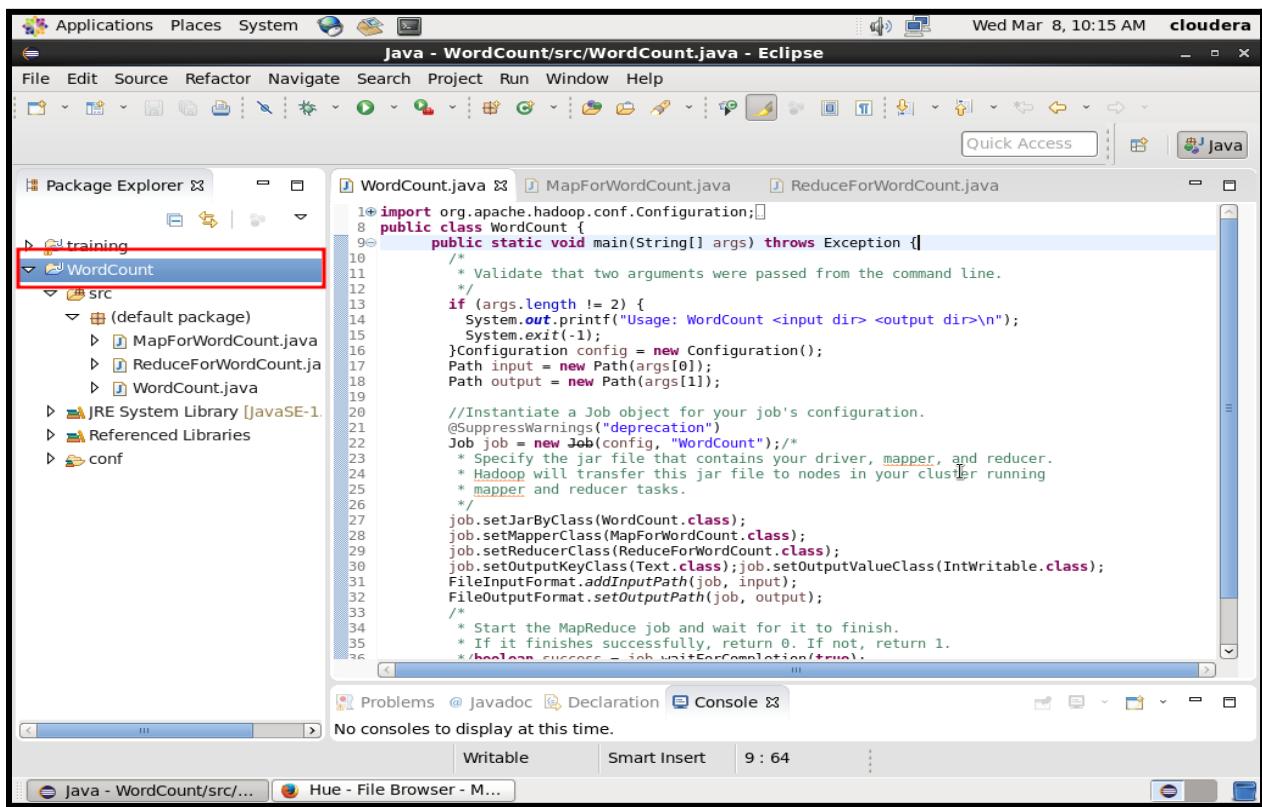
Click Save as to see the directory where the file save

