

Large Scale Data Processing

WINTER SEMESTER 2017-18

INSTRUCTOR: N. MAHESWARI

CHALLA DINESH REDDY
15BCE1274

IAB PROGRAMS

Contents:-

- Lab1 Introduction and basic commands in Hadoop.
- Lab2 basic wordcount Mapreduce program
- Lab3 part1 Find maximum occurring word
- Lab3 part2 total word count
- Lab4 part1 Combiner
- Lab4 part2 Check for certain range of value and output it.(salary>1lakh)
- Lab5 custom partitioning
- Lab6 Working with different i/o formats.
- Lab7 part1 Find top k records based on a value(Tree Map)
- Lab7 part2 User defined Counters
- Lab8 Side Data
- Lab9 part 1 Distributed Cache
- Lab9 part 2 Reduce side join Customer details and transaction details.
- Lab 10 Secondary sorting

Lab 1:- basic hadoop commands.

/////hadoop commands

Start-all.sh:- starts all the demons

Jps:- shows the jobs or daemons running along with the ports.

hadoop dfs -ls <path> :- lists the files in hdfs current

hadoop dfs -put <local> <dst> :- puts the data file from local machine to hdfs

hadoop dfs -cat :- shows the file in terminal.

hadoop dfs -get :- get files from hdfs to local machine

hadoop dfs -rm :- removes files in hdfs

hadoop dfs -fsck :- Runs a HDFS file system checking utility.

hadoop jar :- run Mapreduce Programming

hadoop distcp <src> <dst> :- copy from hdfs to hdfs

///web ui localhost:<port>

default ports:-

Namenode 50070

Datanodes 50075

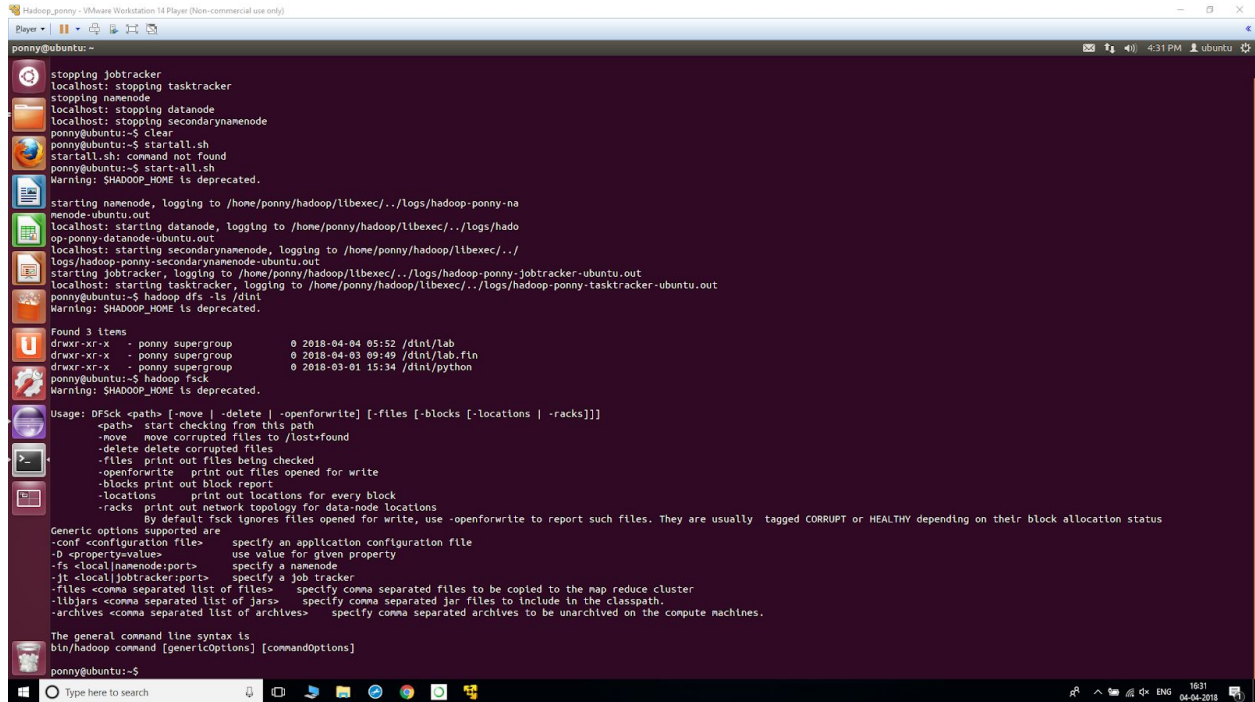
Secondarynamenode 50090

Backup/Checkpoint node 50105

Jobtracker 50030

Tasktrackers 50060

Output:-



```
stopping jobtracker
localhost: stopping tasktracker
stopping namenode
localhost: stopping datanode
localhost: stopping secondarynamenode
ponny@ubuntu:~$ clear
ponny@ubuntu:~$ startall.sh
startall.sh: command not found
ponny@ubuntu:~$ start-all.sh
Warning: SHADOOP_HOME is deprecated.

starting namenode, logging to /home/ponny/hadoop/libexec/../logs/hadoop-ponny-na
menode-ubuntu.out
localhost: starting datanode, logging to /home/ponny/hadoop/libexec/../logs/hado
op-ponny-datanode-ubuntu.out
localhost: starting secondarynamenode, logging to /home/ponny/hadoop/libexec/../
logs/hadoop-ponny-secondarynamenode-ubuntu.out
starting jobtracker, logging to /home/ponny/hadoop/libexec/../logs/hadoop-ponny-jobtracker-ubuntu.out
localhost: starting tasktracker, logging to /home/ponny/hadoop/libexec/../logs/hadoop-ponny-tasktracker-ubuntu.out
ponny@ubuntu:~$ hadoop dfs -ls /dint
Warning: SHADOOP_HOME is deprecated.

Found 3 items
drwxr-xr-x  - ponny supergroup          0 2018-04-04 05:52 /dint/lab
drwxr-xr-x  - ponny supergroup          0 2018-04-03 09:49 /dint/lab.fln
drwxr-xr-x  - ponny supergroup          0 2018-03-01 15:34 /dint/python
ponny@ubuntu:~$ hadoop fsck
Warning: SHADOOP_HOME is deprecated.

Usage: DFSck <paths> [-move | -delete | -openforwrite] [-files [-blocks [-locations | -racks]]]
<paths> start checking from this path
-move move corrupted files to /lost+found
-delete delete corrupted files
-files print out files being checked
-openforwrite print out files opened for write
-blocks print out block report
-locations print out locations for every block
-racks print out network topology for data-node locations
By default fsck ignores files opened for write, use -openforwrite to report such files. They are usually tagged CORRUPT or HEALTHY depending on their block allocation status

Generic options supported are
-conf <configuration file> specify an application configuration file
-D <property=value> use value for given property
-fs <local|namenode:port> specify a namenode
-jt <local|jobtracker:port> specify a job tracker
-files <comma separated list of files> specify comma separated files to be copied to the map reduce cluster
-libjars <comma separated list of jars> specify comma separated jar files to include in the classpath.
-archives <comma separated list of archives> specify comma separated archives to be unarchived on the compute machines.

The general command line syntax is
bin/hadoop command [genericOptions] [commandOptions]
```

Lab2:- First Map Reduce Program. Word count.

Algorithm:

Map Function – It takes a set of data and converts it into another set of data, where individual elements are broken down into tuples (Key-Value pair).

Example – (Map function in Word Count)

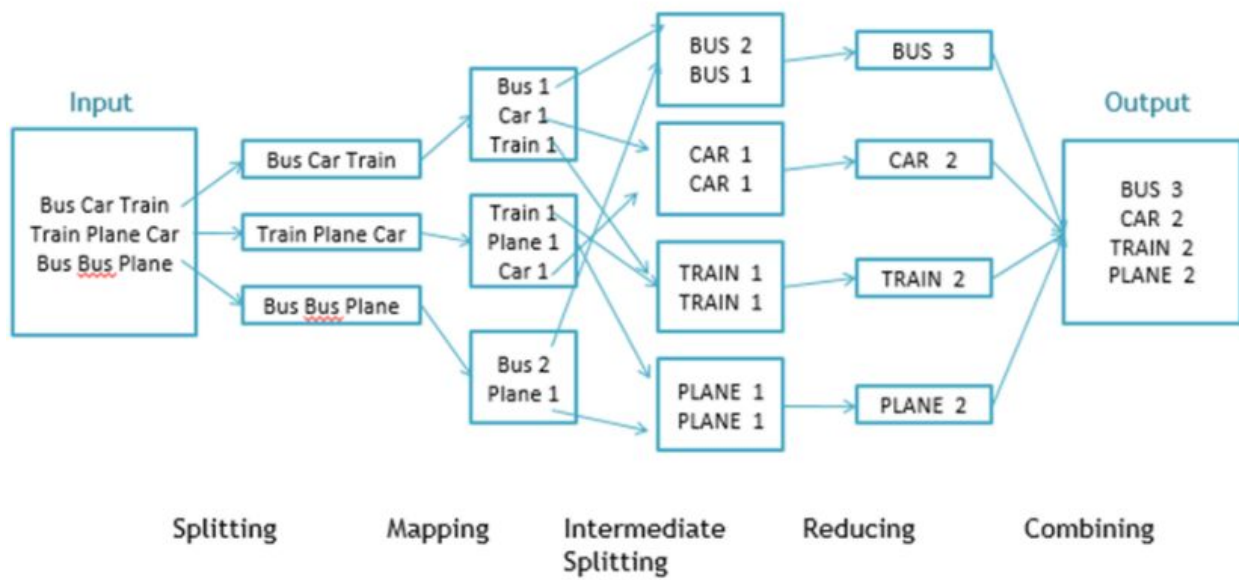
Input	Set of data	Bus, Car, bus, car, train, car, bus, car, train, bus, TRAIN,BUS, buS, caR, CAR, car, BUS, TRAIN
Output	Convert into another set of data (Key,Value)	(Bus,1), (Car,1), (bus,1), (car,1), (train,1), (car,1), (bus,1), (car,1), (train,1), (bus,1), (TRAIN,1),(BUS,1), (buS,1), (caR,1), (CAR,1), (car,1), (BUS,1), (TRAIN,1)

Reduce Function:-Takes the output from Map as an input and combines those data tuples into a smaller set of tuples.

