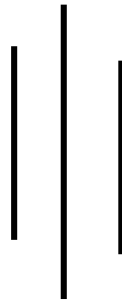# Single Node Cluster Setup with Cloudera Distributed Hadoop (CDH)

**Big Data (CS522)**

**Maharishi University of Management**

Submitted By

Jivan Nepali, 985095

**Jun 06, 2016**

# Download and Install Oracle VirtualBox

Use the following link to download the VirtualBox setup file:

https://www.virtualbox.org/wiki/Downloads



The version, I downloaded, was 5.0.14 –



# Download and Install Cloudera Distributed Hadoop (CDH)

Use the following link:

http://www.cloudera.com/downloads/cdh/5-7-0.html

The VM file is about 5 GB, so have patience before it gets completely downloaded! We can download the latest stable version. I'd already downloaded the CDH5.5 version –



# Configure VM

## Import CDH VM into VirtualBox

In the VirtualBox, go to File > Import Appliance, select correct appliance, then click Next and Import

## Configure Settings

Click on Settings menu and configure as like below:

Boot Order:    Hard Disk, Optical

cloudera-quickstart-vm-5.5.0-0-virtualbox - Settings                          ?    ✕
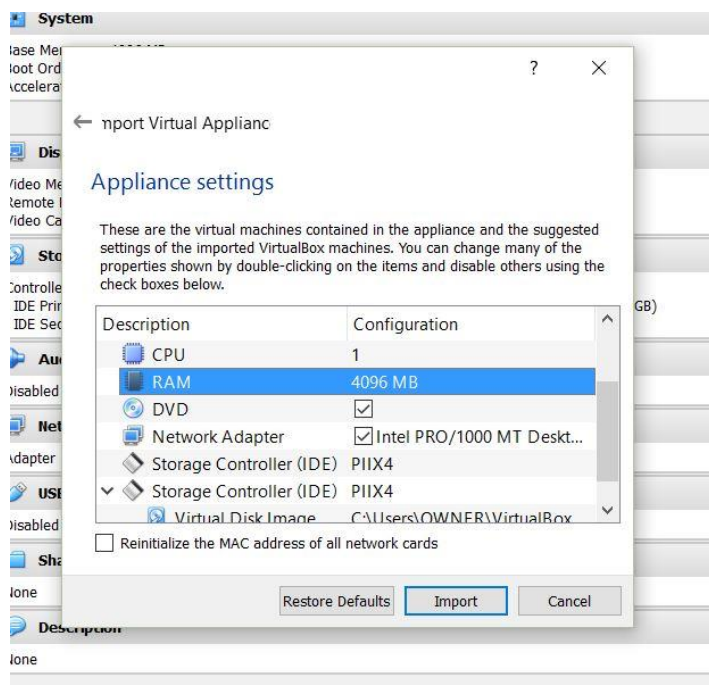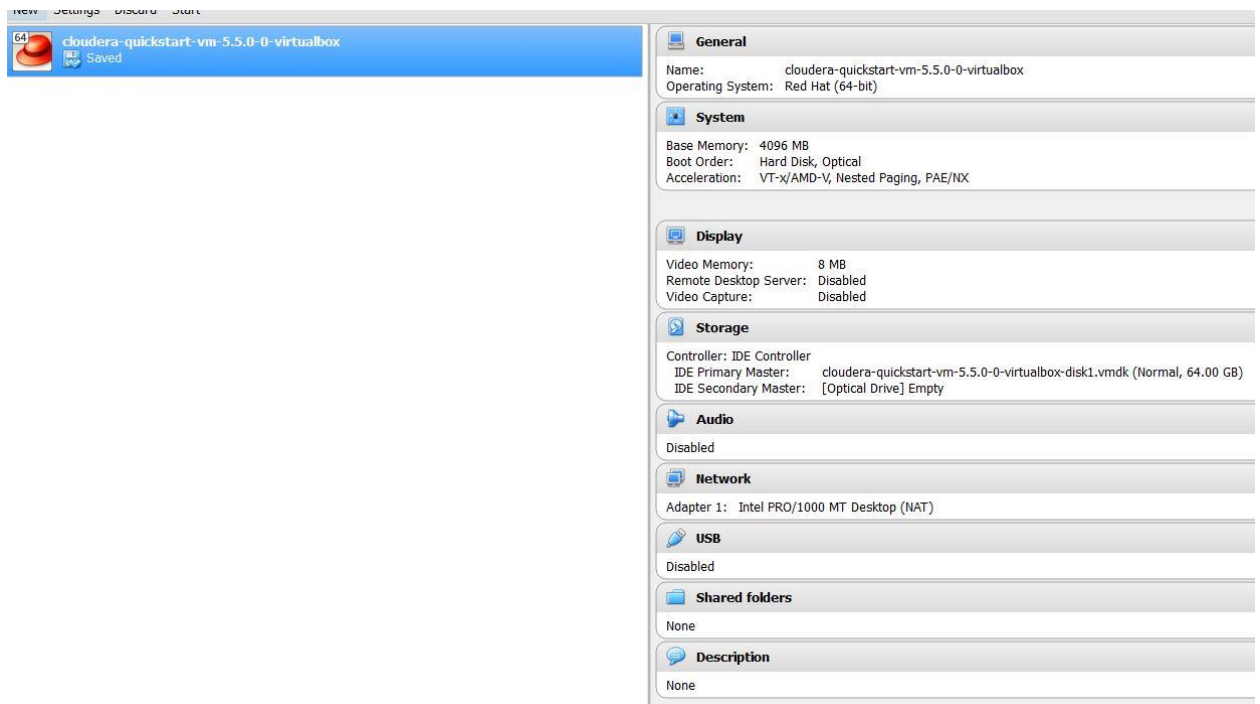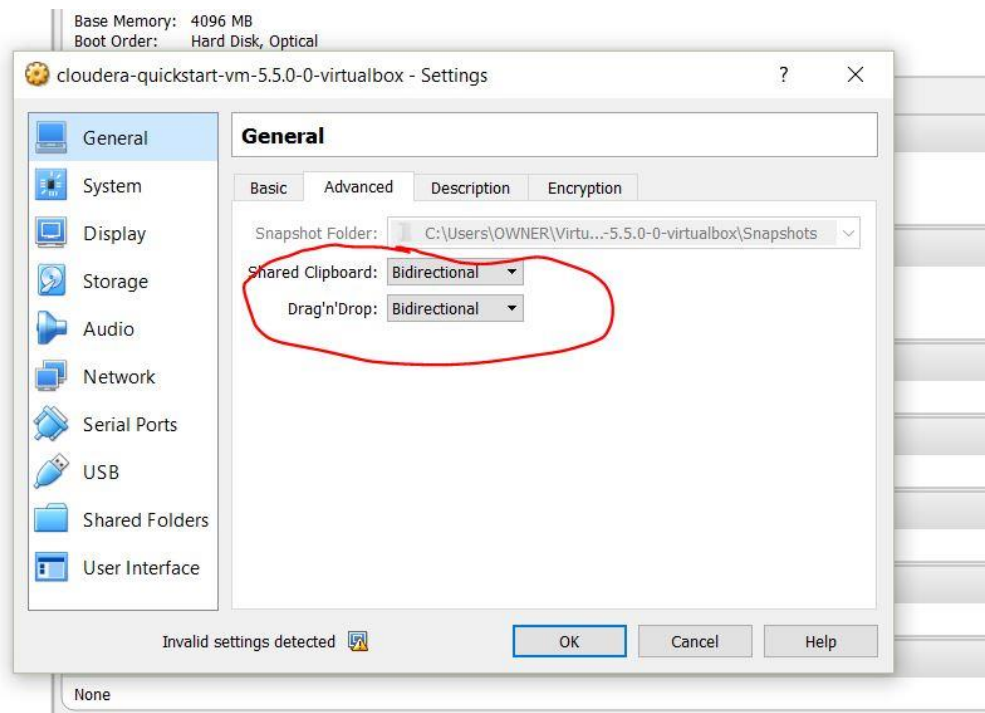
**Network**

| General | | | | |
| System | | | | |
| Display | | | | |
| Storage | | | | |
| Audio | | | | |
| Network | | | | |
| Serial Ports | | | | |
| USB | | | | |
| Shared Folders | | | | |
| User Interface | | | | |

