## CS590RA Tutorial 2

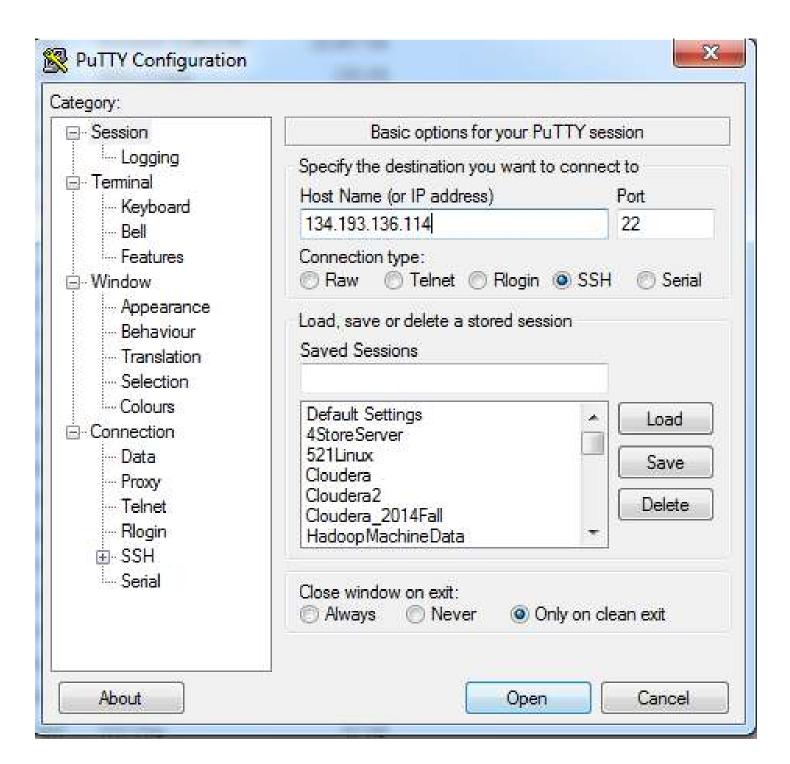
Cloudera/MapReduce

### Introduction

- Platform for Apache Hadoop and its ecosystem
- CDH (Cloudera Distribution Including Apache Hadoop) is its open source Apache Hadoop distribution