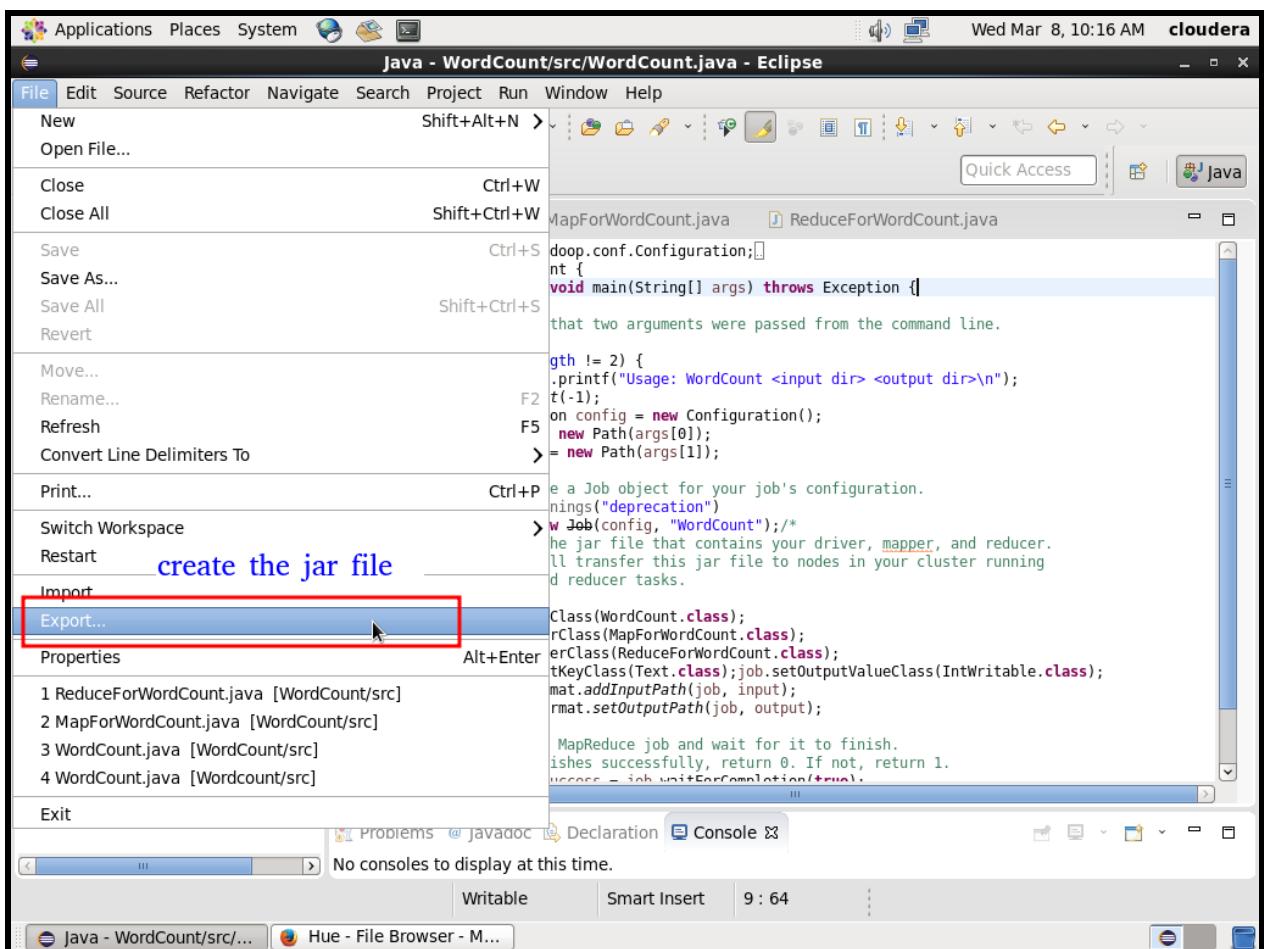
Save Save as