Example – (Reduce function in Word Count)

Input (output of Map function)	Set of Tuples	(Bus,1), (Car,1), (bus,1), (car,1), (train,1), (car,1), (bus,1), (car,1), (train,1), (bus,1), (TRAIN,1),(BUS,1), (buS,1), (caR,1), (CAR,1), (car,1), (BUS,1), (TRAIN,1)
Output	Converts into smaller set of tuples	(BUS,7), (CAR,7), (TRAIN,4)

Overall workflow of Mapreduce program:-



Java code:-

```
//basic 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);

            }

        }

    }

}

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);

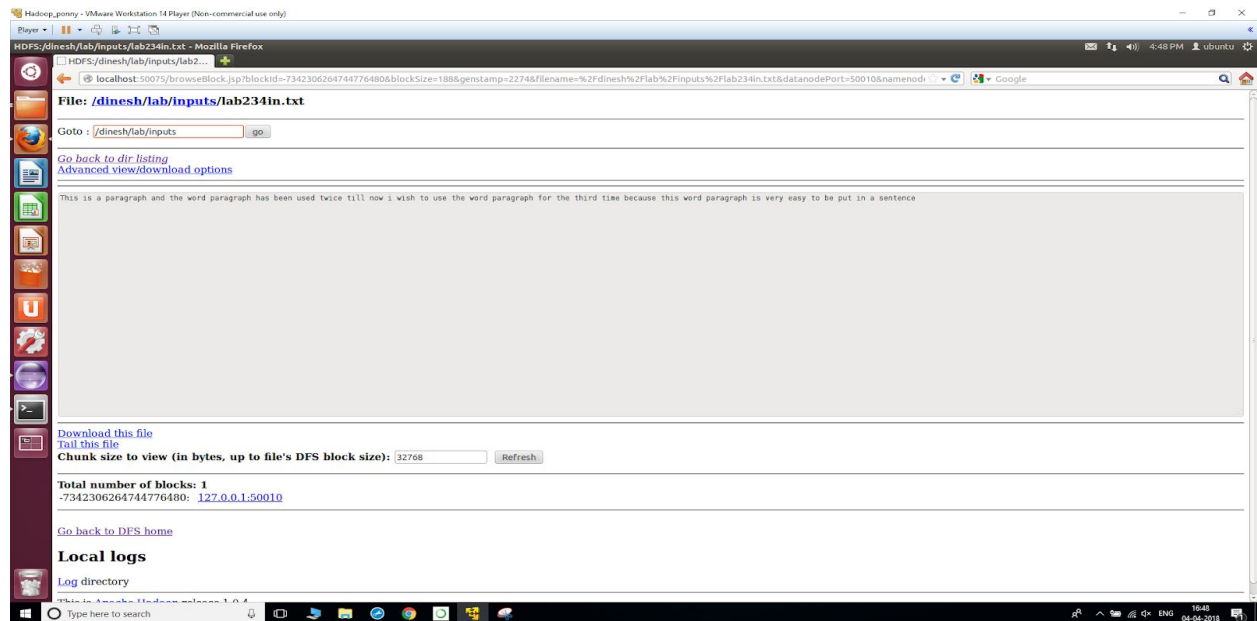
}

}
```

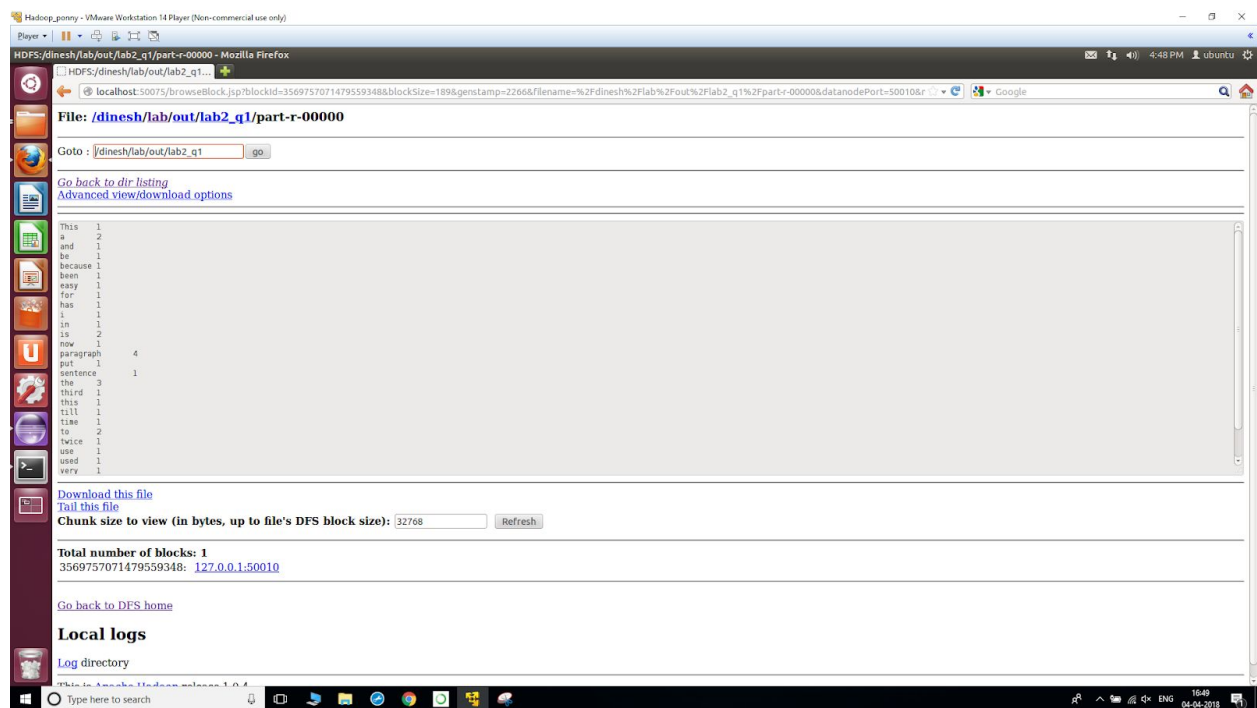
Execution:-

hadoop jar < jar file path <class name> <input path in hdfs> <output path in hdfs>

Input:-



Output:-



Lab3:-

q1) Find the maximum occurring word and write the no of time it occurred.

Algorithm:-

In reducer have a separate variable max which is initialised to 0 and it keeps on checking for all words, replace it the count greater than previous and finally use cleanup to write the output to hdfs.