Adapter 1    Adapter 2    Adapter 3    Adapter 4

☑ Enable Network Adapter

Attached to:  NAT  ▼

Name:  ▼

▼ Advanced

Adapter Type:  Intel PRO/1000 MT Desktop (82540EM)  ▼

Promiscuous Mode:  Deny  ▼

MAC Address:  080027FD90B7  🔄

☑ Cable Connected

Port Forwarding

Invalid settings detected 🔲        OK       Cancel       Help

None

Start CDH

Use the Browser on the Host Machine to browse the Cloudera Quick Start Page
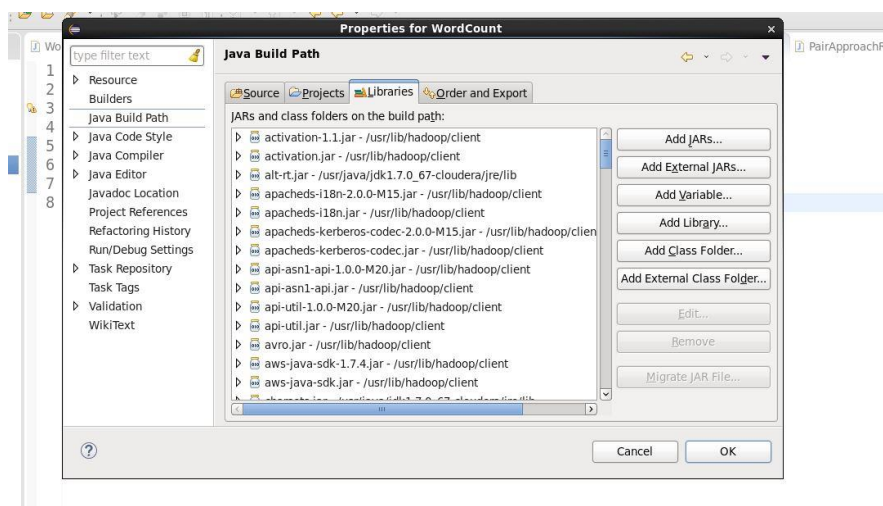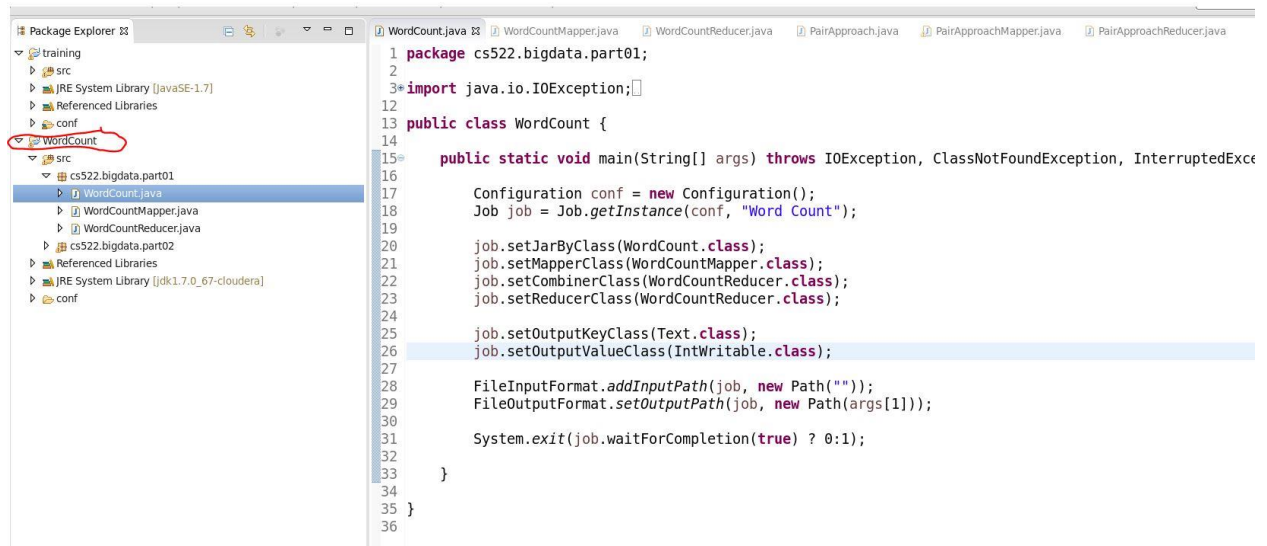


# Eclipse Project Setup

## Create New Java Project

First, create a Java project in Eclipse inside the Host Machine (Linux). And, then add all the Hadoop related libraries as needed, for example –

