**Readme for Assignment 1**
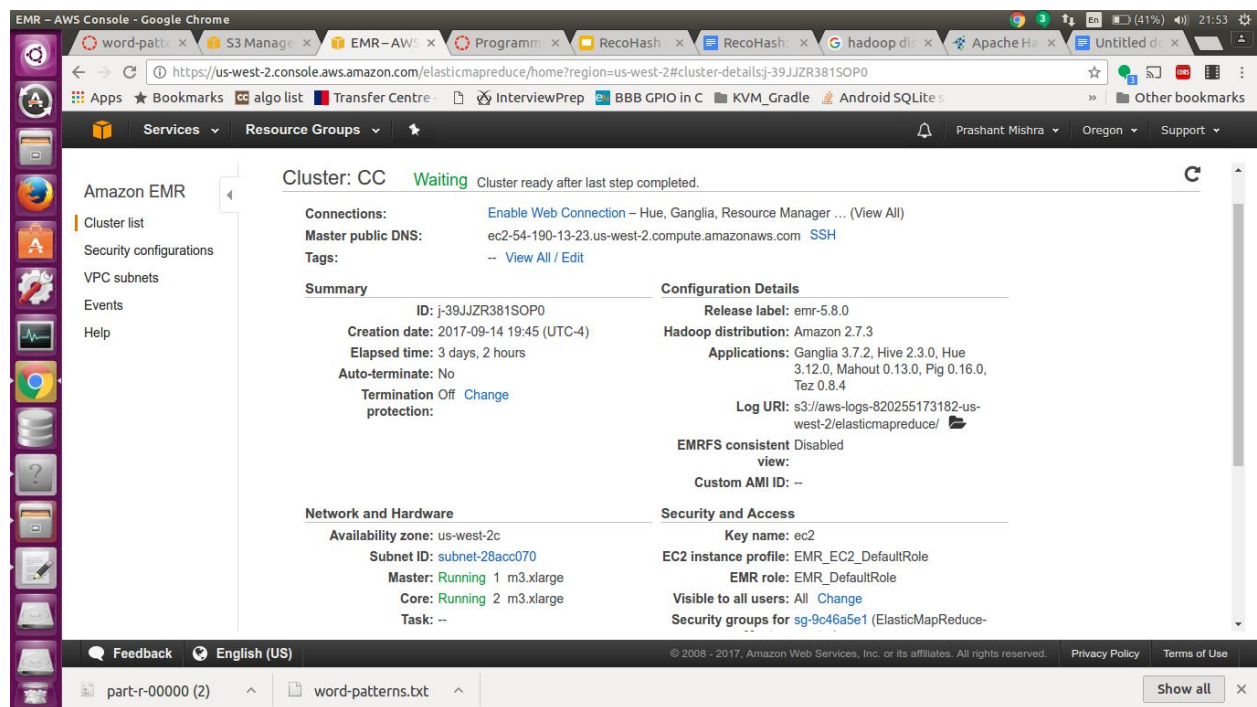
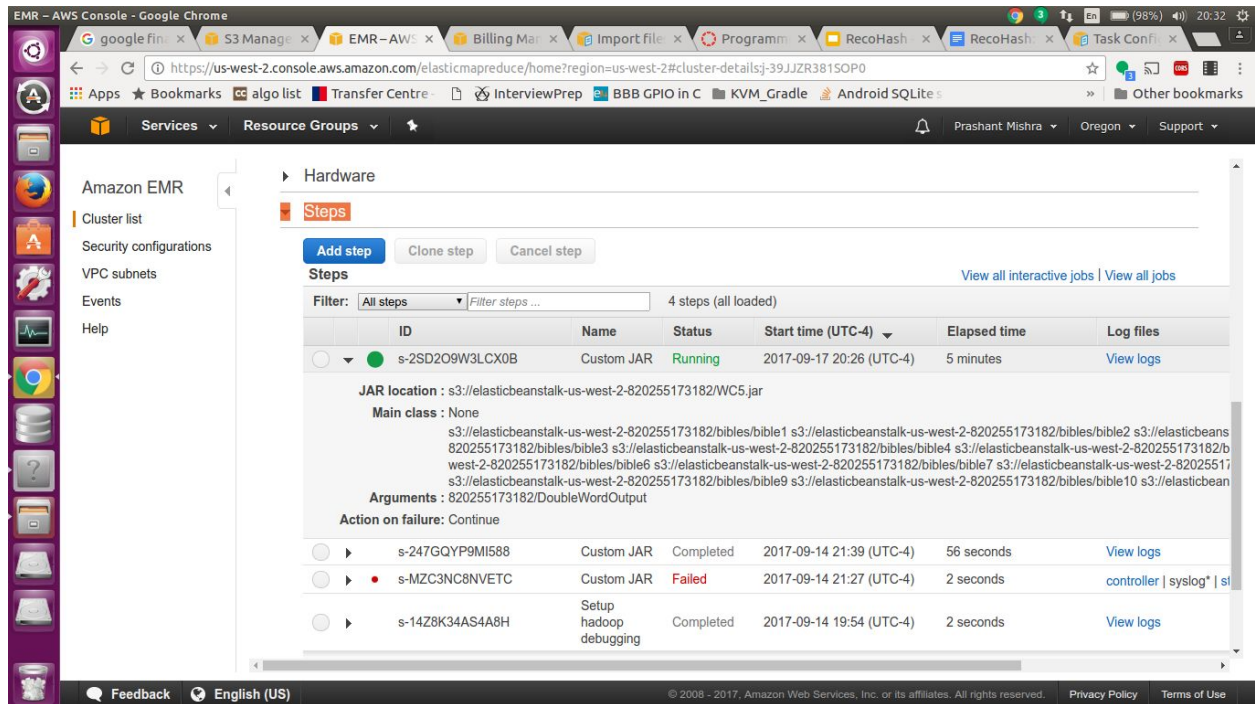**Submission Zip Folder Contents**
1. 3 JAR files: SingleWordCount.jar, DoubleWordCount.jar, DistributedWordCount.jar
2. Readme
3. 3 Output Folders: SingleWordCountOutput, DoubleWordCountOutput, DistributedWordCount

