| Name of Student: Pushkar Sane | |
|---|---|
| Roll Number: 45 | Lab Assignment Number: 2 |
| Title of Lab Assignment: To study functionality of Mapreduce and implement following programs using Mapreduce.<br><br>1. Write a program in Map Reduce for WordCount operation.<br>2. Write a program in Map Reduce for Matrix Multiplication | |
| DOP: 30-08-2024 | DOS: 02-09-2024 |
| CO Mapped:<br>CO1 | PO Mapped:<br>PO1, PO2, PO3, PO4, PO5, PO7, PSO1, PSO2 | Signature: |

## Practical No. 2

**Aim: To study functionality of Mapreduce and implement following programs using Mapreduce.**

    1. **Write a program in Map Reduce for WordCount operation.**

    2. **Write a program in Map Reduce for Matrix Multiplication.**

**Theory:**

What is MapReduce?

MapReduce is a programming paradigm that enables massive scalability across hundreds or thousands of servers in a Hadoop cluster. As the processing component, MapReduce is the heart of Apache Hadoop. The term "MapReduce" refers to two separate and distinct tasks that Hadoop programs perform.

The first is the map job, which takes a set of data and converts it into another set of data, where individual elements are broken down into tuples (key/value pairs). The reduce job takes the output from a map as input and combines those data tuples into a smaller set of tuples.

As the sequence of the name MapReduce implies, the reduce job is always performed after the map job. MapReduce programming offers several benefits to help you gain valuable insights from your big data:

- Scalability: Businesses can process petabytes of data stored in the Hadoop Distributed File System (HDFS).
- Flexibility: Hadoop enables easier access to multiple sources of data and multiple types of data.
- Speed: With parallel processing and minimal data movement, Hadoop offers fast processing of massive amounts of data.
- Simple: Developers can write code in a choice of languages, including Java, C++ and Python.

**Program 1: MapReduce for Average Word Count**

**WordCount.java**

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;

```java
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 TokenizerMapper 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 itr = new StringTokenizer(value.toString());
                        while (itr.hasMoreTokens()) {
                                word.set(itr.nextToken());
                                context.write(word, one);
                        }
                }
        }


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

```
public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "word count");
        job.setJarByClass(WordCount.class);
        job.setMapperClass(TokenizerMapper.class);
        job.setCombinerClass(IntSumReducer.class);
        job.setReducerClass(IntSumReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}
```

**Output:**

**Program 2. To implement Matrix-Multiplication using MapReduce**

**Theory:**

To perform matrix multiplication, 2 sets of MapReduce are required,

1. MapReduce to multiply:
    a. Mapper: Form pairs of Row from matrix A and Column from Matrix B
    b. Reducer: Multiply the results of the above Mapper
2. MapReduce to Calculate:
    a. Mapper: Identity, returns the same Key value pair
    b. Reducer: Calculates the sum of all the same keys and returns the key with the sum.

The diagram below shows an example of Matrix Multiplication using MapReduce.

**Code:**

```
package expt4;

import java.io.DataInput;
import java.io.DataOutput;
import java.io.IOException;
import java.util.ArrayList;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.DoubleWritable;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.io.Writable;
import org.apache.hadoop.io.WritableComparable;
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.*;
import org.apache.hadoop.mapreduce.lib.output.*;
import org.apache.hadoop.util.ReflectionUtils;

class Element implements Writable {
        int tag;
        int index;
        double value;

        Element() {
                tag = 0;
                index = 0;
                value = 0.0;
        }

        Element(int tag, int index, double value) {
                this.tag = tag;
                this.index = index;
                this.value = value;
```

```
        }

        @Override
        public void readFields(DataInput input) throws IOException {
                tag = input.readInt();
                index = input.readInt();
                value = input.readDouble();
        }

        @Override
        public void write(DataOutput output) throws IOException {
                output.writeInt(tag);
                output.writeInt(index);
                output.writeDouble(value);
        }
}

class Pair implements WritableComparable<Pair> {

        int i;
        int j;