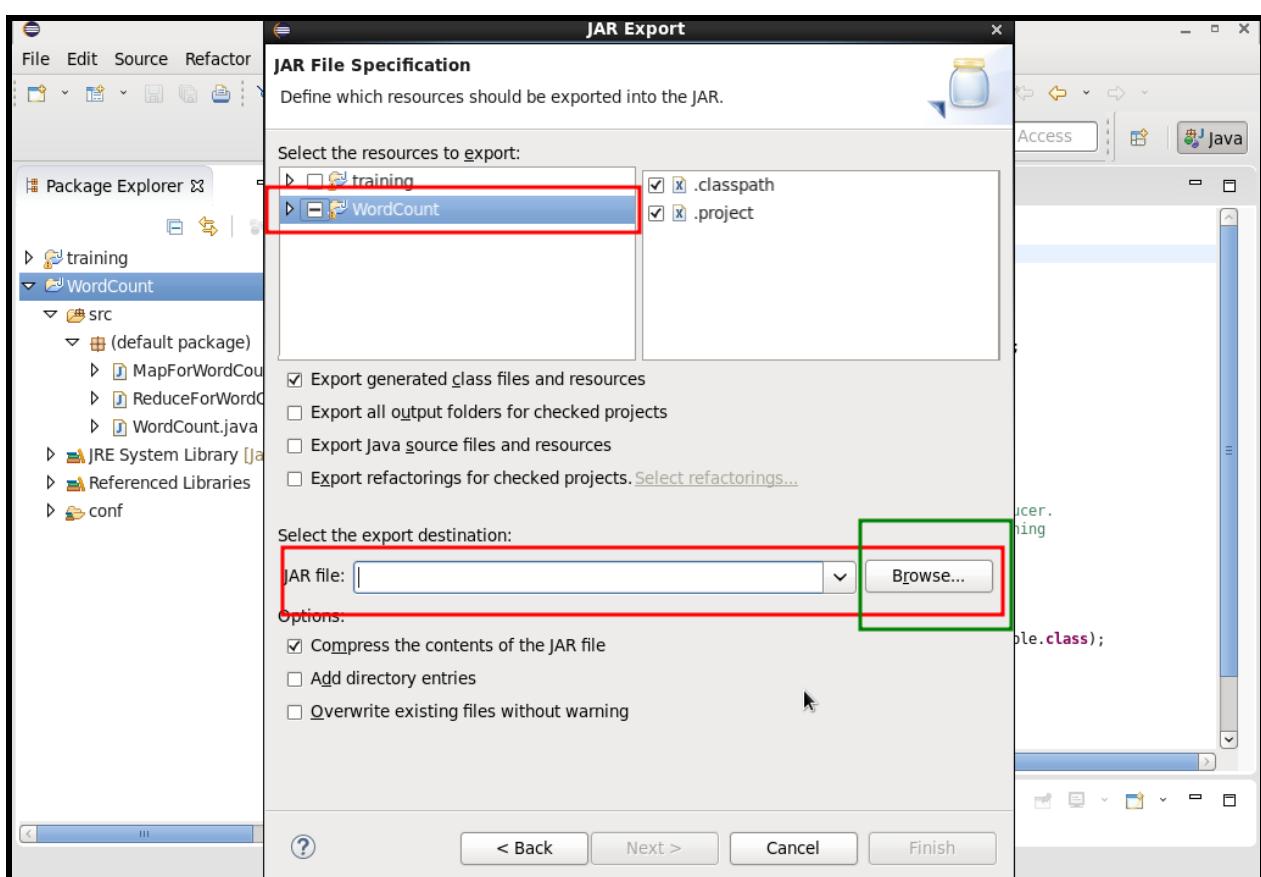
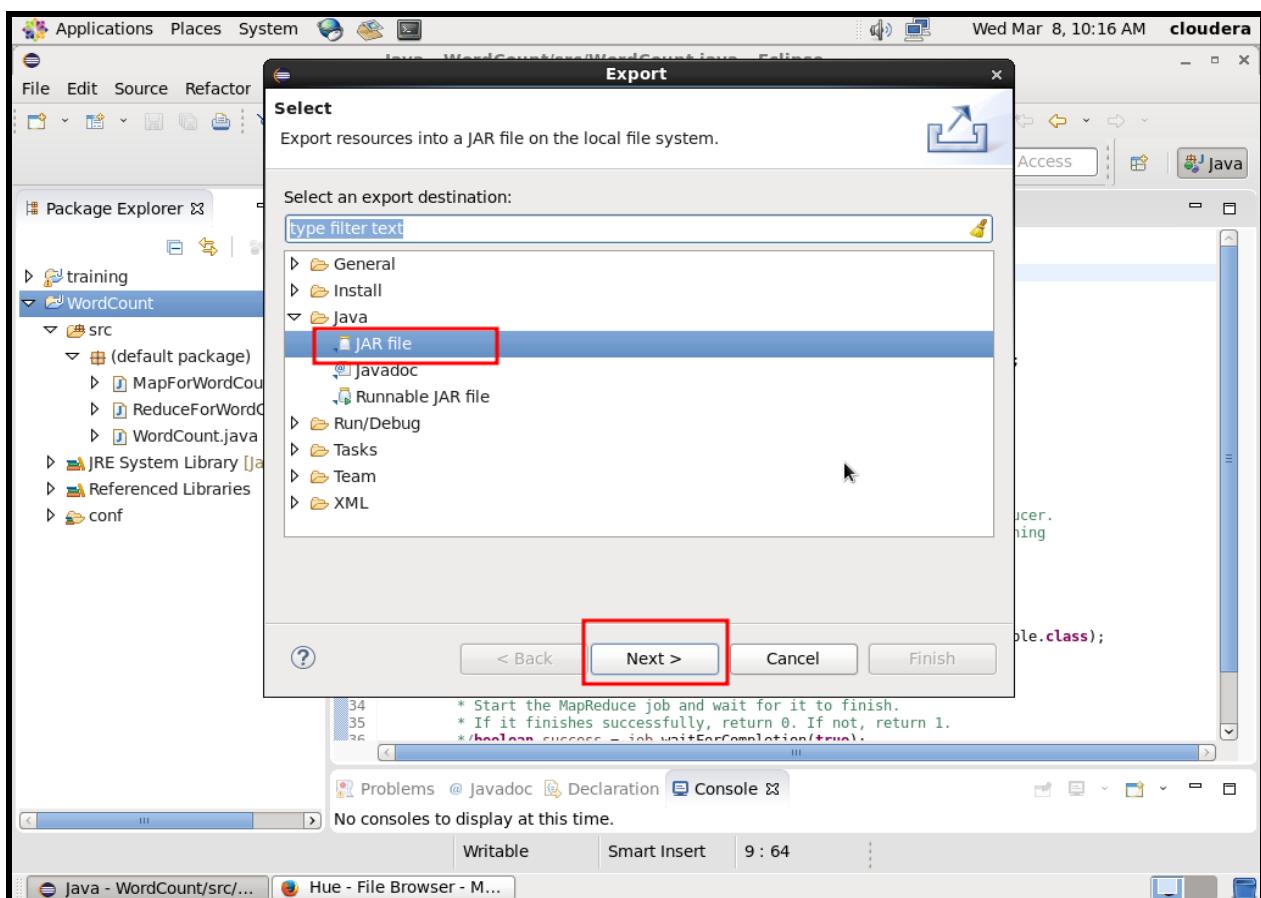