**Cluster Setup**
I have used Amazon AWS EMR ( Elastic Map-Reduce ) to set up a cluster of 3 machines (One master and two slave nodes). Screenshots of the setup below:



Once this cluster is set up, we can add a map-reduce job by adding a new step with a custom jar file that contains all the map-reduce JAVA code. A screenshot of a particular step is show below. All the parameters to the job is stored in AWS S3. As you can see we can specify the location of the custom jar, and location of all the 10 copies of bible and the output folder as arguments to this jar.

## JAR Files
For single word count, I got the code from
https://github.com/uzresk/aws-emr-examples/tree/master/src/main/java/aws/emr/wordcount. For
the other parts I have made slight modifications to the Mapper class to work for the new cases.

Here in the code we have 3 classes, WordCount.java, WordCountMapper.java and
WordCountReducer.java

## WordCount.java
The main method of this class, initiates a new hadoop job and sets the input and output
expected file types and their path from arguments.
Specifically in distributed cache example, this method also sets up the code to add the cache
file to the distributed cache, and because of this, each data node can access this file as if from
its local file system and processing becomes faster.

## WordCountMapper.java
The mapper class has the method map which is common in all three cases. Given a line of text
as the input argument, it splits the line into a list of words, and for each of the words issues a
map with key as the word and its frequency as 1.

For double word count, I used the key as word[currentIndex] + "," + word[currentIndex + 1] to
get pairs of words as key. The rest of the code remains the same.

For distributed cache example, the mapper class has two additional methods, setup and readFile. Setup is an overidden method that always runs before the map method. It gets all the files present in the distributed cache, and then calls readFile on each of those files. ReadFile method basically parses the file line by line and adds all distinct words in a Hash called importantWords. Now the map method has access to this Hash and it issues a map command only for those words that are found in importantWords.

**Output Generation**
Each of the nodes in the hadoop cluster in AWS EMR has some cores and for fast processing, each cores writes to its own output file. This option can be modified to write all output to a single file but it's slower, so my output is contained in multiple output part files. The output files themselves are sorted by key.

The location of the output files for this assignment can be found in the 3 different folders: SingleWordCountOutput, DoubleWordOutput and DistributedWordCountOutput. Each of these folders contains 7 part files which contain the frequency for different distinct words.