# Generic Explanation:

This assignment required multiple things such as file handling, string parsing etc. I am using file handling for *positive.txt and negative.txt.* Whereas the input file (*Comments.txt*) is being read by the mapper function itself. I am converting the whole line into a string and then parsing the string by storing the string into a different variable. For this I hardcoded the parser code as I saw the input file and before every comment, there is always a "T".

Another implementation was that in order to see if the reducer has ended, I am using a built-in function Cleaner (). This function is basically called when the reducer has finished its tasks.

Furthermore, for this assignment I am using HDFS as for some reason when I installed HDFS, my Hadoop-standalone stopped working.

Rest of the implementation is requirement specific explained below:

# Average Length of Comments:

For this task, I started finding the length of each comment using the built-in functions. Then I added them in a variable and was counting the iterations. Once the iterations ended, I calculated the average length of all comments using the formula below:

Average = Sum of lengths of comments / Total number of comments

I ran it on a sample dataset first and then verified my answer with an online compiler and the answer was correct. Attached are the screenshots of the sample code, its output and the output of an online compiler:



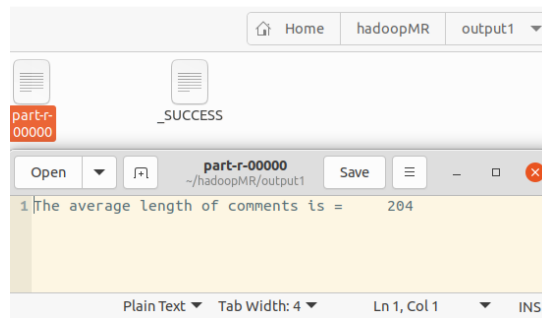*Fig 2.1: Contents of Sample.txt*

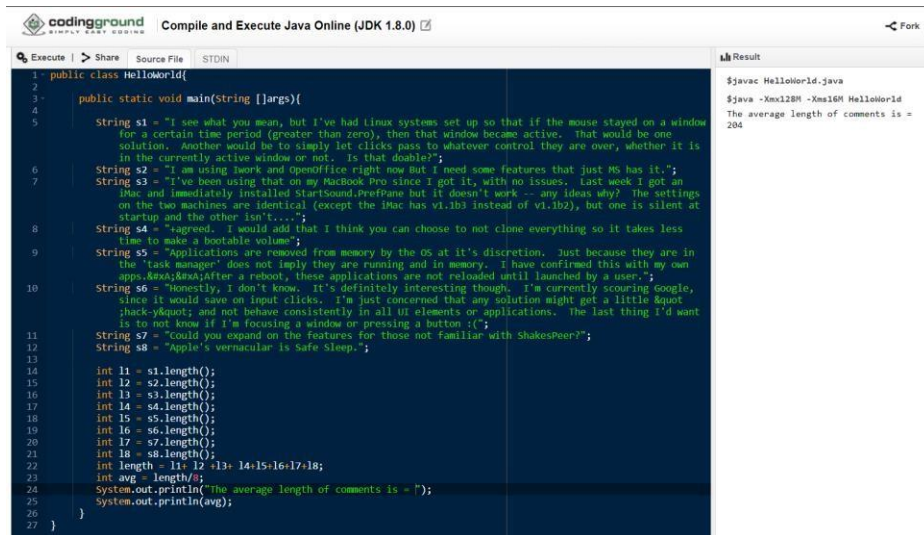*Fig 2.2: Output of Sample.txt*



*Fig 2.3: Output of Sample.txt on an online compiler*

After verifying my output, I ran the same code on the whole Comments.txt file. For that the code used and the outputs are given below:
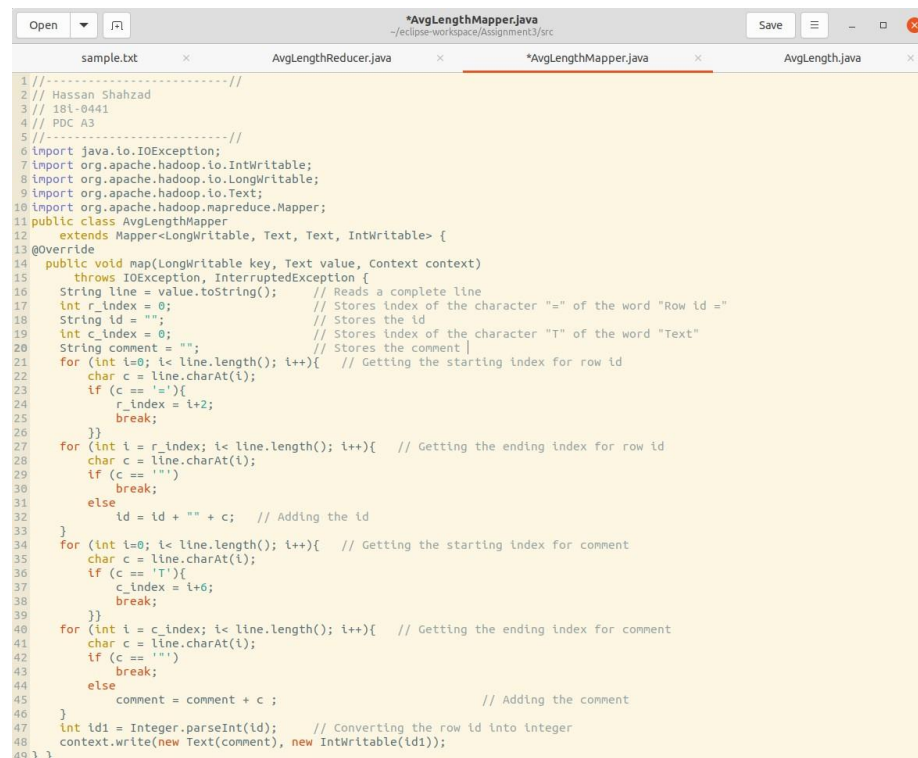


*Fig 2.4: Mapper Class*

```java
//------------------------//
// Hassan Shahzad
// 18i-0441
// PDC A3
//------------------------//

import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

public class AvgLengthReducer                      // last comment of the file
    extends Reducer<Text, IntWritable, Text, IntWritable> {
    int length = 0;      // Length of a comment
    int count = 0;       // Number of comments

    @Override
    public void cleanup(Context context) throws IOException, InterruptedException {    // If last comment has been read