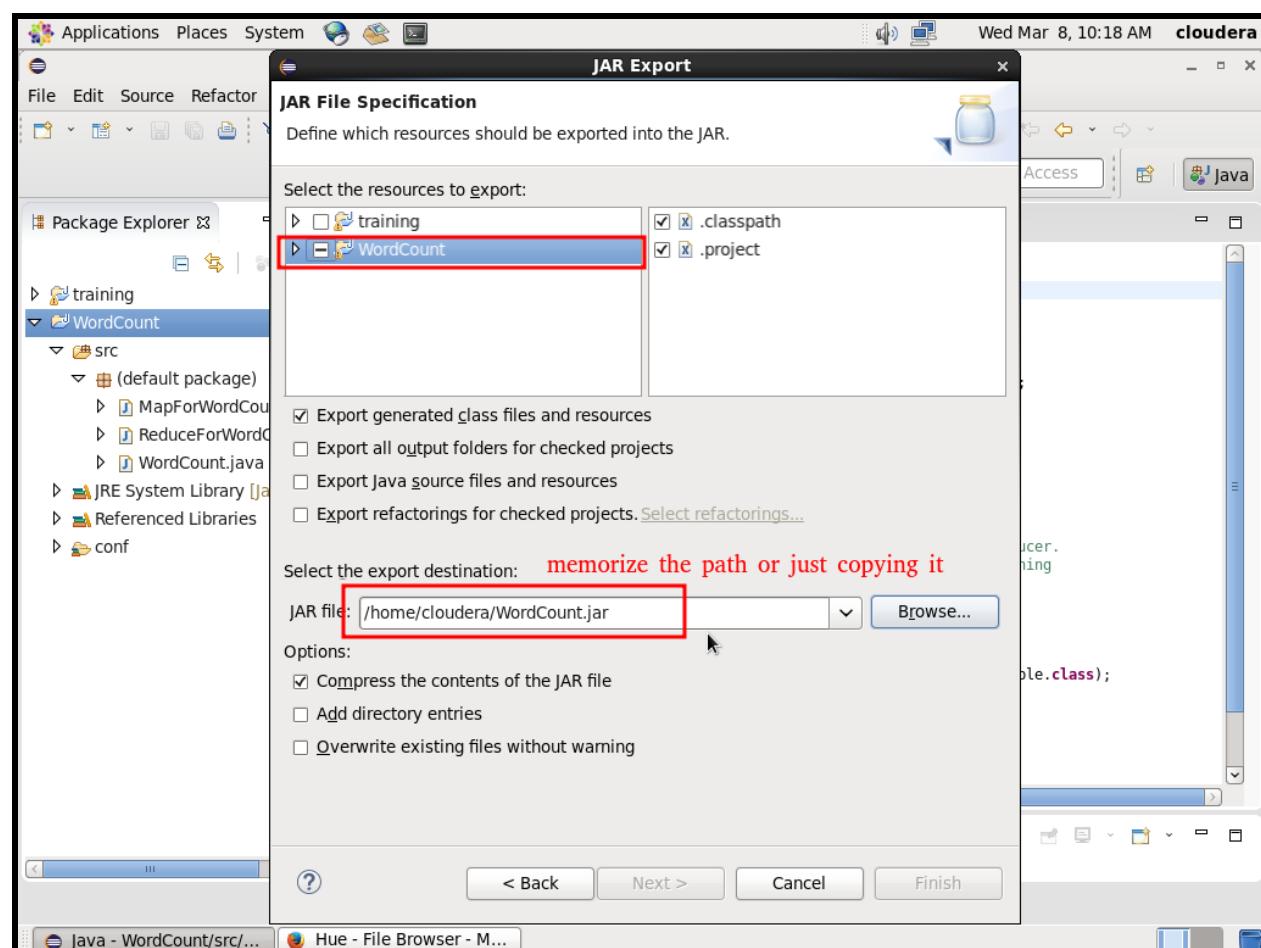
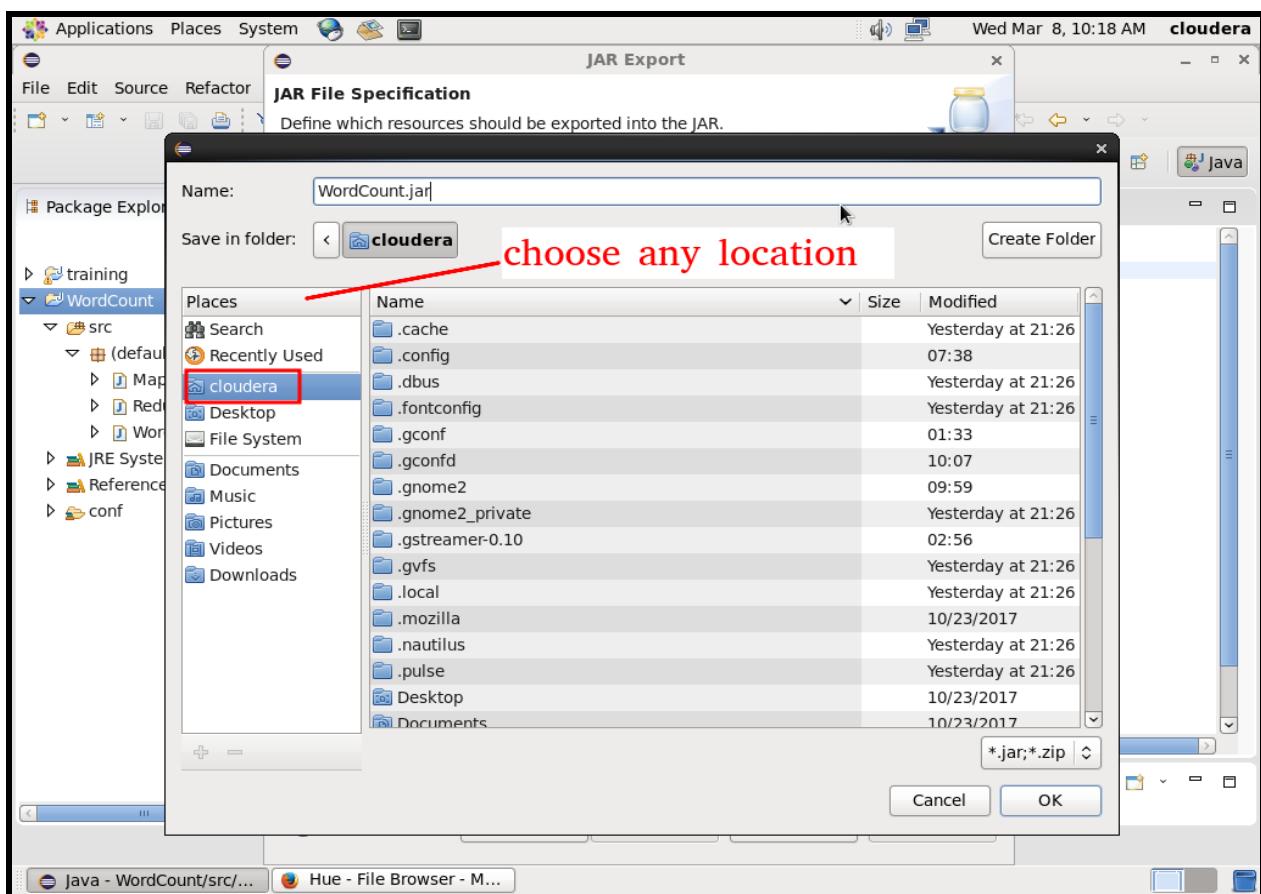