        Pair() {
                i = 0;
                j = 0;
        }

        Pair(int i, int j) {
                this.i = i;
                this.j = j;
        }

        @Override
        public void readFields(DataInput input) throws IOException {
                i = input.readInt();
                j = input.readInt();
```

```java
        }

        @Override
        public void write(DataOutput output) throws IOException {
                output.writeInt(i);
                output.writeInt(j);
        }

        @Override
        public int compareTo(Pair compare) {

                if (i > compare.i) {
                        return 1;
                } else if ( i < compare.i) {
                        return -1;
                } else {
                        if(j > compare.j) {
                                return 1;
                        } else if (j < compare.j) {
                                return -1;
                        }
                }
                return 0;
        }

        public String toString() {
                return i + " " + j + " ";
        }

}

public class MatrixMulitply {

        public static class MatriceMapperM extends
Mapper<Object,Text,IntWritable,Element> {
```

```java
        @Override
        public void map(Object key, Text value, Context context)
                    throws IOException, InterruptedException {
            String readLine = value.toString();
            String[] stringTokens = readLine.split(",");

            int index = Integer.parseInt(stringTokens[0]);
            double elementValue = Double.parseDouble(stringTokens[2]);

            Element e = new Element(0, index, elementValue);

            IntWritable keyValue = new
IntWritable(Integer.parseInt(stringTokens[1]));
            context.write(keyValue, e);
        }
    }

    public static class MatriceMapperN extends
Mapper<Object,Text,IntWritable,Element> {

        @Override
        public void map(Object key, Text value, Context context)
                    throws IOException, InterruptedException {
            String readLine = value.toString();
            String[] stringTokens = readLine.split(",");

            int index = Integer.parseInt(stringTokens[1]);
            double elementValue = Double.parseDouble(stringTokens[2]);

            Element e = new Element(1,index, elementValue);

            IntWritable keyValue = new
IntWritable(Integer.parseInt(stringTokens[0]));
            context.write(keyValue, e);
        }
    }
```

```java
        public static class ReducerMxN extends Reducer<IntWritable,Element, Pair,
DoubleWritable> {

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

                ArrayList<Element> M = new ArrayList<Element>();
                ArrayList<Element> N = new ArrayList<Element>();

                Configuration conf = context.getConfiguration();

                for(Element element : values) {

                        Element tempElement = ReflectionUtils.newInstance(Element.class,
conf);
                        ReflectionUtils.copy(conf, element, tempElement);

                        if (tempElement.tag == 0) {
                                M.add(tempElement);
                        } else if(tempElement.tag == 1) {
                                N.add(tempElement);
                        }
                }

                for(int i=0;i<M.size();i++) {
                        for(int j=0;j<N.size();j++) {

                                Pair p = new Pair(M.get(i).index,N.get(j).index);
                                double multiplyOutput = M.get(i).value * N.get(j).value;

                                context.write(p, new DoubleWritable(multiplyOutput));
                        }
                }
        }
```

```java
        }

        public static class MapMxN extends Mapper<Object, Text, Pair, DoubleWritable> {
            @Override
            public void map(Object key, Text value, Context context)
                        throws IOException, InterruptedException {

                    String readLine = value.toString();
                    String[] pairValue = readLine.split(" ");

                    Pair p = new
Pair(Integer.parseInt(pairValue[0]),Integer.parseInt(pairValue[1]));
                    DoubleWritable val = new
DoubleWritable(Double.parseDouble(pairValue[2]));

                    context.write(p, val);
            }
        }

        public static class ReduceMxN extends Reducer<Pair, DoubleWritable, Pair,
DoubleWritable> {
            @Override
            public void reduce(Pair key, Iterable<DoubleWritable> values, Context
context)
            throws IOException, InterruptedException {

                    double sum = 0.0;
                    for(DoubleWritable value : values) {
                        sum += value.get();
                    }

                    context.write(key, new DoubleWritable(sum));
            }
        }

        public static void main(String[] args) throws Exception {
```

```
Job job = Job.getInstance();
job.setJobName("MapIntermediate");
job.setJarByClass(MatrixMulitply.class);