Code:-

```
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 MaxCount{

    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);

}

}

}

}

public static class Reduce extends Reducer<Text, IntWritable, Text, IntWritable>{

    int maxSum = 0;

    Text maxOccuredKey = new Text();

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

        int sum = 0;

        for(IntWritable val: values){

            sum += val.get();

        }

        if(sum > maxSum){

            maxSum = sum;

            maxOccuredKey.set(key);

        }

    }

}
```

```
        @Override

        protected void cleanup(Context context) throws IOException, InterruptedException{

            context.write(maxOccuredKey, new IntWritable(maxSum));

        }

    }

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

        Configuration conf = new Configuration();

        Job job = new Job(conf,"maxcount");

        job.setJarByClass(MaxCount.class);

        job.setMapperClass(Map.class);

        job.setOutputKeyClass(Text.class);

        job.setOutputValueClass(IntWritable.class);

        job.setReducerClass(Reduce.class);

        FileInputFormat.addInputPath(job, new Path(args[0]));

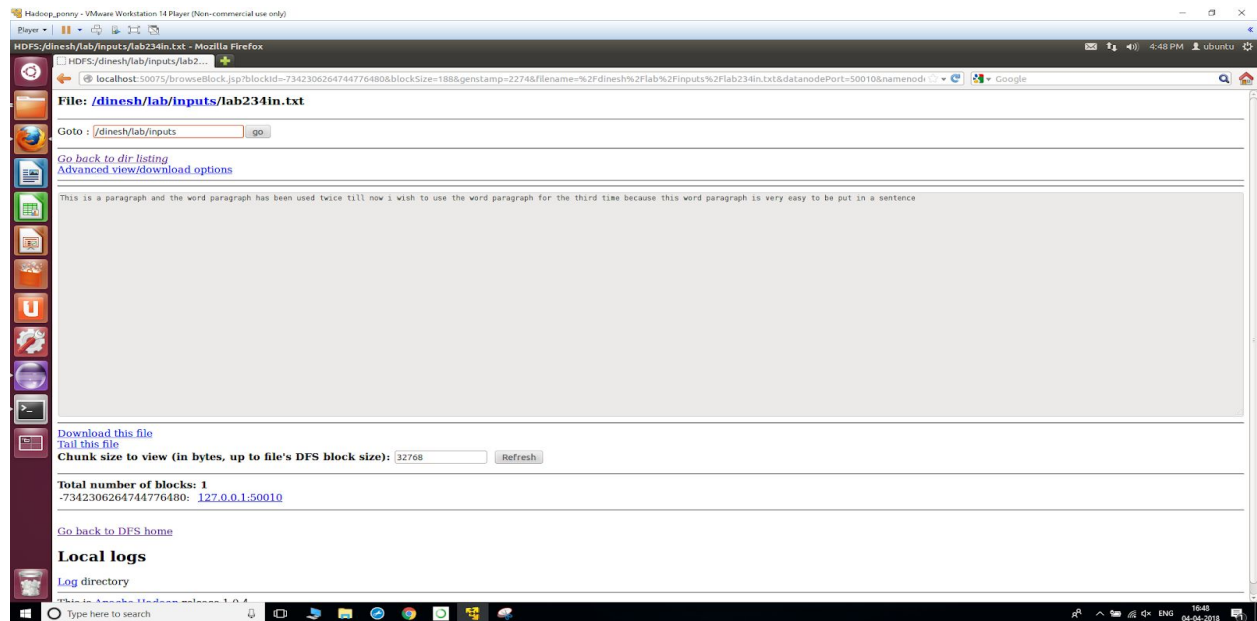
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.waitForCompletion(true);

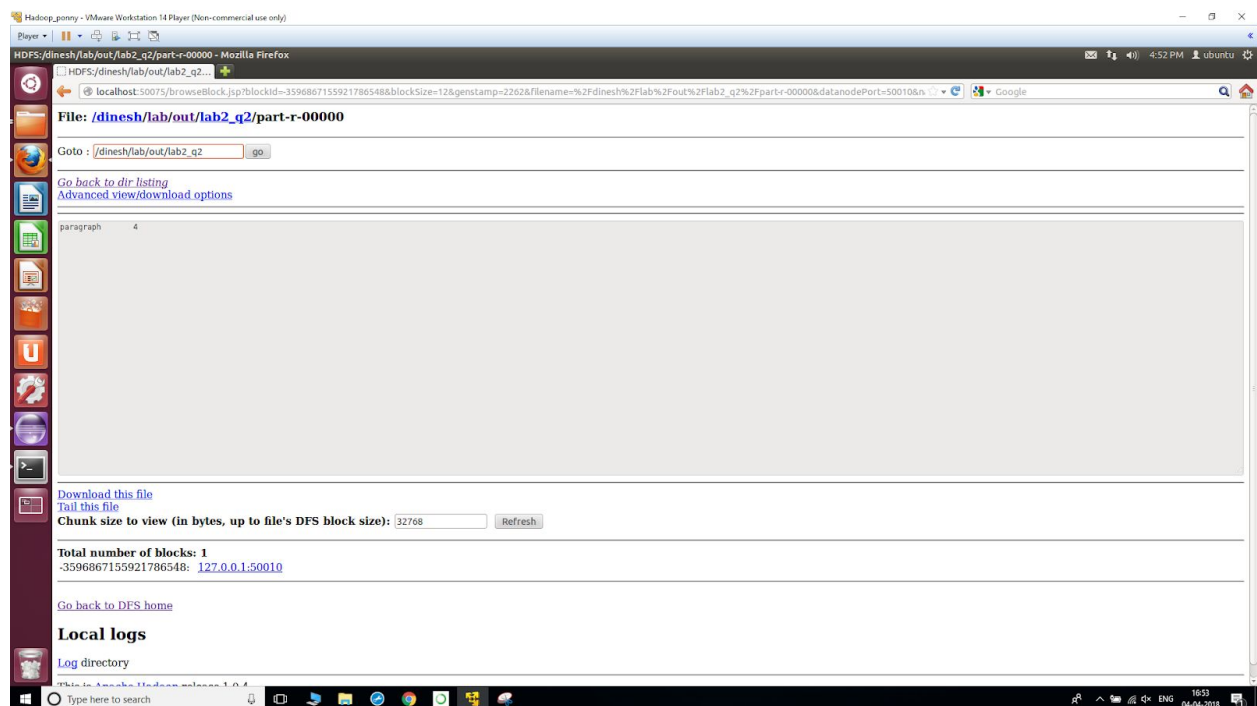
    }

}
```

Input:- Same as that of lab2.



Output:-



q2) Find the total number of words. Total word count.

Algorithm:- In reducer, add up the individual counts which are in the value of key value pair.

Code:-

```
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 totalCount{

    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);

    }

}

}

public static class Reduce extends Reducer<Text, IntWritable, Text, IntWritable>{

    int maxSum = 0;

    Text totalKey = new Text();

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

        int sum = 0;

        for(IntWritable val: values){

            sum += val.get();

        }

        totalSum += sum;

    }

    @Override

    protected void cleanup(Context context) throws IOException, InterruptedException{

        context.write(totalKey, new IntWritable(totalSum));

    }

}
```

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

    Configuration conf = new Configuration();

    Job job = new Job(conf,"totalcount");

    job.setJarByClass(totalCount.class);

    job.setMapperClass(Map.class);

    job.setOutputKeyClass(Text.class);

    job.setOutputValueClass(IntWritable.class);

    job.setReducerClass(Reduce.class);

    FileInputFormat.addInputPath(job, new Path(args[0]));

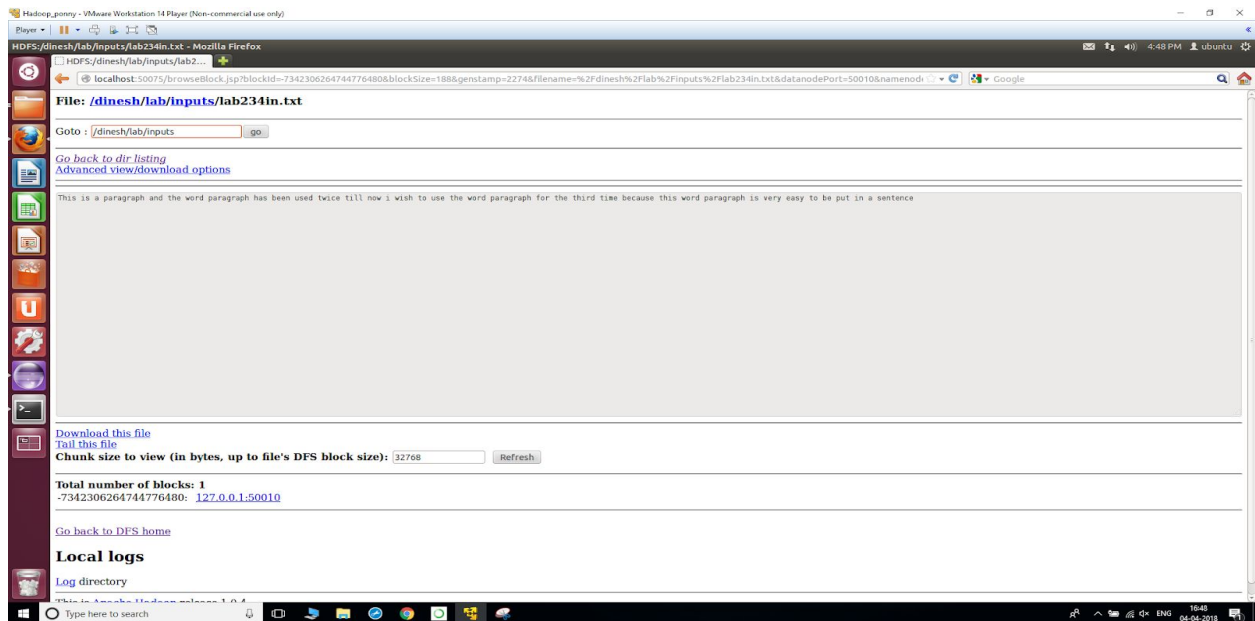
    FileOutputFormat.setOutputPath(job, new Path(args[1]));

    job.waitForCompletion(true);