5. Jar file creation

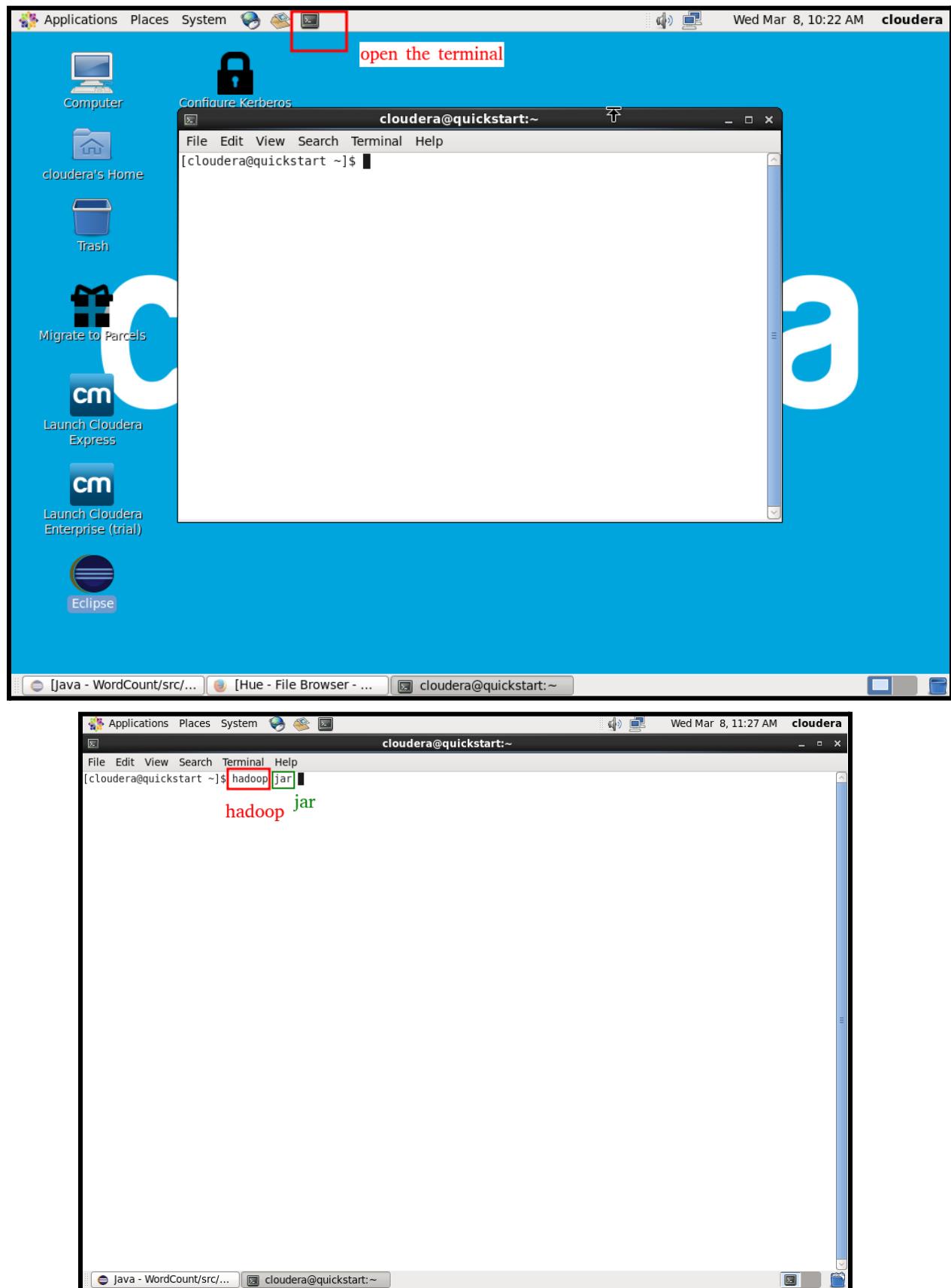


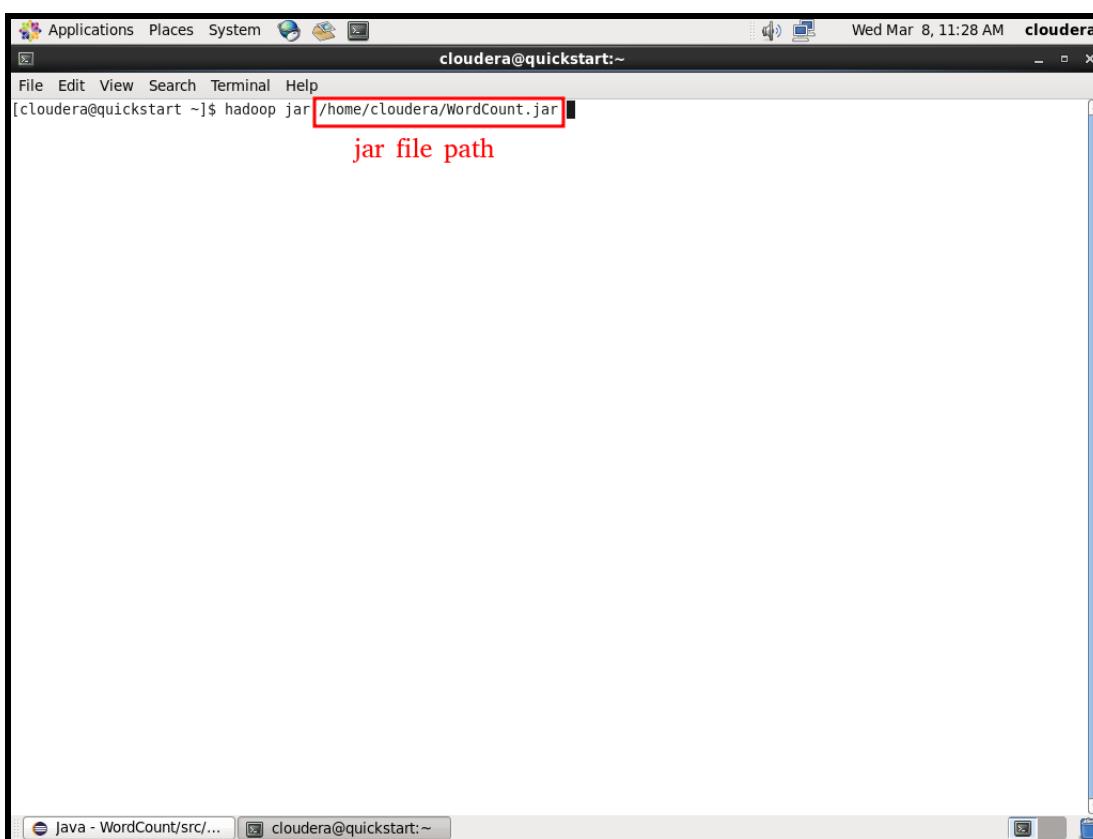
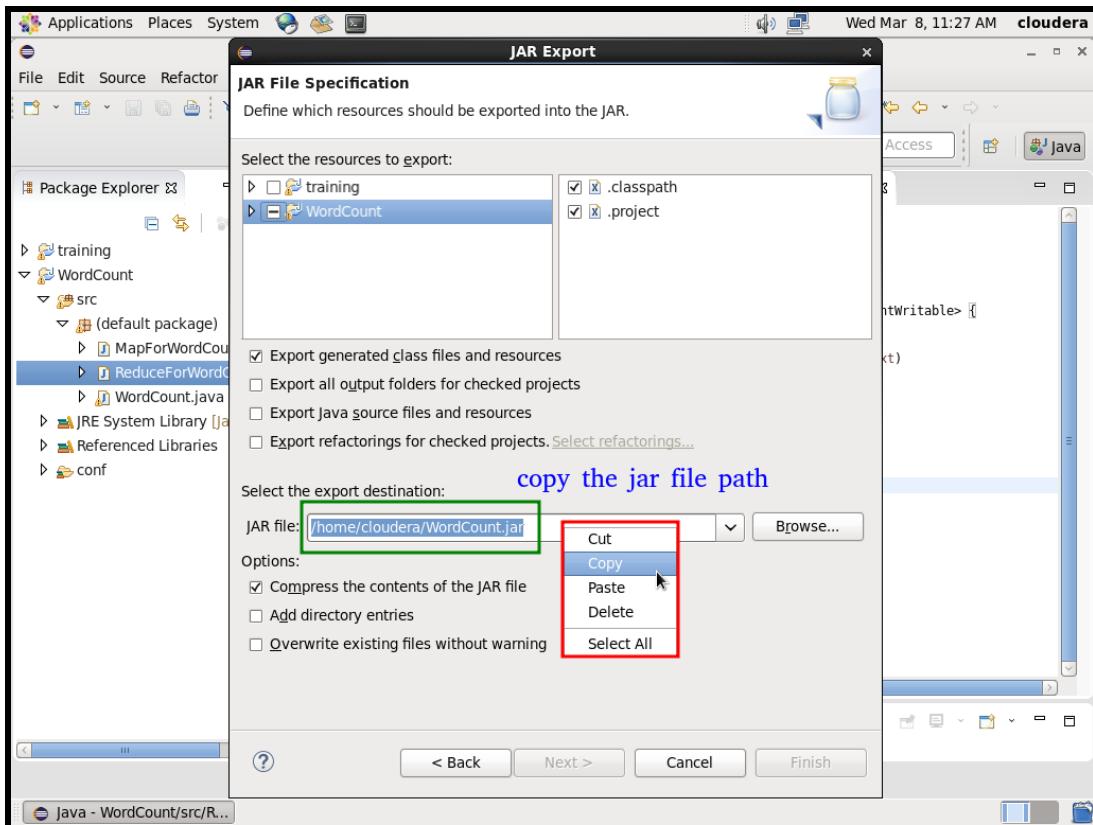