        Text key = new Text();
        int avglen;                 // Average length of comments
        avglen = length/count;      // Calculating average
        key.set("The average length of comments is = ");
        context.write(key, new IntWritable(avglen));
    }


@Override
  public void reduce(Text key, Iterable<IntWritable> values, Context context)
      throws IOException, InterruptedException {

    String str = key.toString();    // Converting the comment into string
    length += str.length();         // Storing length of the comment
    count++;                        // Increasing count after each comment

  }
}
```

*Fig 2.5: Reducer Class*

```java
//------------------------//
// Hassan Shahzad
// 18i-0441
// PDC A3
//------------------------//

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.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class AvgLength {
  public static void main(String[] args) throws Exception {
    if (args.length != 2) {
      System.err.println("Usage: AvgTemperature <input path> <output path>");
      System.exit(-1);
    }

    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "Average Temprature");
    job.setJarByClass(AvgLength.class);
    job.setJobName("Average Length of Comments");
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path    (args[1]));
    job.setMapperClass(AvgLengthMapper.class);
    job.setReducerClass(AvgLengthReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);
    System.exit(job.waitForCompletion(true) ? 0 : 1);
  }
}
```

*Fig 2.6: Main Java File*

```
                Combine output records=0
                Reduce input groups=367977
                Reduce shuffle bytes=67142422
                Reduce input records=370371
                Reduce output records=1
                Spilled Records=740742
                Shuffled Maps =1
                Failed Shuffles=0
                Merged Map outputs=1
                GC time elapsed (ms)=111
                Total committed heap usage (bytes)=1715470336
        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=116132712
        File Output Format Counters
                Bytes Written=41
hxn@hxn:~$ hadoop fs -copyToLocal outputttt /home/hxn/hadoopMR/outputttt
hxn@hxn:~$
```

*Fig 2.7: Execution of Code via HDFS*

```
1 Hassan Shahzad
2 i180441
3 |
4 The average length of comments is =    172
```

*Fig 2.8: Output of File*

4