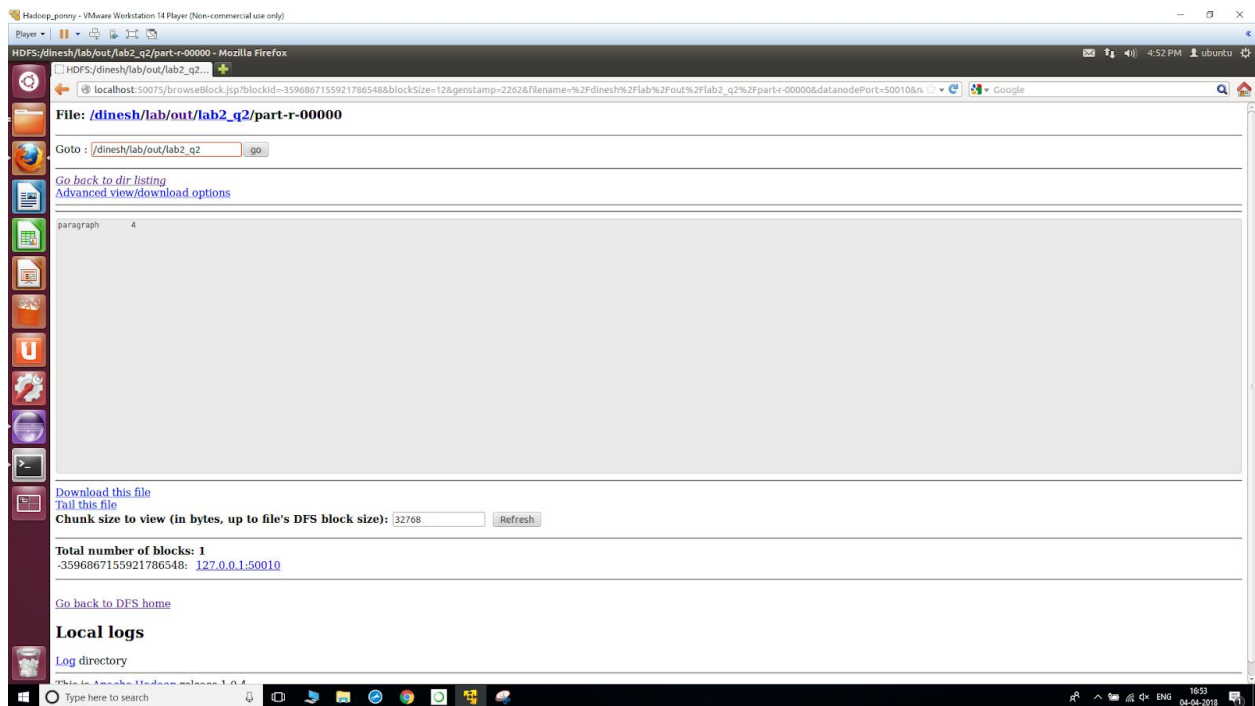
}

}
```

Input:-same as lab2.



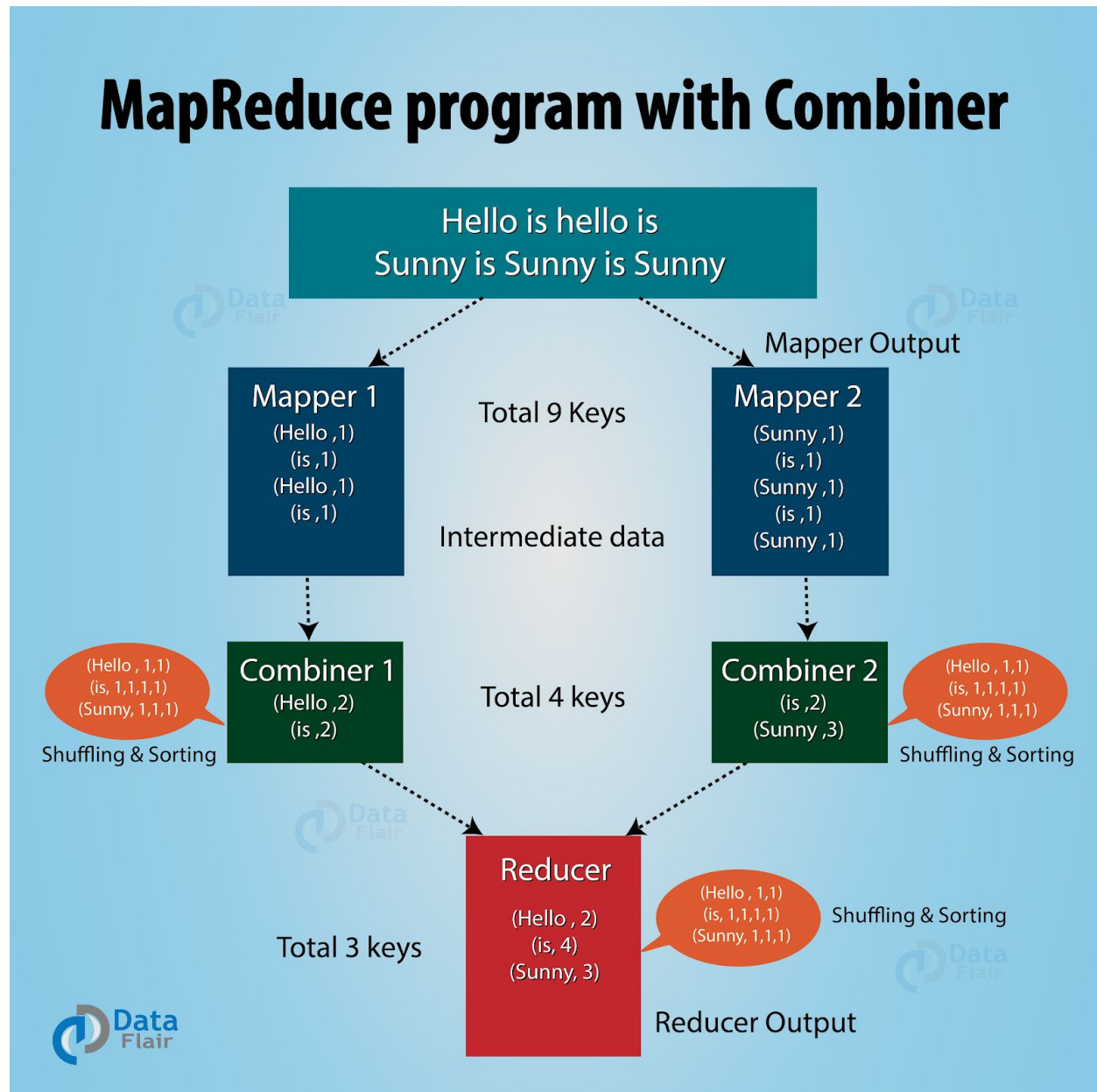
Output:-



lab4:-COMBINER

Algorithm:-

The combiner in MapReduce is also known as 'Mini-reducer'. The primary job of Combiner is to process the output data from the Mapper, before passing it to Reducer. It runs after the mapper and before the Reducer and its use is optional.



Code:-

*****lab4-combiner

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 WordCountCombiner{

 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);

    }

}

}

}

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(WordCountCombiner.class);

    job.setMapperClass(Map.class);

    job.setOutputKeyClass(Text.class);
```

```
        job.setOuttputValueClass(IntWritable.class);

        job.setReducerClass(Reduce.class);

        job.setCombinerClass(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);

    }

}

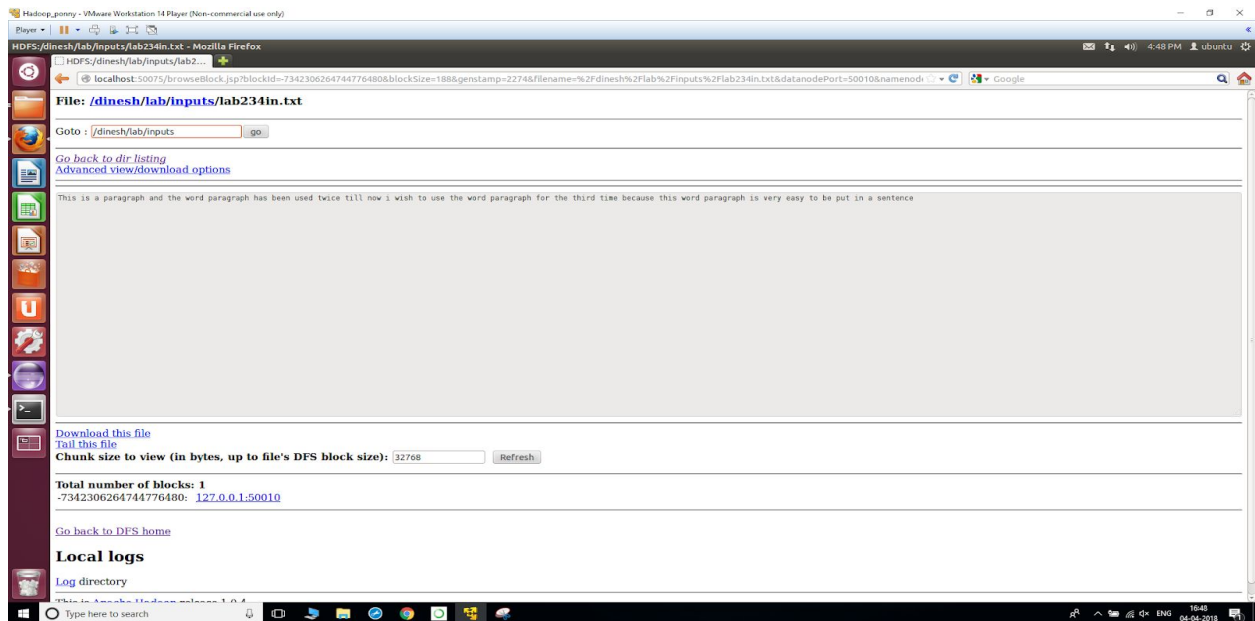
/*
***

to execute:-

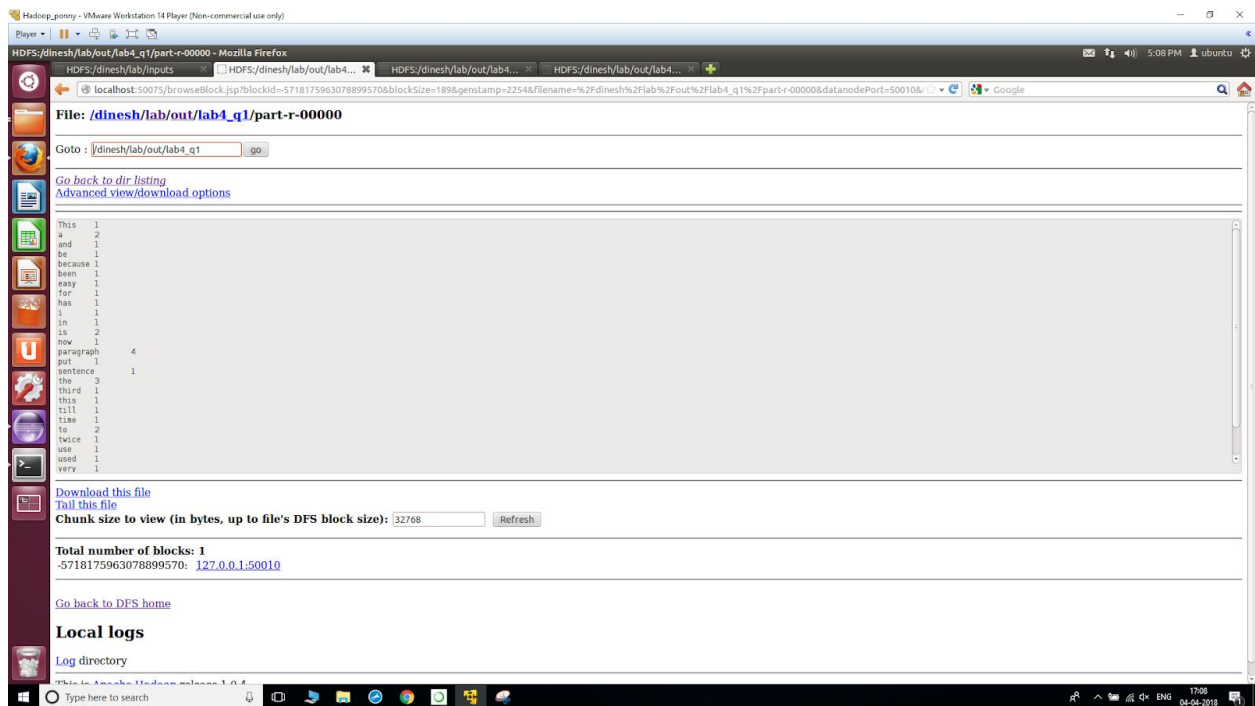
hadoop jar /home/ponny/*jar WordCountCombiner.java /user/testWords.txt /user/dini/out4

*/
```

Input:-
Same as lab 2



Output:-



q2) Write output of only a certain range.