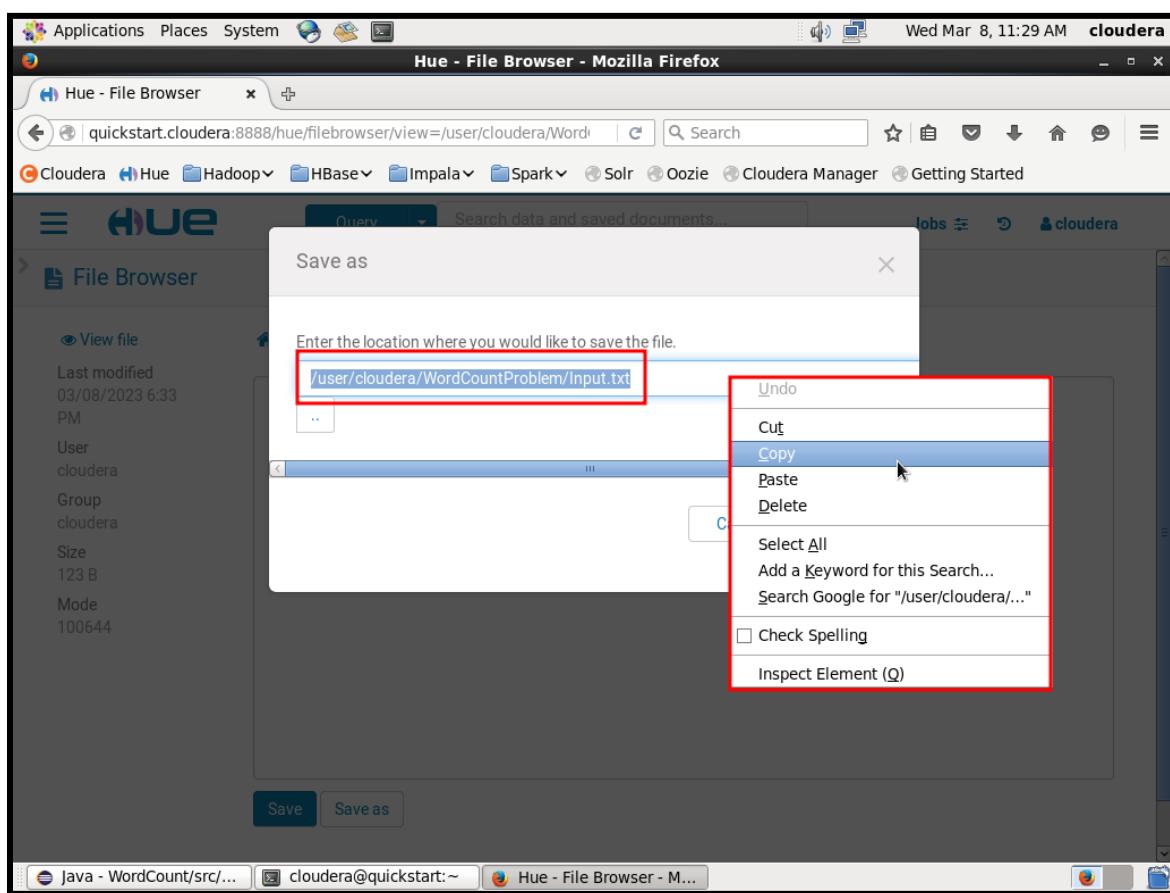
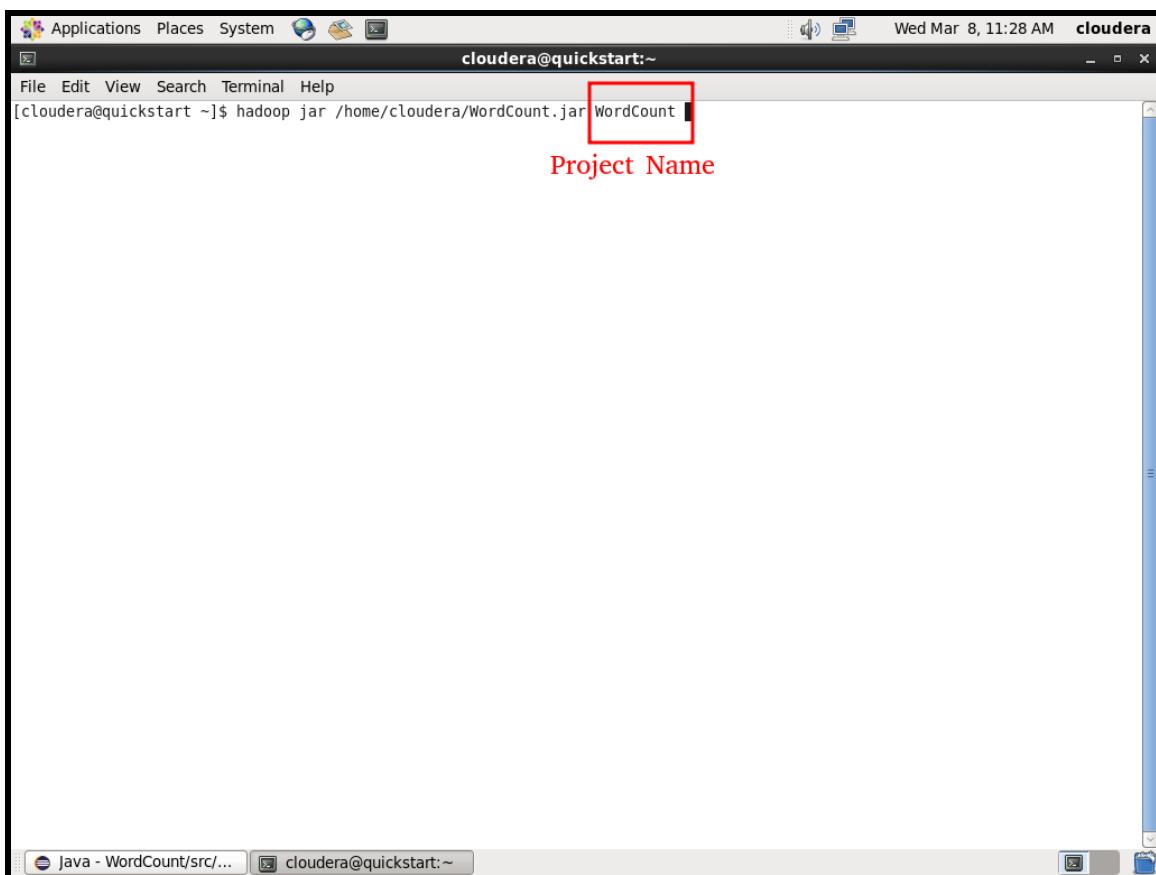