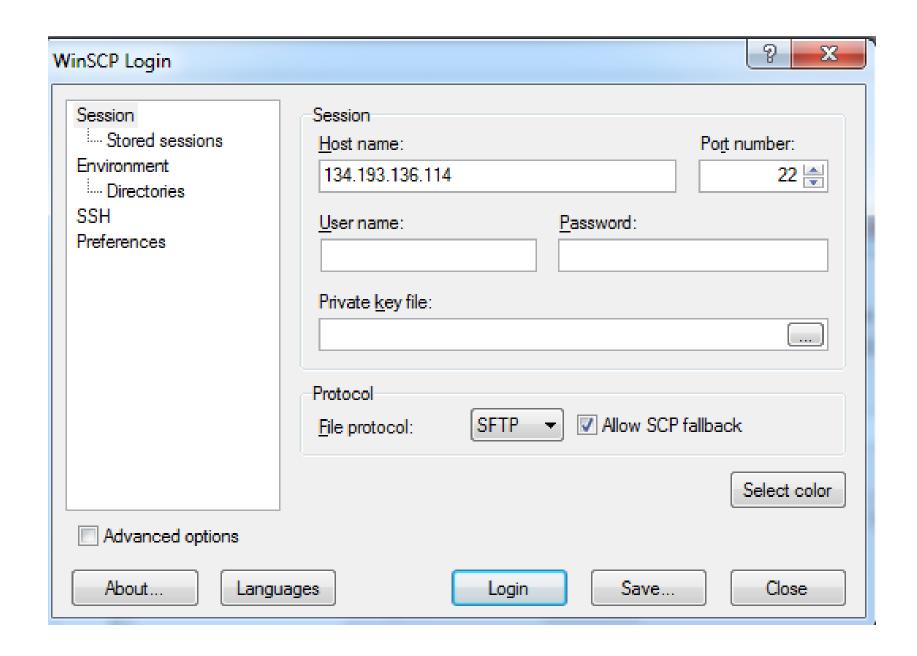
#### Remote Linux Machine

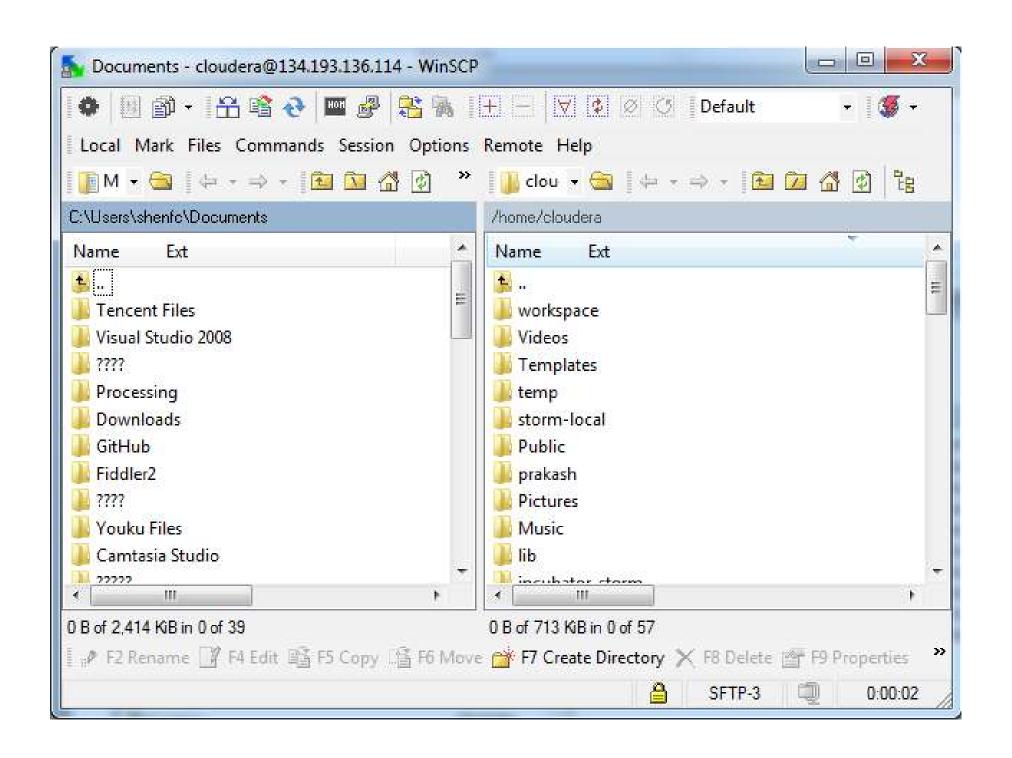
- How to login?
- Download putty
   http://www.chiark.greenend.org.uk/~sgtatham/p
   utty/
- Host name is 134.193.136.114
- Each group will have an account, user name and password are same
- For example, for group1, user name is group1, password is also group1
- You can change your password by typing passwd



#### SSH & FTP tools

- WinSCP provids a very easy way to transfer files remotely from local windows machine to remote linux machine
- http://winscp.net/eng/download.php





### SSH & FTP tools

 For Linux and Mac users, Filezilla is also a good choice. You can download it from <a href="https://filezilla-project.org/download.php?type=client">https://filezilla-project.org/download.php?type=client</a>

How to use Filezilla:

https://www.youtube.com/watch?v=MCGctsHKi
zw

#### 1. Run WordCount

Basic Hadoop command:

```
Hadoop fs —Is (list everything under HDFS)
Hadoop fs —rm xxx (remove file from HDFS)
Hadoop fs —rm —r xxx (remove dir from HDFS)
Hadoop fs —put xxx xxx (put file from local to HDFS)
Hadoop fs —cat xxx (concatenate message to screen)
Hadoop jar xxx.jar (run jar file in Hadoop)
```