Check for salary greater than 1 lakh and write if name if salary greater than 1 lakh.

Code:-

```
//if salary is more than 1 lakh write the data to output file else ignore.
```

```
import java.io.IOException;
```

```
import org.apache.hadoop.conf.*;
```

```
import org.apache.hadoop.fs.Path;
```

```
import org.apache.hadoop.io.*;
```

```
import org.apache.hadoop.mapreduce.*;
```

```
import org.apache.hadoop.mapreduce.lib.output.*;
```

```
import org.apache.hadoop.mapreduce.lib.input.*;
```

```
public class Salary {
```

```
    public static void main(String[] args) throws IOException, ClassNotFoundException, InterruptedException {
```

```
        Configuration conf = new Configuration();
```

```
        Job job = new Job(conf, "wordcount");
```

```
        job.setJarByClass(Salary.class);
```

```
        job.setOutputKeyClass(Text.class);
```

```
        job.setOutputValueClass(IntWritable.class);
```

```
        job.setMapperClass(Map.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);
```

```
}
```

```
public static class Map extends Mapper<LongWritable, Text, Text, IntWritable>{
```

```
    int count = 0;
```

```
    public Text word = new Text("Total Number of Employees: ");
```

```
    public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
```

```
        String[] line = value.toString().split(",");
```

```
        int i = Integer.parseInt(line[1]);
```

```
        final IntWritable sal = new IntWritable(i);
```

```
        if(i > 100000) {
```

```
            context.write(new Text(line[3]+", "+line[0]),sal);
```

```
            //sal as value and rest two as paired key.
```

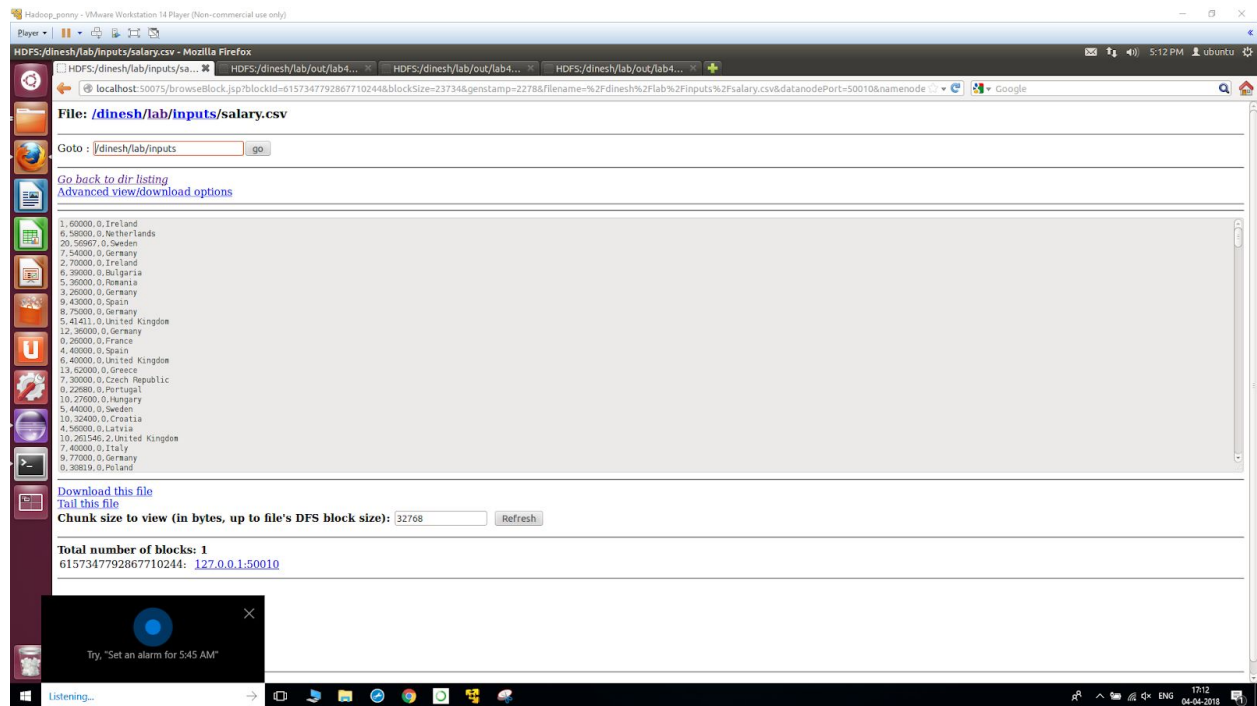
```
            count++;
```

```
        }
```

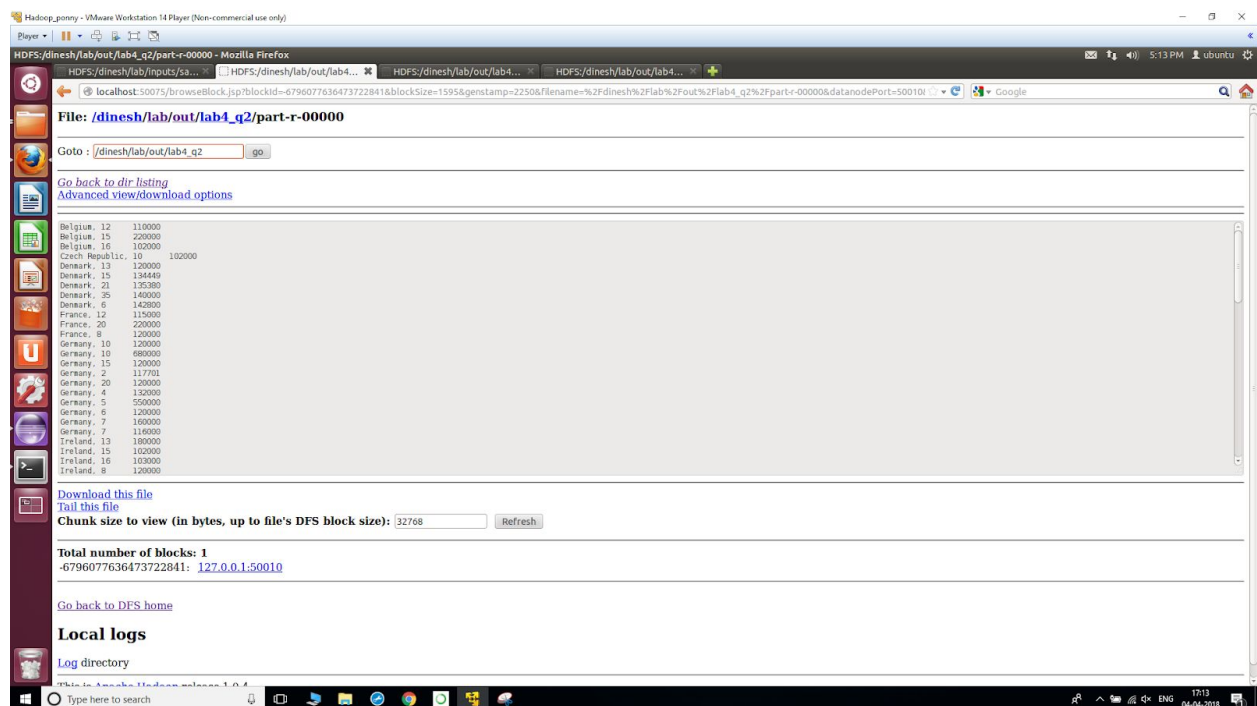
```
    }
```

```
        public void cleanup(Context context) throws IOException, InterruptedException {  
            context.write(word, new IntWritable(count));  
        }  
    }  
}
```

Input:-salary.csv



Output:-



lab5:-custom partitioning

Output data in different files based on a criterion.

Ex:-

salary less than 30000 in 1 file

Salary less than 50000 in another

And rest in another file. Output.