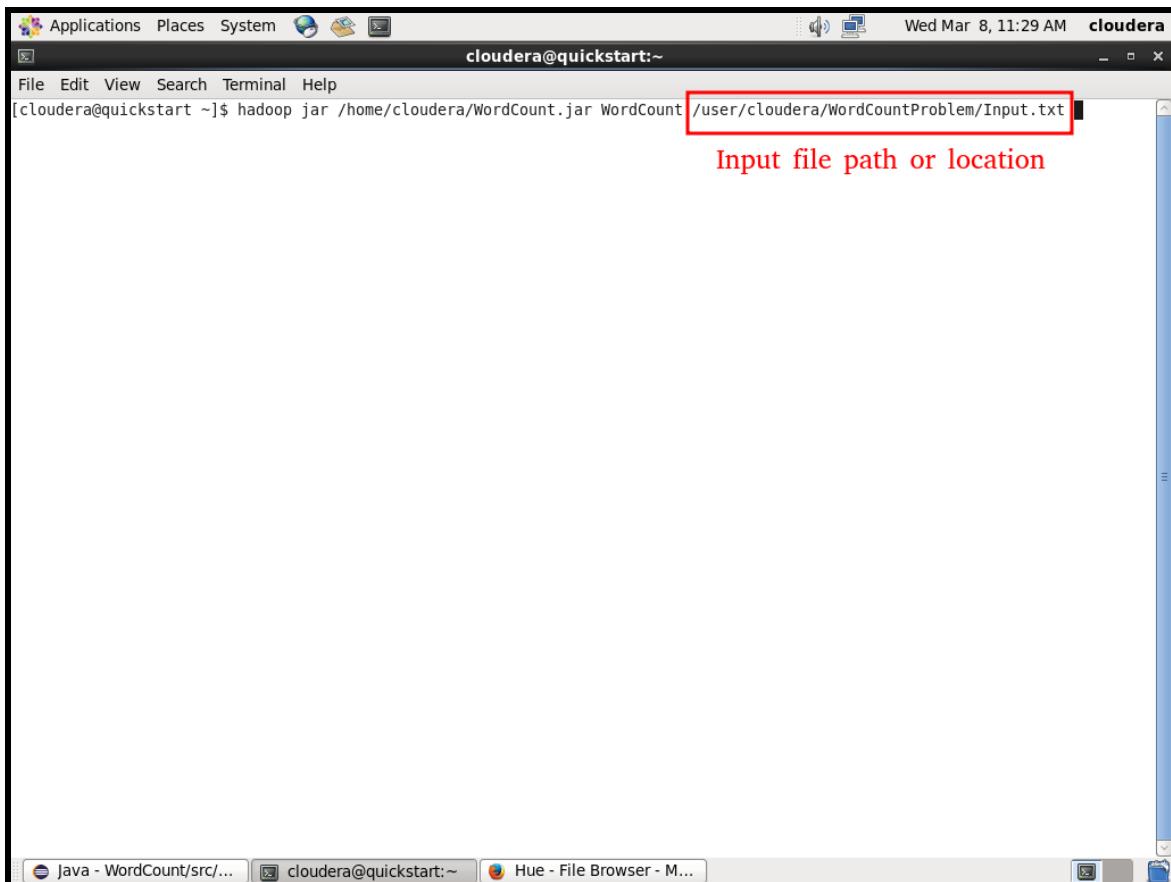


6. Now it's time to opening the terminal and writing some commands to showing the Output for our Input.txt file



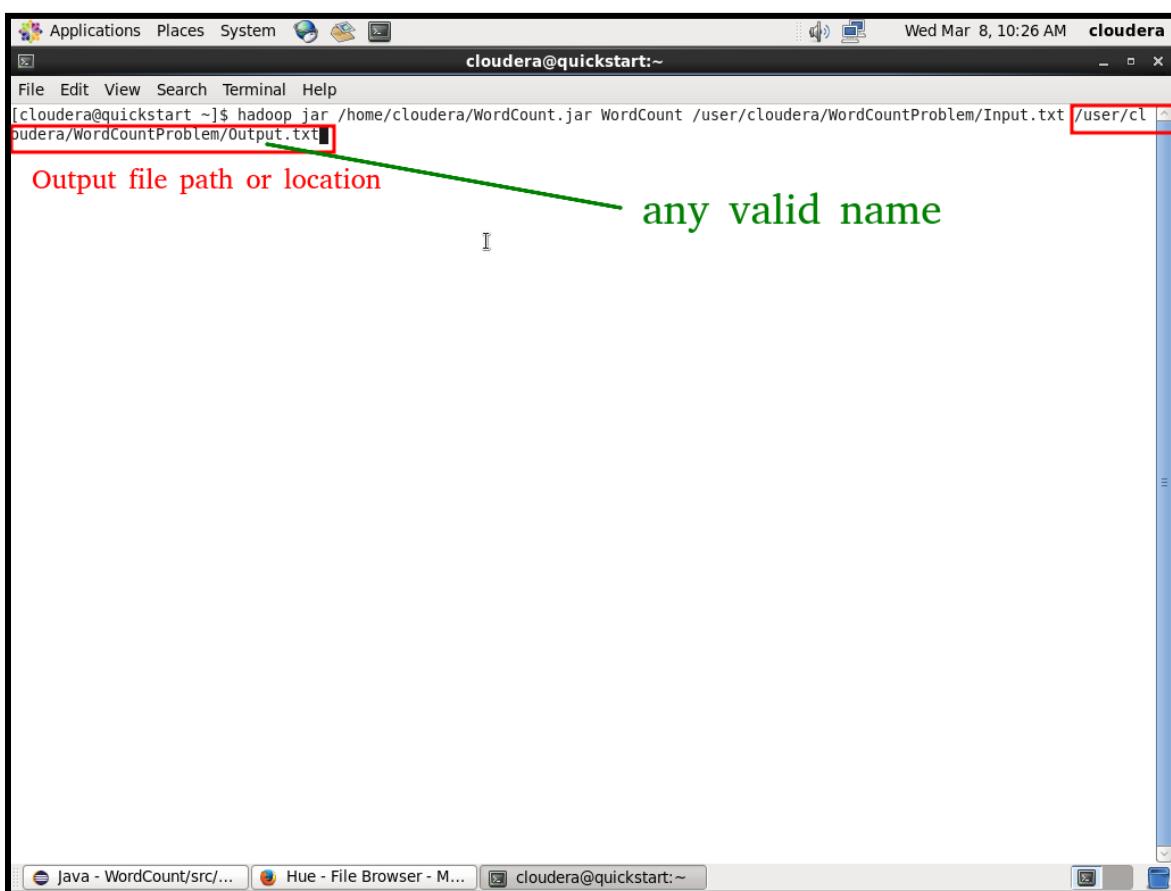






Applications Places System cloudera@quickstart:~
File Edit View Search Terminal Help
[cloudera@quickstart ~]\$ hadoop jar /home/cloudera/WordCount.jar WordCount /user/cloudera/WordCountProblem/Input.txt

Input file path or location



Applications Places System cloudera@quickstart:~
File Edit View Search Terminal Help
[cloudera@quickstart ~]\$ hadoop jar /home/cloudera/WordCount.jar WordCount /user/cloudera/WordCountProblem/Input.txt /user/cloudera/WordCountProblem/Output.txt

Output file path or location

any valid name

Now press enter button on keyboard

A screenshot of a terminal window titled "cloudera@quickstart:~". The window shows the command "hadoop jar /home/cloudera/WordCount.jar WordCount /user/cloudera/WordCountProblem/Input.txt /user/cloudera/WordCountProblem/Output.txt" being typed into the terminal. The background of the window is white, and the text is black. The window has a standard Linux-style title bar and a scroll bar on the right side.