MultipleInputs.addInputPath(job, new Path(args[0]), TextInputFormat.class,
MatriceMapperM.class);
MultipleInputs.addInputPath(job, new Path(args[1]), TextInputFormat.class,
MatriceMapperN.class);
job.setReducerClass(ReducerMxN.class);

job.setMapOutputKeyClass(IntWritable.class);
job.setMapOutputValueClass(Element.class);

job.setOutputKeyClass(Pair.class);
job.setOutputValueClass(DoubleWritable.class);

job.setOutputFormatClass(TextOutputFormat.class);

FileOutputFormat.setOutputPath(job, new Path(args[2]));

job.waitForCompletion(true);

Job job2 = Job.getInstance();
job2.setJobName("MapFinalOutput");
job2.setJarByClass(MatrixMulitply.class);

job2.setMapperClass(MapMxN.class);
job2.setReducerClass(ReduceMxN.class);

job2.setMapOutputKeyClass(Pair.class);
job2.setMapOutputValueClass(DoubleWritable.class);

job2.setOutputKeyClass(Pair.class);
job2.setOutputValueClass(DoubleWritable.class);
job2.setInputFormatClass(TextInputFormat.class);
job2.setOutputFormatClass(TextOutputFormat.class);
```
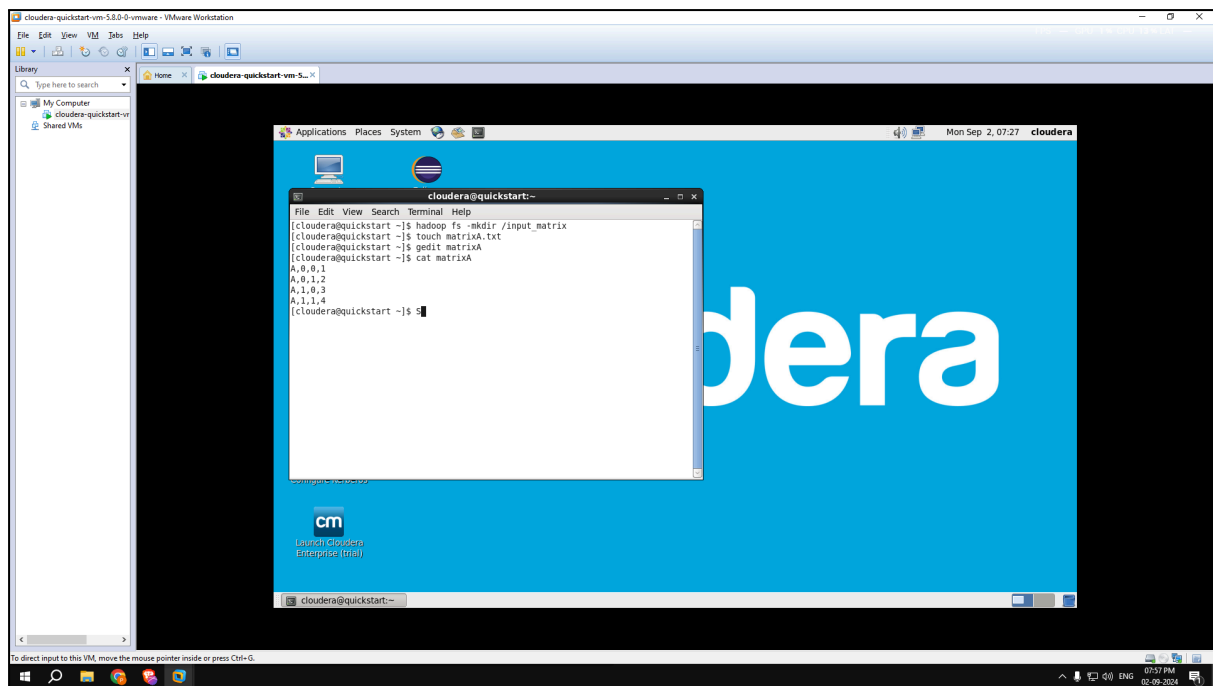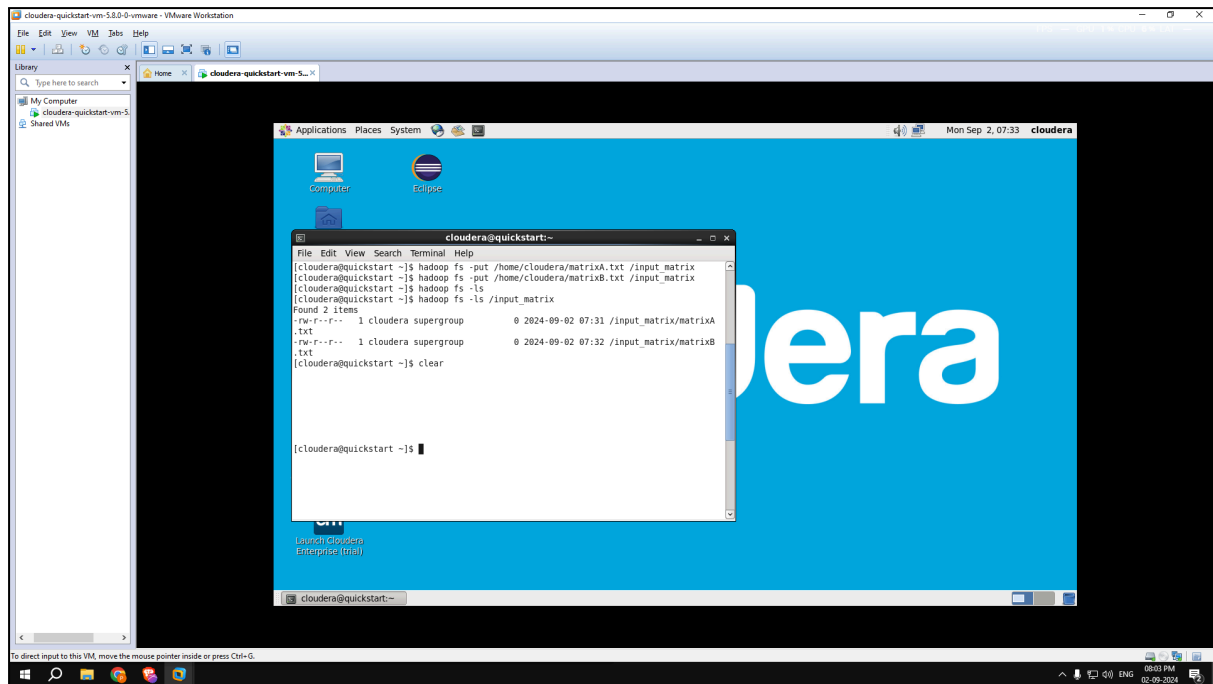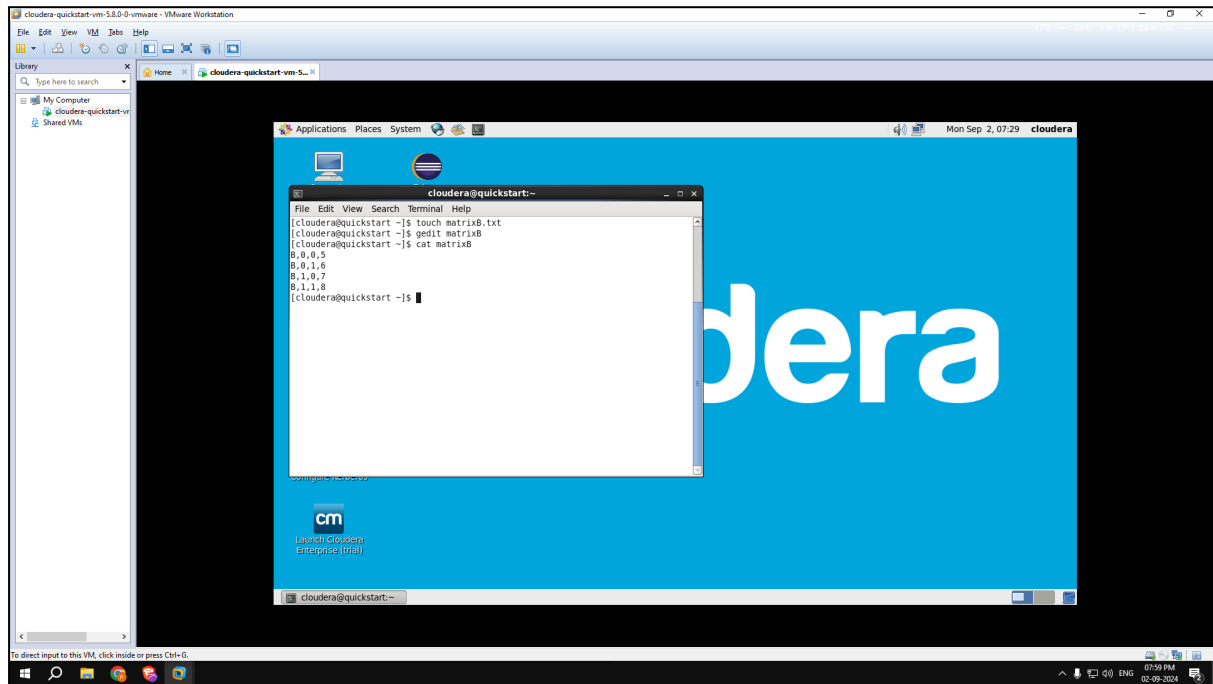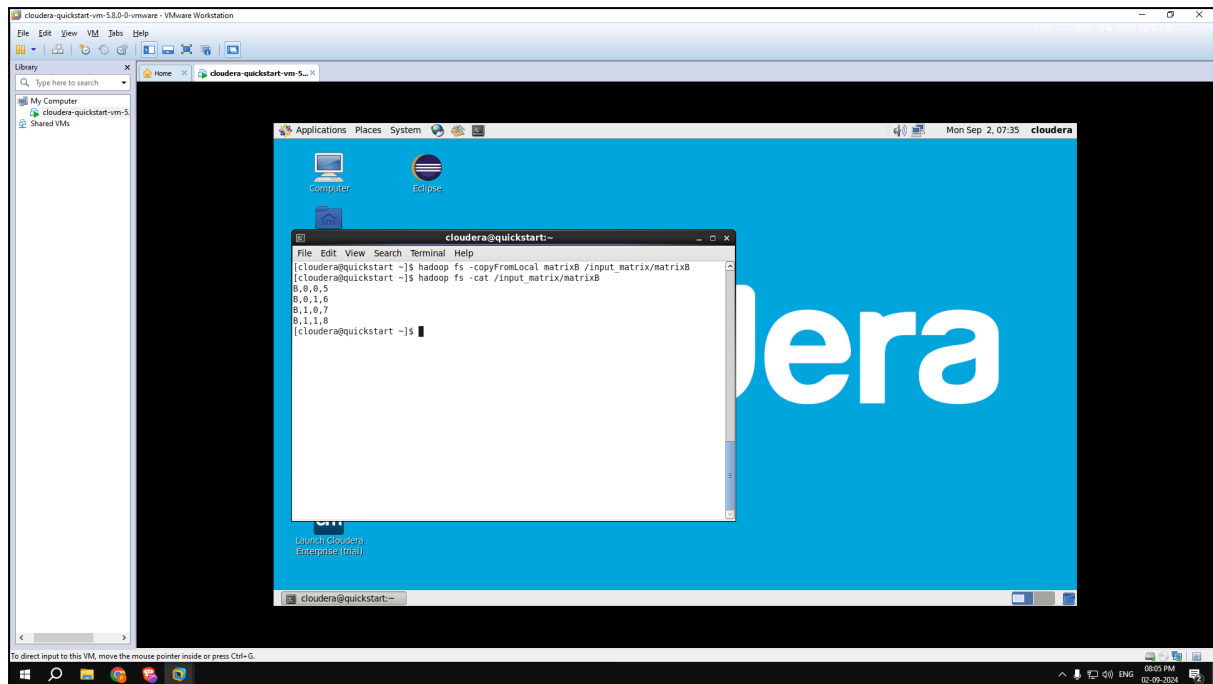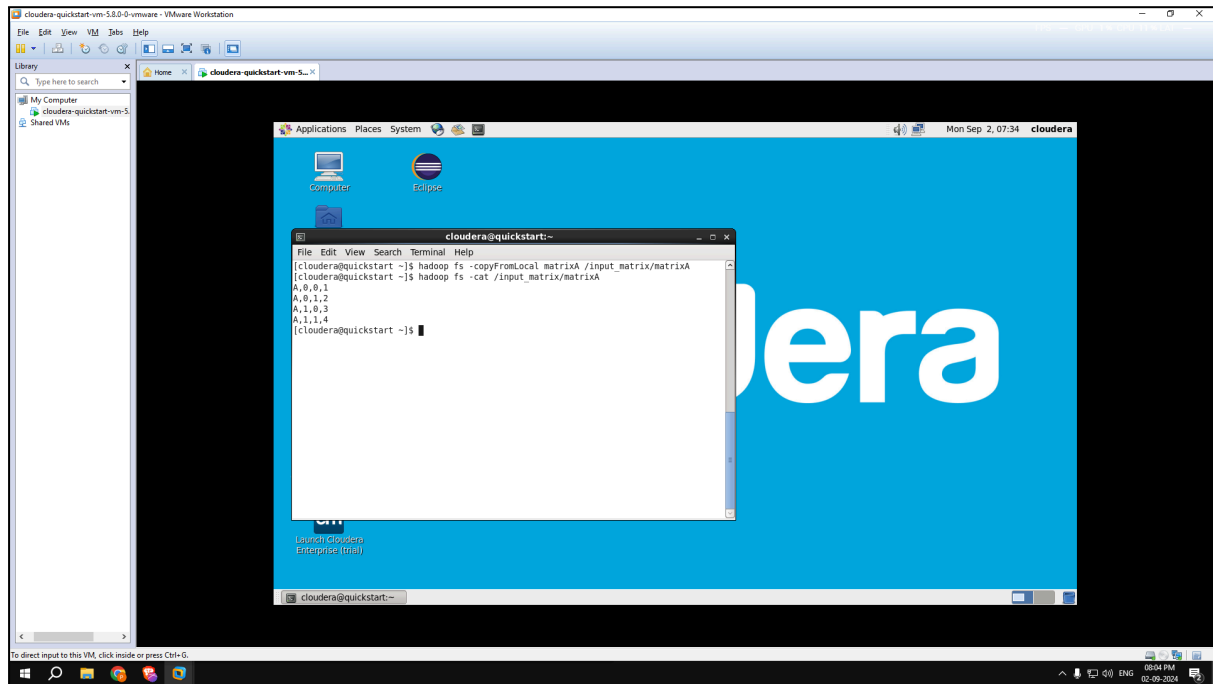
```
                FileInputFormat.setInputPaths(job2, new Path(args[2]));