Code:-

```
import java.io.*;

import java.util.*;

import org.apache.hadoop.fs.Path;

import org.apache.hadoop.io.*;

import org.apache.hadoop.conf.*;

import org.apache.hadoop.mapreduce.*;

import org.apache.hadoop.mapreduce.Reducer.Context;

import org.apache.hadoop.mapreduce.lib.input.*;

import org.apache.hadoop.mapreduce.lib.output.*;

public class customP{

    public static class Map extends Mapper<LongWritable,Text,Text,IntWritable>{

        public void map(LongWritable key,Text value,Context context)throws IOException,InterruptedException{

            String[] line=value.toString().split(",");

            int i=Integer.parseInt(line[1]);

            context.write(new Text(line[3]), new IntWritable(i));
```

```
    }

}

public static class dpart extends Partitioner<Text,IntWritable>{

    public int getPartition(Text key,IntWritable value,int nr){

        if(value.get()<30000)

            return 0;

        if(value.get() < 50000)

            return 1;

        else

            return 2;

    }

}

public static class Reduce extends Reducer<Text,IntWritable,Text,IntWritable>{

    public void reduce(Text key,IntWritable value,Context context)throws IOException,InterruptedException{

        context.write(key,value);

    }

}

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

    Configuration conf=new Configuration();

    Job job=new Job(conf,"customP");

    job.setJarByClass(customP.class);
```

```
        job.setOutputKeyClass(Text.class);

        job.setOutputValueClass(IntWritable.class);

        job.setMapperClass(Map.class);

        job.setPartitionerClass(dpart.class);

        job.setNumReduceTasks(3);

        job.setInputFormatClass(TextInputFormat.class);

        job.setOutputFormatClass(TextOutputFormat.class);

        FileInputFormat.addInputPath(job,new Path(args[0]));

        FileOutputFormat.setOutputPath(job,new Path(args[1]));

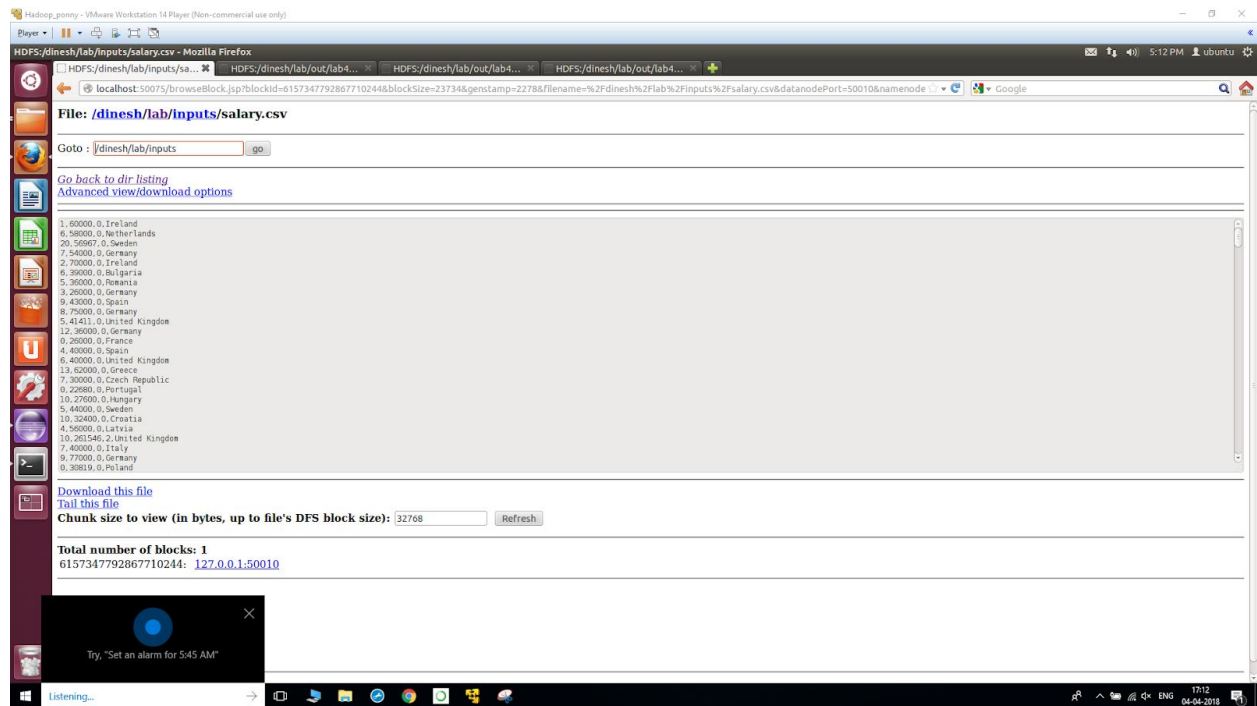
        job.waitForCompletion(true);

    }

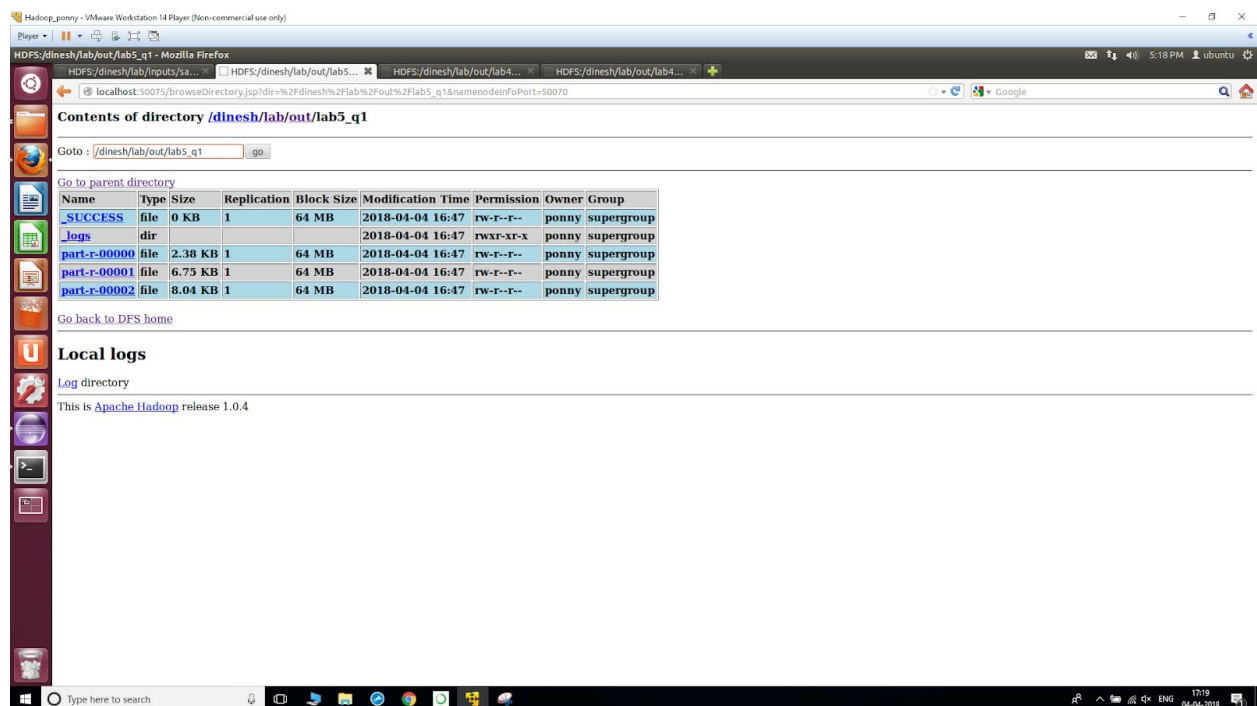
}
```

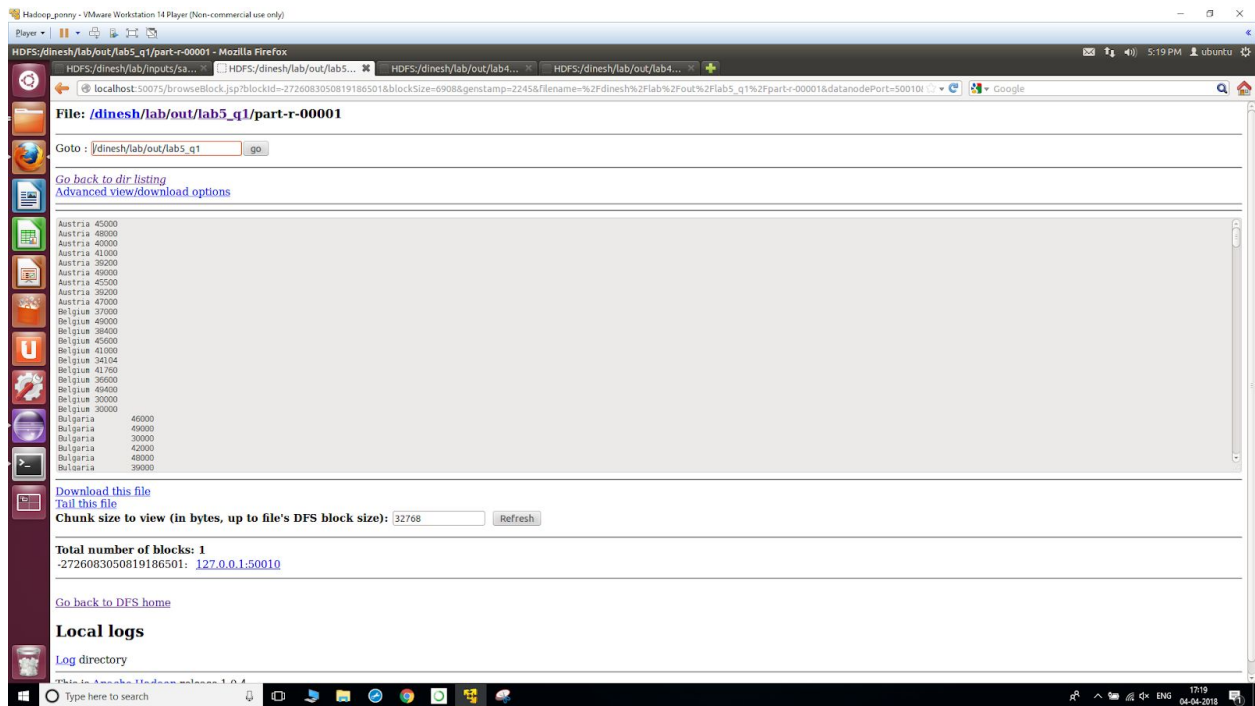
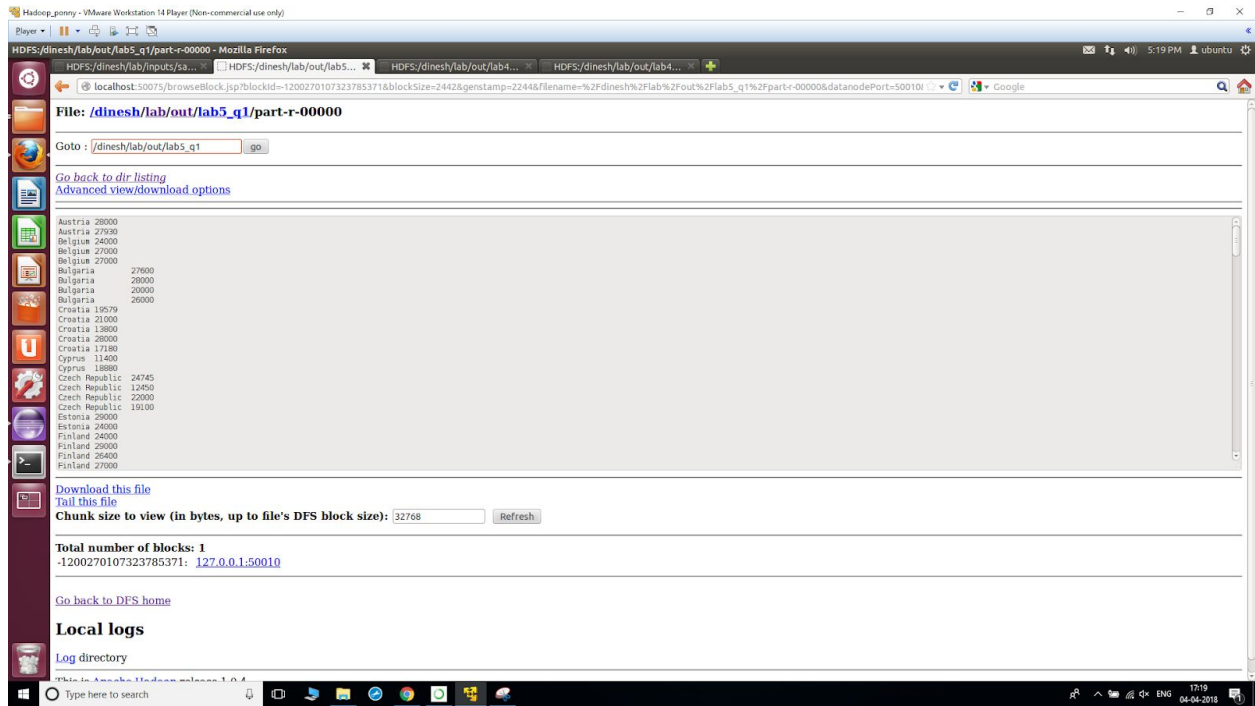
//results in 3 different output files.