### 1. Run WordCount

1. Download wordcount example from this link.

<a href="https://portal.futuregrid.org/manual/hadoop-wordcount">https://portal.futuregrid.org/manual/hadoop-wordcount</a>

1. Download and unzip WordCount under \$HADOOP\_HOME

Assuming vou start SalsaHadoop/Hadoop with setting \$HADOOP\_HOME=~/hadoop-0.20.203.0, and are running the master node on i55, download and unzip the WordCount source code from Big Data for Science tutorial under \$HADOOP\_HOME.

• 2. Use FTP tool to transfer the project

#### 1. Run WordCount

- 3. Unzip the wordcount.zip (unzip xxx.zip)
- 4. Go to wordcount folder
- 5. put local input file to the hadoop input directory: hadoop fs –put LocalInput HDFSInput
- 6. Run hadoop:

hadoop jar wordcount.jar WordCount input output3

7. View result from output3

hadoop fs -cat output3/\*

 Instead of running hadoop command hadoop jar wordcount.jar WordCount input output3

we can write a jar file to complete the same job and specify input and ouput dir in the java code

- To finish the job, you have to include hadoop library:
  - hadoop-core-1.0.0.jar

```
public class WordCount {
public static void main(String[] args) {
JobClient client = new JobClient();
JobConf conf = new JobConf(WordCount.class);
FileInputFormat.addInputPath(conf, new Path("input"));/// The dir here is the Hadoop directory
FileOutputFormat.setOutputPath(conf, new Path("out7"));
conf.setOutputKeyClass(Text.class);
conf.setOutputValueClass(LongWritable.class);
conf.setMapperClass(TokenCountMapper.class);
conf.setCombinerClass(LongSumReducer.class);
conf.setReducerClass(LongSumReducer.class);
client.setConf(conf);
try {
JobClient.runJob(conf);
} catch (Exception e) {
e.printStackTrace();
```

### Hadoop map/reduce algorithm

#### Map:

```
String line = value.toString();
List[];
///when met a space, consider word count++
String Arr[] = line.split ("");
For (int i=0;i<Arr.length;i++)
   String word = Arr[i];
   List.put(word, 1);
                                                                             1
                                                                      а
                                                                             1
                                                                      a
                            laamand
                                                                             1
                                                                      m
                                                                             1
                                                                      a
                                                                             1
                                                                      n
                                                                      d
                                                                             1
```

## Hadoop map/reduce algorithm

Reduce:

int sum = 0;

1	1
а	3
m	1
n	1
d	1

```
/// Iterator the hashmap, calculate the total number of a word
```

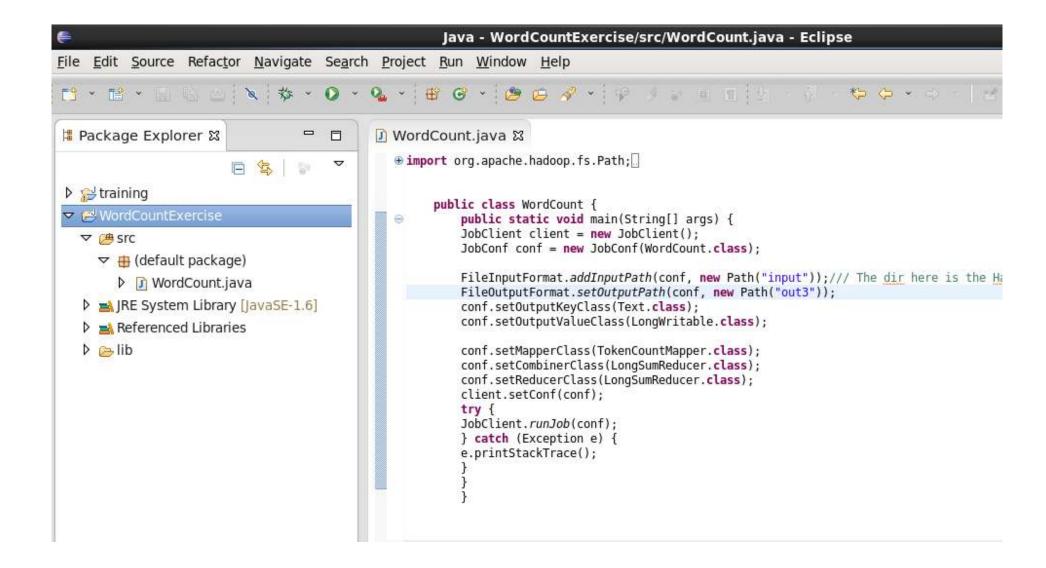
```
for (IntWritable val : values) { sum += val.get(); }
context.write(key, new IntWritable(sum));
```