                FileOutputFormat.setOutputPath(job2, new Path(args[3]));


                job2.waitForCompletion(true);
        }
}
```
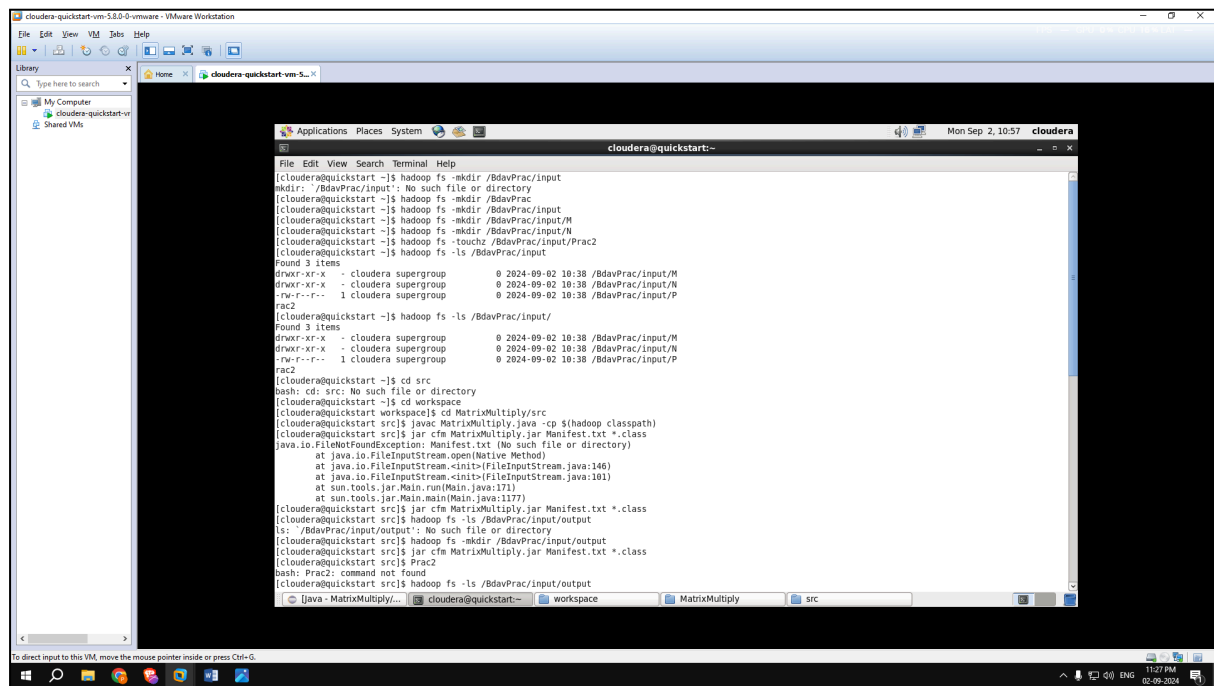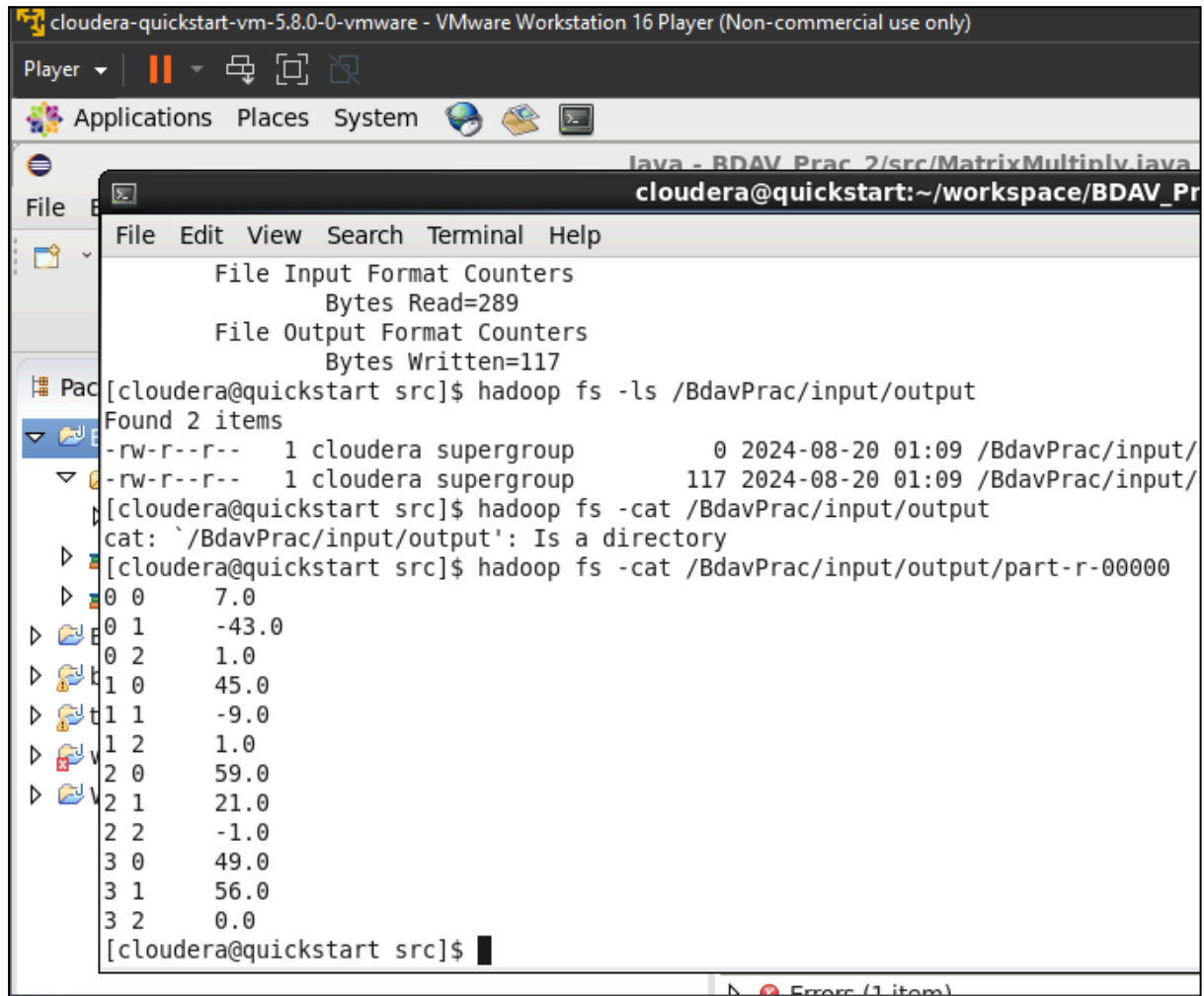
**Output:**

hadoop fs –mkdir /input_matrix
[cloudera@quickstart ~]$ hadoop fs -put /home/cloudera/matrixA.txt /input_matrix
[cloudera@quickstart ~]$ hadoop fs -put /home/cloudera/matrixB.txt /input_matrix
[cloudera@quickstart ~]$ hadoop fs -copyFromLocal matrixA /input_matrix/matrixA
[cloudera@quickstart ~]$ hadoop fs -cat /input_matrix/matrixA
[cloudera@quickstart ~]$ hadoop fs -copyFromLocal matrixB /input_matrix/matrixB
[cloudera@quickstart ~]$ hadoop fs -cat /input_matrix/matrixB
[cloudera@quickstart ~]$ hadoop jar MatrixMulti.jar OneStepMatrixMultiplication
/input_matrix/* /output
[cloudera@quickstart ~]$ hadoop fs -cat /output/*

**Conclusion:** Successfully implemented Word Count and Matrix Multiplication in hadoop using MapReduce.