Input:-same salary.csv file.



Output:-





q2) Specify the data that is to be assigned to each mapper by using splittble function.

Only the number of mappers will change, Output will be the same

Code:

results in 3 different output files.

*****splittable

```
import java.io.*;
```

```
import java.util.*;
```

```
import org.apache.hadoop.fs.Path;
```

```
import org.apache.hadoop.io.*;
```

```
import org.apache.hadoop.conf.*;
```

```
import org.apache.hadoop.mapreduce.*;
```

```
import org.apache.hadoop.mapreduce.Reducer.Context;
```

```
import org.apache.hadoop.mapreduce.lib.input.*;
```

```
import org.apache.hadoop.mapreduce.lib.output.*;
```

```
public class splittable{
```

```
    public static class Map extends Mapper<LongWritable,Text,Text,IntWritable>{
```

```
        public void map(LongWritable key,Text value,Context context)throws IOException,InterruptedException{
```

```
            String[] line=value.toString().split(",");
```

```
            int i=Integer.parseInt(line[1]);
```

```
            context.write(new Text(line[3]), new IntWritable(i));
```

```
        }
```

```
    }
```

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

    Configuration conf=new Configuration();

    //setting the split size before the job set as this is the configuration section.

    conf.set("mapred.max.split.size","10000");

    Job job=new Job(conf,"splitable");

    job.setJarByClass(splitable.class);

    job.setOutputKeyClass(Text.class);

    job.setOutputValueClass(IntWritable.class);

    job.setMapperClass(Map.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);

}

}
```

//open job tracker or task tracker to see the way job is done job tracker:localhost:50030 and task tracker :50060

Code for no split, i.e. All data to be processed in a single mapper.

```
*****split false

import java.io.*;

import org.apache.hadoop.fs.Path;

import org.apache.hadoop.io.*;

import org.apache.hadoop.conf.*;

import org.apache.hadoop.mapreduce.*;

import org.apache.hadoop.mapreduce.lib.input.*;

import org.apache.hadoop.mapreduce.lib.output.*;

public class split{

    public static class Map extends Mapper<LongWritable,Text,Text,IntWritable>{

        public void map(LongWritable key,Text value,Context context)throws IOException,InterruptedException{

            String[] line=value.toString().split(",");

            int i=Integer.parseInt(line[1]);

            context.write(new Text(line[3]), new IntWritable(i));

        }

    }

}

//all the input file will be done in a single map task

public class splitfalse extends TextInputFormat{

    protected boolean isSplittable(JobContext context, Path filename) {

        return false;

    }

}
```

```
    }  
}  
  
public static void main(String[] args) throws Exception  
{  
    Configuration conf=new Configuration();  
    Job job=new Job(conf,"split");  
    job.setInputFormatClass(splitfalse.class);  
    job.setJarByClass(split.class);  
    job.setOutputKeyClass(Text.class);  
    job.setOutputValueClass(IntWritable.class);  
    job.setMapperClass(Map.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);  
}  
  
}
```

Lab6:-

Different input and output formats

1) Key value input format

Code:-

q1)working with keyval input pair

```
import java.io.IOException;

import org.apache.hadoop.conf.*;

import org.apache.hadoop.fs.Path;

import org.apache.hadoop.io.Text;

import org.apache.hadoop.mapreduce.*;

import org.apache.hadoop.mapreduce.Mapper.Context;

import org.apache.hadoop.mapreduce.lib.output.*;

import org.apache.hadoop.mapreduce.lib.input.*;

public class keyval {

    public static void main(String[] args) throws IOException, ClassNotFoundException, InterruptedException {

        Configuration conf = new Configuration();

        //setting to keyval pair

        conf.set("mapreduce.input.keyvalinelinerecordreader.key.value.separator", ",");

        Job job = new Job(conf, "keyval");
```

```

        job.setJarByClass(keyval.class);

        job.setOutputKeyClass(Text.class);

        job.setOutputValueClass(Text.class);

        job.setMapperClass(Map.class);


        //

        job.setInputFormatClass(KeyValueTextInputFormat.class);

        //

        job.setOutputFormatClass(TextOutputFormat.class);

        FileInputFormat.addInputPath(job, new Path(args[0]));

        FileOutputFormat.setOutputPath(job, new Path(args[1]));


        job.waitForCompletion(true);

    }

    //all the arguments are text.

    public static class Map extends Mapper<Text, Text, Text, Text>{

        String c="ram";

        public void map(Text key, Text value, Context context) throws

            IOException, InterruptedException

        {

            String[] line=key.toString().split(",");

            if(c.equalsIgnoreCase(line[0]))

```

```
        context.write(key,value);
    }
}
}
```

Input will be taken considered as pair of key and value.

2)Nline input format

No.of lines are set per each mapper. So based on number of lines of data input file, no of mappers are decided.

Code:-

```
import java.io.*;

import java.util.*;

import org.apache.hadoop.fs.Path;

import org.apache.hadoop.io.*;

import org.apache.hadoop.conf.*;

import org.apache.hadoop.mapreduce.*;

import org.apache.hadoop.mapreduce.Reducer.Context;

import org.apache.hadoop.mapreduce.lib.input.*;

import org.apache.hadoop.mapreduce.lib.output.*;

public class nline{

public static class Map extends Mapper<Object,Text,Text,Text>{

String c = "ram";
```

```
public void map(Text key,Text value,Context context)throws IOException,InterruptedException{
```

```
String line=key.toString();
```

```
if(c.equalsIgnoreCase(line))
```

```
context.write(key,value);
```

```
}
```

```
}
```

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

```
Configuration conf=new Configuration();
```

```
//set no of lines per each mapper.
```

```
conf.setInt(NLIneInputFormat.LINES_PER_MAP, 3);
```

```
Job job = new Job(conf,"nline");
```

```
job.setJarByClass(nline.class);
```

```
job.setOutputKeyClass(Text.class);
```

```
job.setOutputValueClass(Text.class);
```

```
job.setMapperClass(Map.class);
```

```
job.setNumReduceTasks(0);//setting no of reducers to 0
```

```
job.setInputFormatClass(NLIneInputFormat.class);
```

```
FileInputFormat.addInputPath(job,new Path(args[0]));
```

```
FileOutputFormat.setOutputPath(job,new Path(args[1]));
```

```
job.waitForCompletion(true);
```

```
}
```

```
}
```

3) Sequential file input and output.

A file that contains records or other elements that are stored in a chronological order based on account number or some other identifying data.

Code:-

```
*****output part

import java.io.*;

import java.util.*;

import org.apache.hadoop.fs.Path;

import org.apache.hadoop.io.*;

import org.apache.hadoop.conf.*;

import org.apache.hadoop.mapreduce.*;

import org.apache.hadoop.mapreduce.Reducer.Context;

import org.apache.hadoop.mapreduce.lib.input.*;

import org.apache.hadoop.mapreduce.lib.output.*;

public class output{

    public static class Map extends Mapper<Object,Text,LongWritable,Text>{

        public void map(LongWritable key,Text value,Context context)throws IOException,InterruptedException{

            context.write(key,value);

        }

    }

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

        Configuration conf=new Configuration();
```

```
conf.setInt(NLineInputFormat.LINES_PER_MAP, 3);

Job job = new Job(conf,"output");

job.setJarByClass(output.class);

job.setOutputKeyClass(LongWritable.class);

job.setOutputValueClass(Text.class);

job.setMapperClass(Map.class);

job.setNumReduceTasks(0);

//setting output as sequential file

job.setOutputFormatClass(SequenceFileOutputFormat.class);

//

FileInputFormat.addInputPath(job,new Path(args[0]));

FileOutputFormat.setOutputPath(job,new Path(args[1]));

job.waitForCompletion(true);