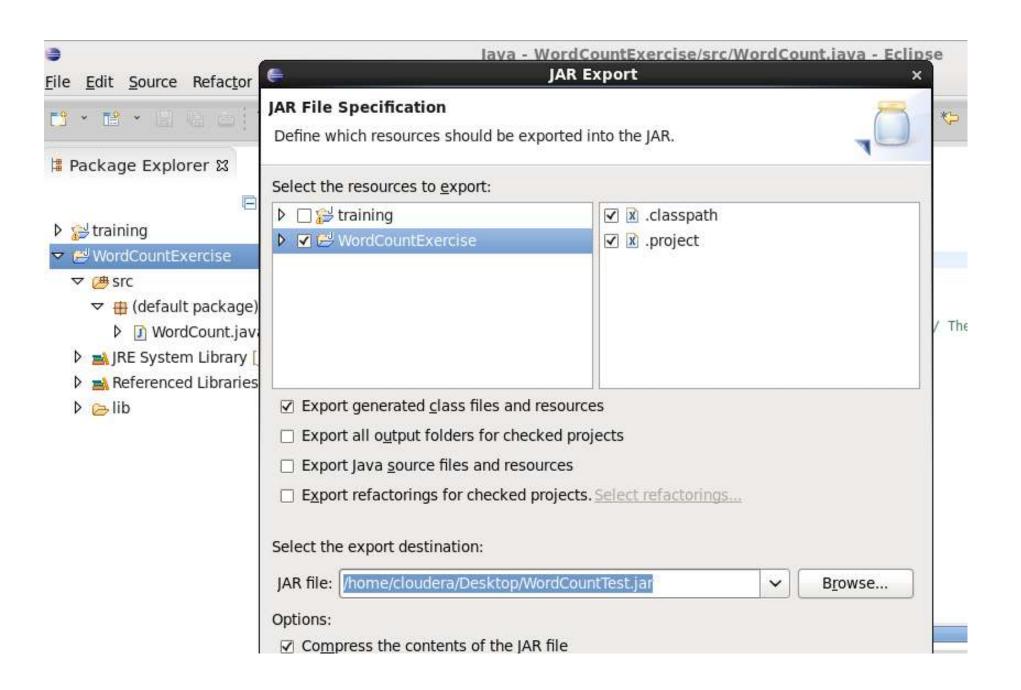
 1). Import the WordCountExercise project into the Eclipse in your laptop

