# CSE3025 – Large Scale Data Processing
## Lab - 2 – 11/12/2017

PREPARED BY - HARGUR PARTAP SINGH BEDI (15BCE1257)

## Simple MapReduce program in Hadoop

1. Start the hadoop cluster by typing the command

**$ bash hadoop/bin/start-all.sh**

2. Copy a text file for counting the count of each distinct word in it, using command

**$ hadoop fs -copyFromLocal <source file absolute location> <destination file location>**

### Program: WordCount.java

```java
//package com.wordCount;
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount{
    public static class Map extends Mapper<Object, Text, Text, IntWritable>{
            private final static IntWritable one = new IntWritable(1);
            private Text word = new Text();

            public void map(Object key, Text value, Context context) throws
IOException, InterruptedException{
                    StringTokenizer st = new StringTokenizer(value.toString());
                    while(st.hasMoreTokens()){
                            word.set(st.nextToken());
                            context.write(word,one);
                    }
            }
    }
```

```java
        public static class Reduce extends Reducer<Text, IntWritable, Text,
IntWritable>{
                public void reduce(Text key, Iterable<IntWritable> values, Context
context) throws IOException,InterruptedException{
                        int sum = 0;
                        for(IntWritable val: values){
                                sum += val.get();
                        }
                        context.write(key, new IntWritable(sum));
                }
        }

        public static void main(String[] args) throws Exception{
                Configuration conf = new Configuration();
                Job job = new Job(conf,"wordcount");
                job.setJarByClass(WordCount.class);
                job.setMapperClass(Map.class);
                job.setOutputKeyClass(Text.class);
                job.setOutputValueClass(IntWritable.class);
                job.setReducerClass(Reduce.class);
                //job.setInputFormatClass(TextInputFormat.class);
                //job.setOutputFormatClass(TextOutputFormat.class);
                FileInputFormat.addInputPath(job, new Path(args[0]));
                FileOutputFormat.setOutputPath(job, new Path(args[1]));
                job.waitForCompletion(true);
        }
}
```

3.  Now run the WordCount.java file using command

**$ java WordCount.java**

4.  Make a jar file out of it using command

**$ jar cf wc.jar WordCount*.class**

5.  Now run the wordcount mapreduce program on the text file uploaded on Hadoop
    using command

**$ hadoop jar wc.jar WordCount <input textfile hadoop location> <output file
location>**