A screenshot of a terminal window titled "cloudera@quickstart:~". The window displays the output of a Hadoop job execution. The log starts with INFO messages from the client.RMProxy and mapreduce.JobSubmitter, followed by a warning about command-line option parsing. It then shows the input paths, number of splits, tokens submitted, application ID, and the URL to track the job. The log continues with INFO messages for the Job, including the start of the job, running mode, map and reduce percentages, and the successful completion of the job with 49 counters. At the bottom, it lists File System Counters and Job Counters, including task counts and total times spent in occupied slots. The background of the window is white, and the text is black. The window has a standard Linux-style title bar and a scroll bar on the right side.

```
[cloudera@quickstart ~]$ hadoop jar /home/cloudera/WordCount.jar WordCount /user/cloudera/WordCountProblem/Input.txt /user/cloudera/WordCountProblem/Output.txt
23/03/08 10:33:48 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
23/03/08 10:33:49 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
23/03/08 10:33:49 INFO input.FileInputFormat: Total input paths to process : 1
23/03/08 10:33:49 INFO mapreduce.JobSubmitter: number of splits:1
23/03/08 10:33:50 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1678267961831_0004
23/03/08 10:33:50 INFO impl.YarnClientImpl: Submitted application application_1678267961831_0004
23/03/08 10:33:50 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1678267961831_0004/
23/03/08 10:33:50 INFO mapreduce.Job: Running job: job_1678267961831_0004
23/03/08 10:34:01 INFO mapreduce.Job: Job job_1678267961831_0004 running in uber mode : false
23/03/08 10:34:01 INFO mapreduce.Job: map 0% reduce 0%
23/03/08 10:34:08 INFO mapreduce.Job: map 100% reduce 0%
23/03/08 10:34:21 INFO mapreduce.Job: map 100% reduce 100%
23/03/08 10:34:21 INFO mapreduce.Job: Job job_1678267961831_0004 completed successfully
23/03/08 10:34:21 INFO mapreduce.Job: Counters: 49
  File System Counters
    FILE: Number of bytes read=166
    FILE: Number of bytes written=286747
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=260
    HDFS: Number of bytes written=67
    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)=5905
    Total time spent by all reduces in occupied slots (ms)=8412
    Total time spent by all map tasks (ms)=5905
    Total time spent by all reduce tasks (ms)=8412
    Total vcore-milliseconds taken by all map tasks=5905
    Total vcore-milliseconds taken by all reduce tasks=8412
```

A screenshot of a terminal window titled "cloudera@quickstart:~". The window displays the output of a Hadoop WordCount job. The output includes various performance metrics such as total time spent by map and reduce tasks, map output records, and shuffle errors. It also shows file input and output counters like bytes read and written.

```
Total time spent by all map tasks (ms)=5905
Total time spent by all reduce tasks (ms)=8412
Total vcore-milliseconds taken by all map tasks=5905
Total vcore-milliseconds taken by all reduce tasks=8412
Total megabyte-milliseconds taken by all map tasks=6046720
Total megabyte-milliseconds taken by all reduce tasks=8613888
Map-Reduce Framework
  Map input records=6
  Map output records=6
  Map output bytes=148
  Map output materialized bytes=166
  Input split bytes=137
  Combine input records=0
  Combine output records=0
  Reduce input groups=3
  Reduce shuffle bytes=166
  Reduce input records=6
  Reduce output records=3
  Spilled Records=12
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=169
  CPU time spent (ms)=1400
  Physical memory (bytes) snapshot=354127872
  Virtual memory (bytes) snapshot=3015303168
  Total committed heap usage (bytes)=240390144
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=123
File Output Format Counters
  Bytes Written=67
[clo...]
```

7. Now going back to the browser Hue section and then respectively the directory which we created before (WordCountProblem) we can see one new file added (Output.txt) or another valid name which we provide before.

Hue - File Browser - Mozilla Firefox

quickstart.cloudera:8888/hue/filebrowser/view=/user/cloudera/WordCountProblem

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

File Browser

Search for file name Actions Move to trash Upload New

Home / user / cloudera / WordCountProblem

Name	Size	User	Group	Permissions	Date
.		cloudera	cloudera	drwxr-xr-x	March 08, 2023 10:09 AM
.		cloudera	cloudera	drwxr-xr-x	March 08, 2023 10:26 AM
Input.txt	0 bytes	cloudera	cloudera	-rw-r--r--	March 08, 2023 10:14 AM
Output.txt		cloudera	cloudera	drwxr-xr-x	March 08, 2023 10:26 AM

Show 45 of 2 items Page 1 of 1

Java - WordCount/src/... Hue - File Browser - M... cloudera@quickstart:~

Hue - File Browser - Mozilla Firefox

quickstart.cloudera:8888/hue/filebrowser/view=/user/cloudera/WordCountProblem/Output.txt

Cloudera Hue Hadoop HBase Impala Spark Solr Oozie Cloudera Manager Getting Started

File Browser

Search for file name Actions Move to trash Upload New

Home / user / cloudera / WordCountProblem / Output.txt

Name	Size	User	Group	Permissions	Date
.		cloudera	cloudera	drwxr-xr-x	March 08, 2023 10:33 AM
.		cloudera	cloudera	drwxr-xr-x	March 08, 2023 10:34 AM
_SUCCESS	0 bytes	cloudera	cloudera	-rw-r--r--	March 08, 2023 10:34 AM
part-r-00000	67 bytes	cloudera	cloudera	-rw-r--r--	March 08, 2023 10:34 AM

Show 45 of 2 items Page 1 of 1

Java - WordCount/src/... Hue - File Browser - M... cloudera@quickstart:~