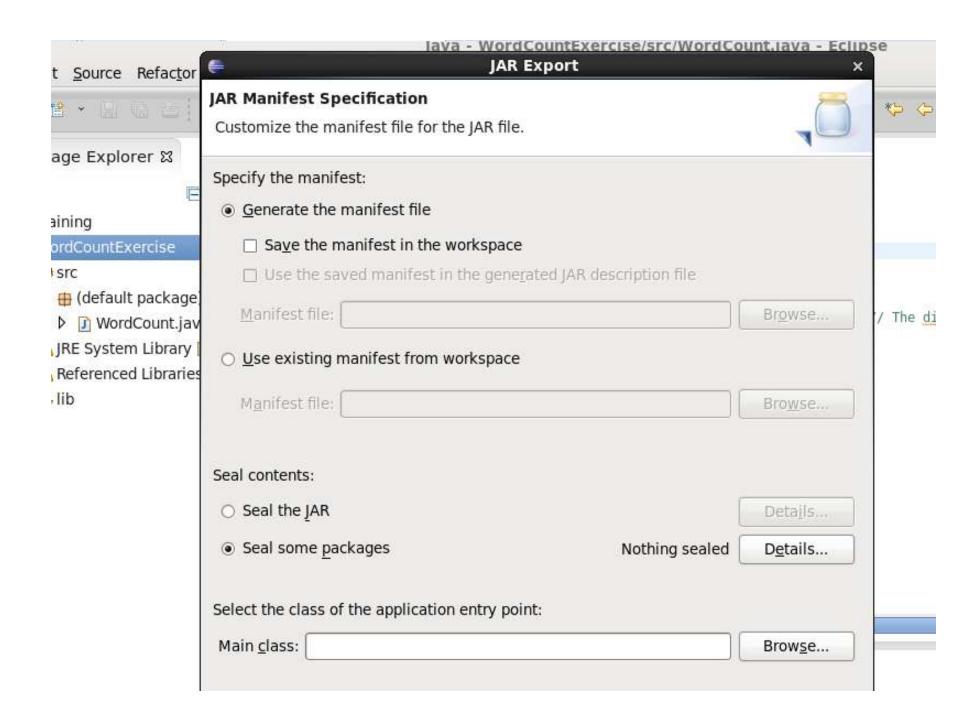
File->import->Existing Projects into Workspace-> Choose "WordCountExercise" project

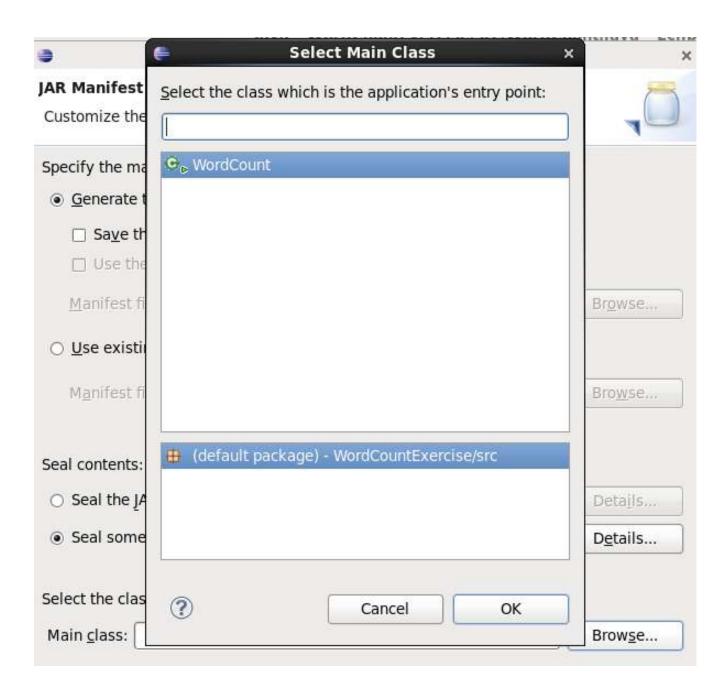
• 2). Export the project as jar file

- 3). How to export:
- a. Right click project
- b. Click export
- c. Select JAR file
- d. Click next, next
- e. Don't forget to select main class by clicking browse to select the main class you want to select.









- 4). Run jar by the command: hadoop jar WordCountExercise.jar
- 5). Check result:
   hadoop fs –cat out3/\*

### Tips

- If you want to run this code again, in java code, you must change the output directory to another new value, like out4 and out5, and export to jar again:
- ☐ FileOutputFormat.setOutputPath(conf, new Path("out4"));

Otherwise, hadoop will report the error like: output directory already exists.

### Thanks