*/usr/lib/Hadoop/Hadoop\*.jar*
*/usr/lib/Hadoop/client/\*.jar*
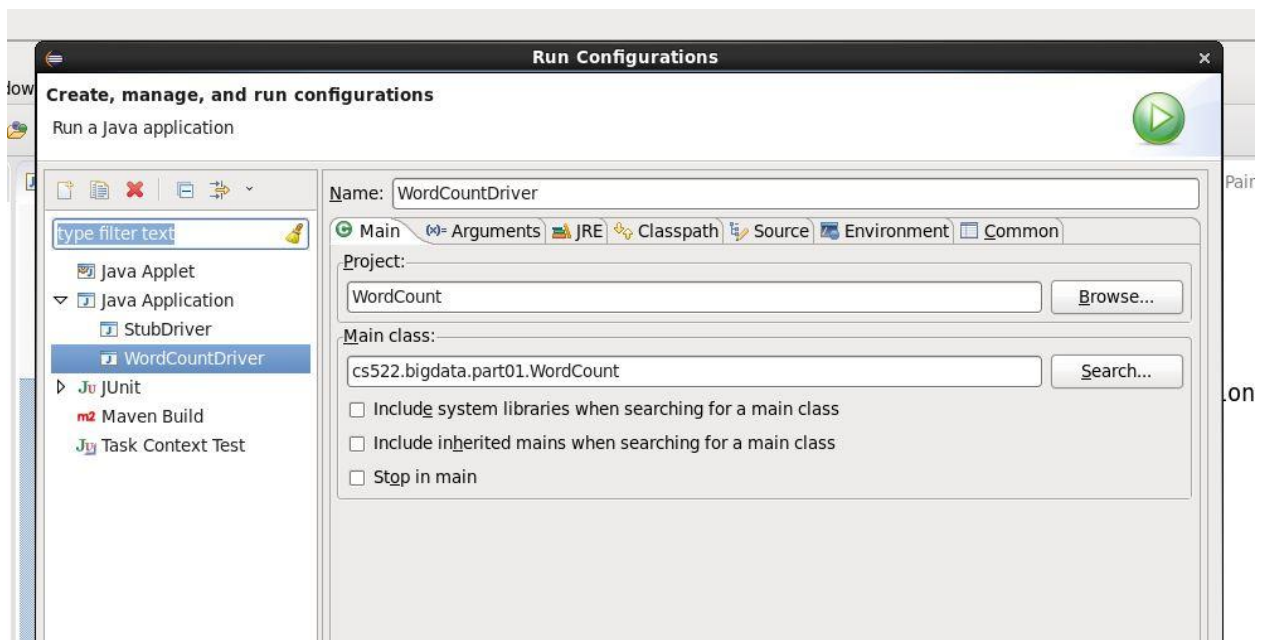
# Write the WordCount Mapper and Reducer classes along with main() Method



```
package cs522.bigdata.part01;

import java.io.IOException;

public class WordCount {

    public static void main(String[] args) throws IOException, ClassNotFoundException, InterruptedExce

        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "Word Count");

        job.setJarByClass(WordCount.class);
        job.setMapperClass(WordCountMapper.class);
        job.setCombinerClass(WordCountReducer.class);
        job.setReducerClass(WordCountReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        FileInputFormat.addInputPath(job, new Path(""));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        System.exit(job.waitForCompletion(true) ? 0:1);

    }
}
```
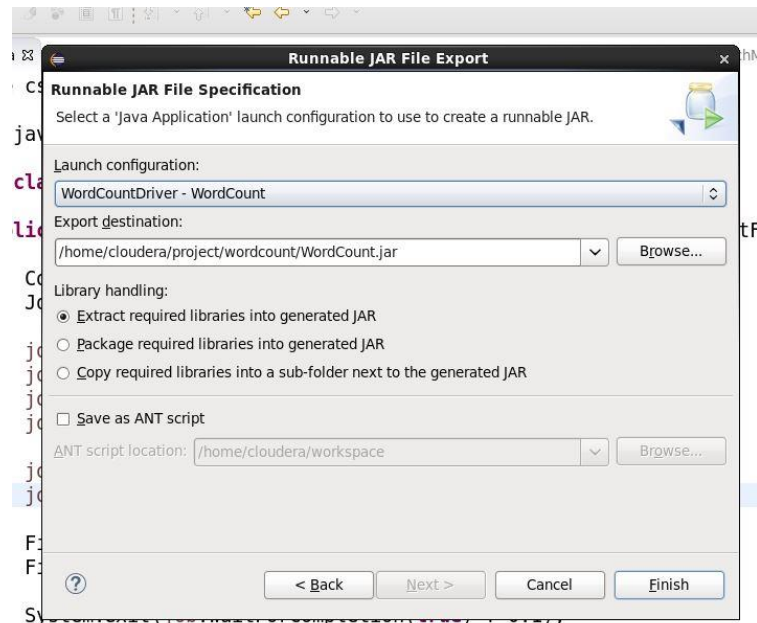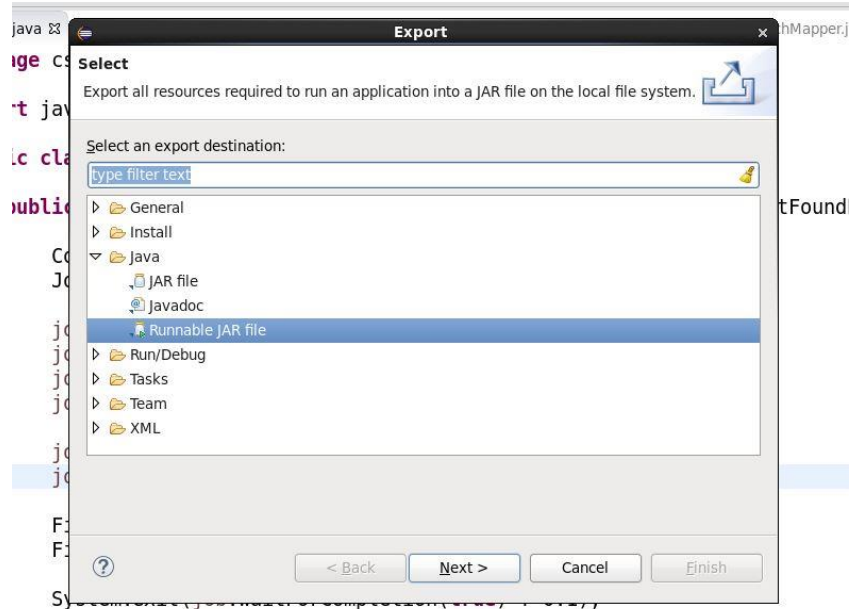
# Export the Jar File

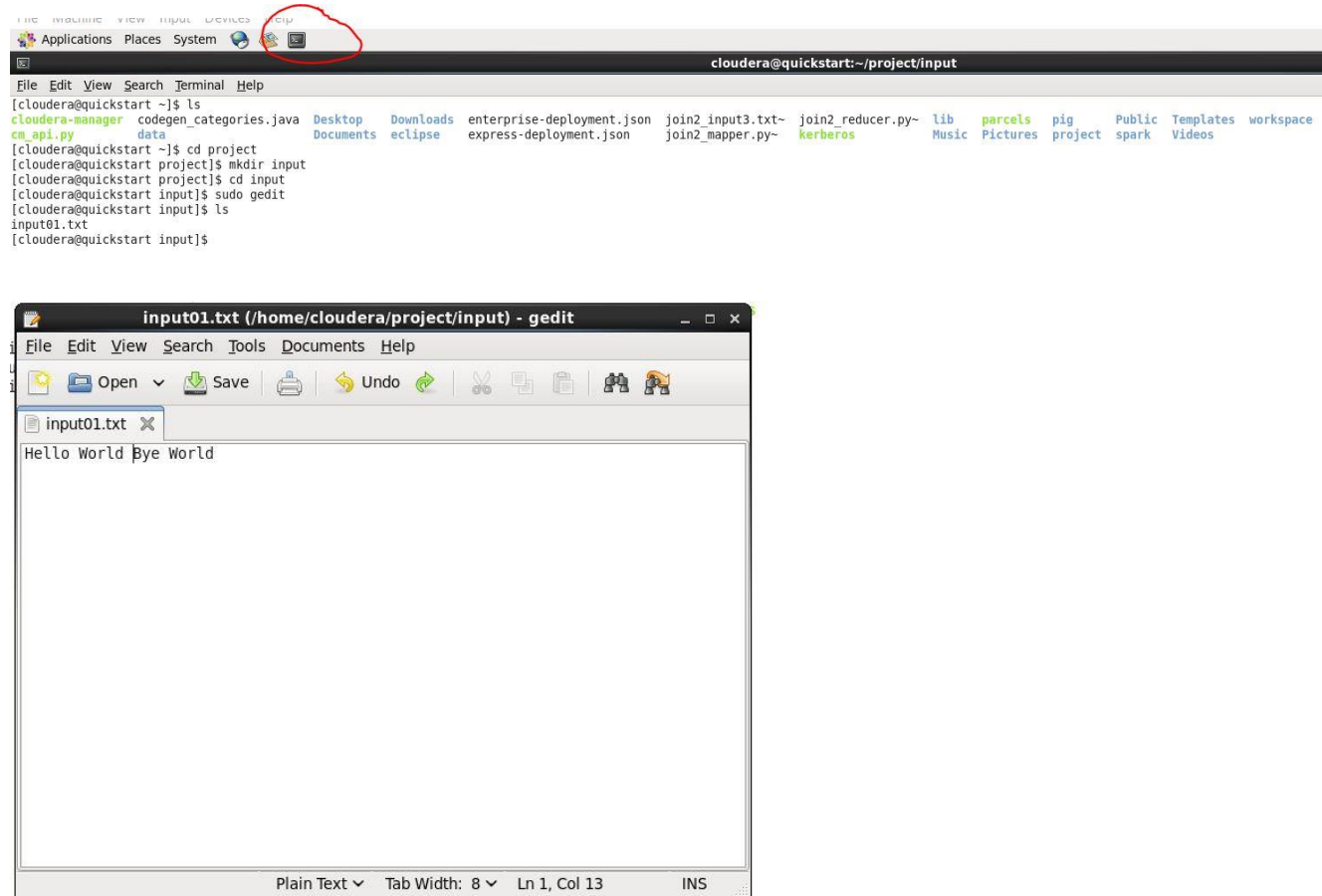First, setup the Run Configurations as –

Right click on the project and click on Export –



And, then click on finish button. It will finally export the jar file into the specified directory.

# Run the Job

First, create an input file for the wordcount as follows –





Copy the file from local Linux FileSystem to the HDFS using the following commands –
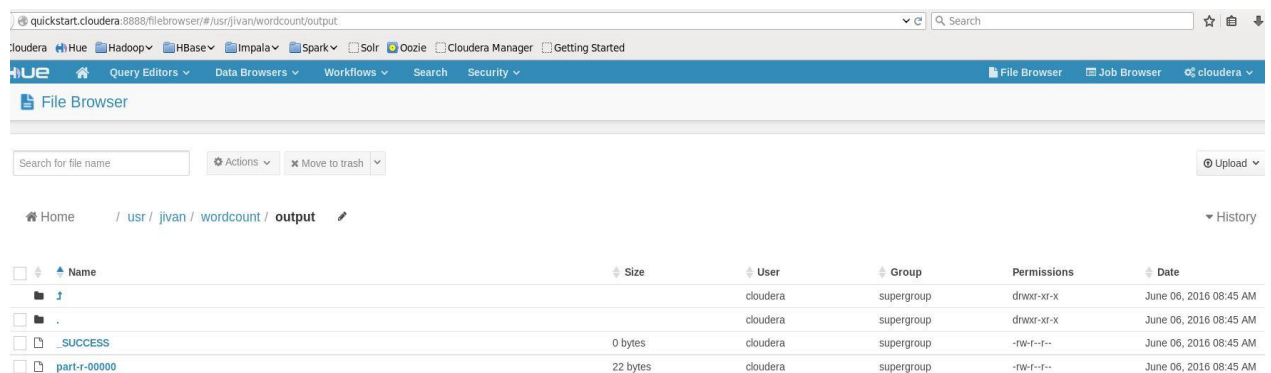


Now, run the Jar file as follows, make sure the output directory does not exists as it will be created during execution –

This job execution takes a little bit time and when the map task and reduce task complete to 100%, then check for the output file in the output directory as specified in the above command –

```
File  Edit  View  Search  Terminal  Help
[cloudera@quickstart wordcount]$ hadoop fs -ls /usr/jivan/wordcount/output
Found 2 items
-rw-r--r--   1 cloudera supergroup          0 2016-06-06 08:45 /usr/jivan/wordcount/output/_SUCCESS
-rw-r--r--   1 cloudera supergroup         22 2016-06-06 08:45 /usr/jivan/wordcount/output/part-r-00000
[cloudera@quickstart wordcount]$ hadoop fs -cat /usr/jivan/wordcount/output/part-r-00000
Bye     1
Hello   1
World   2
[cloudera@quickstart wordcount]$
```

And, it verifies the word count result! We can upload the input/check output files to/in HDFS using the Hue Interface from the browser as follows –



(The default credentials are – user: cloudera, pass: cloudera)

Next – we'll be working on the Pair Approach implementation for finding the relative frequencies.