}

}

*****input part

import java.io.*;

import java.util.*;

import org.apache.hadoop.fs.Path;

import org.apache.hadoop.io.*;

import org.apache.hadoop.conf.*;

import org.apache.hadoop.mapreduce.*;
```

```
import org.apache.hadoop.mapreduce.Reducer.Context;

import org.apache.hadoop.mapreduce.lib.input.*;

import org.apache.hadoop.mapreduce.lib.output.*;

public class input{

    public static class Map extends Mapper<LongWritable,Text,LongWritable,Text>{

        public void map(LongWritable key,Text value,Context context)throws IOException,InterruptedException{

            context.write(key,value);

        }

    }

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

        Configuration conf=new Configuration();

        //conf.setInt(NLineInputFormat.LINES_PER_MAP, 3);

        Job job = new Job(conf,"input");

        job.setJarByClass(input.class);

        job.setOutputKeyClass(LongWritable.class);

        job.setOutputValueClass(Text.class);

        job.setMapperClass(Map.class);

        job.setNumReduceTasks(0);

        //input as sequential data.

        job.setInputFormatClass(SequenceFileInputFormat.class);

        //

        FileInputFormat.addInputPath(job,new Path(args[0]));
```

```
FileOutputFormat.setOutputPath(job,new Path(args[1]));
```

```
job.waitForCompletion(true);
```

```
}
```

```
}
```

Lab7 :- Finding the top k records using tree map.

Algorithm:- general method of sorting data and resulting the top k won't be a efficient method for big data as data is large and parallel computations are incorporated. So we use tree map which sorts the data as it gets entered.

Code:-

```
import java.io.IOException;
```

```
import java.util.TreeMap;
```

```
import org.apache.hadoop.conf.*;
```

```
import org.apache.hadoop.fs.Path;
```

```
import org.apache.hadoop.io.*;
```

```
import org.apache.hadoop.mapreduce.*;
```

```
import org.apache.hadoop.mapreduce.lib.output.*;
```

```
import org.apache.hadoop.mapreduce.lib.input.*;
```

```
public class topk {
```

```
    public static void main(String[] args) throws IOException, ClassNotFoundException, InterruptedException {
```

```
        Configuration conf = new Configuration();
```

```
        Job job = new Job(conf, "topk");
```

```
        job.setJarByClass(topk.class);

        //

        job.setOutputKeyClass(NullWritable.class);

        job.setOutputValueClass(Text.class);

        job.setNumReduceTasks(1);

        //

        /*job.setOutputKeyClass(Text.class);

        job.setOutputValueClass(IntWritable.class);

        */

        job.setMapperClass(Map.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);

    }

    public static class Map extends Mapper<LongWritable,
    Text, NullWritable, Text>

    {
```

```

private TreeMap<Integer, Text> salary = new    TreeMap<Integer, Text>();

public void map(LongWritable key, Text value, Context
context) throws IOException, InterruptedException

{

String line = value.toString();

String[] elements=line.split(",");

int

i= Integer.parseInt(elements[1]);

salary.put(i,new Text(value));

//specifying the k value as 10.

if (salary.size() > 10) {

salary.remove(salary.firstKey());

}

}

protected void cleanup(Context context) throws

IOException, InterruptedException

{

for ( Text name : salary.values() ) {

context.write(NullWritable.get(), name);

}

}

}

public static class Reduce extends

```

```
Reducer<NullWritable, Text, NullWritable, Text>

{

    public void reduce(NullWritable key, Iterable<Text>

    values, Context context) throws IOException,

    InterruptedException {

        TreeMap<Integer, Text> salary = new TreeMap< Integer,Text>();

        for (Text value : values) {

            String line = value.toString();

            String[] elements=line.split(",");

            int      i= Integer.parseInt(elements[1]);

            salary.put(i, new Text(value));

            if (salary.size() > 10) {

                salary.remove(salary.firstKey());

            }

        }

        for (Text t : salary.values()) {

            context.write(NullWritable.get(), t);

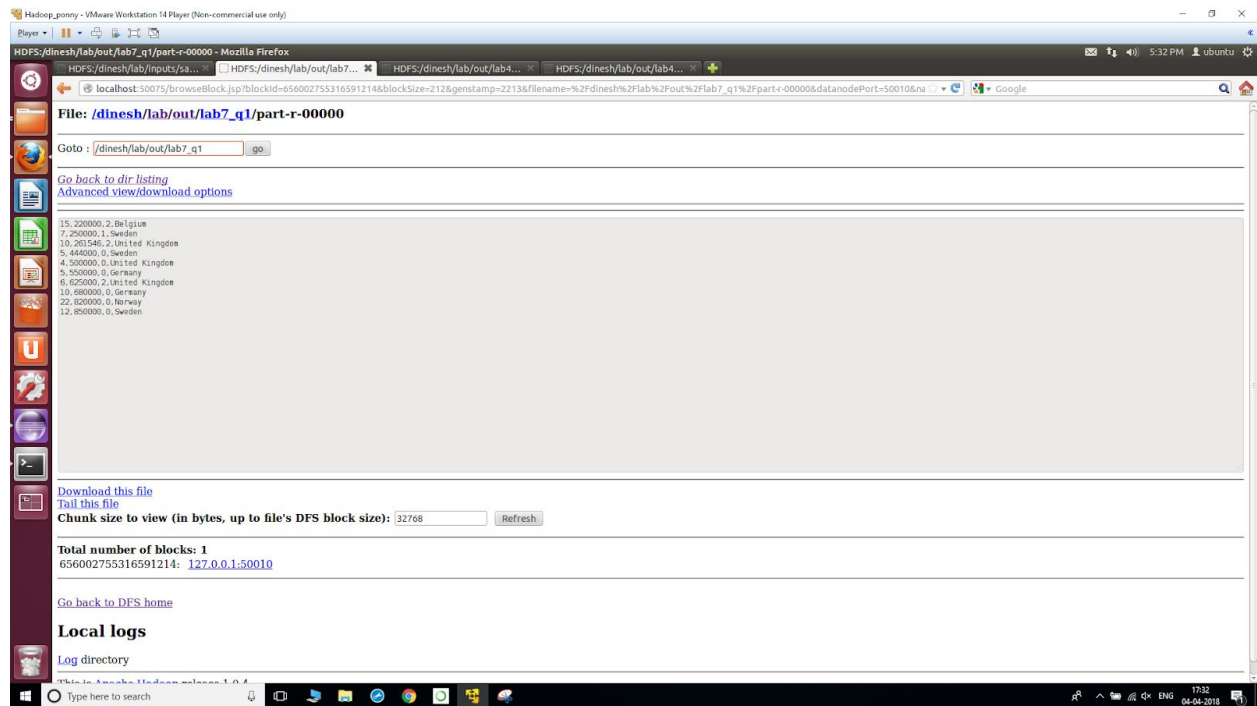
        }

    }

}

}
```

output:-



q2) User defined counters

code:

```
import java.io.IOException;
```

```
import java.util.TreeMap;
```

```
import org.apache.hadoop.conf.*;
```

```
import org.apache.hadoop.fs.Path;
```

```
import org.apache.hadoop.io.*;
```

```
import org.apache.hadoop.mapreduce.*;
```

```
import org.apache.hadoop.mapreduce.lib.output.*;
```

```
import org.apache.hadoop.mapreduce.lib.input.*;
```

```
public class count {
```

```
    public static void main(String[] args) throws IOException, ClassNotFoundException, InterruptedException {
```

```
        Configuration conf = new Configuration();
```

```
        Job job = new Job(conf, "topk");
```

```
        job.setJarByClass(count.class);
```

```
        //
```

```
        job.setNumReduceTasks(0);
```

```
        //
```

```
        job.setOutputKeyClass(Text.class);
```

```
        job.setOutputValueClass(IntWritable.class);
```

```
        job.setMapperClass(Map.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);
```

```
        Counters cn=job.getCounters();
```

```
        cn.findCounter(ct.cnt).getValue();

        cn.findCounter(ct.nt).getValue();

    }
}
```

```
public enum ct

{

    cnt,nt

};

public static class Map extends Mapper<LongWritable,

Text, Text, IntWritable>

{

    public void map(LongWritable key, Text value, Context

context) throws IOException, InterruptedException {

        String line = value.toString();

        String[] elements=line.split(",");

        int i= Integer.parseInt(elements[1]); //i takes the salary values

        float exp= Float.parseFloat(elements[0]);

        Text tt=new Text(elements[3]); //tt takes the country names.

        if(exp==0.0)
```

```
{  
  
    //incrementing the counter  
    context.getCounter(ct.cnt).increment(1);  
  
    context.write(tt,new IntWritable(i));  
}  
  
if(i > 50000)  
{  
    context.getCounter(ct.nt).increment(1);  
}  
}  
}
```

```
}
```

Input:- salary.csv

output:-

```
Warning: $HADOOP_HOME is deprecated.
18/04/04 17:36:06 WARN mapred.JobClient: Use GenericOptionsParser for parsing the arguments. Applications should implement Tool for the same.
18/04/04 17:36:06 INFO input.FileInputFormat: Total input paths to process : 1
18/04/04 17:36:06 INFO util.NativeCodeLoader: Loaded the native-hadoop library
18/04/04 17:36:06 WARN snappy.LoadSnappy: Snappy native library not loaded
18/04/04 17:36:06 INFO mapred.JobClient: Running job: job_201804041630_0003
18/04/04 17:36:07 INFO mapred.JobClient: map 0% reduce 0%
18/04/04 17:36:23 INFO mapred.JobClient: map 100% reduce 0%
18/04/04 17:36:28 INFO mapred.JobClient: Job complete: job_201804041630_0003
18/04/04 17:36:28 INFO mapred.JobClient: Counters: 21
18/04/04 17:36:28 INFO mapred.JobClient:   Job Counters
18/04/04 17:36:28 INFO mapred.JobClient:     SLOTS_MILLIS_MAPS=14572
18/04/04 17:36:28 INFO mapred.JobClient:     Total time spent by all reduces waiting after reserving slots (ms)=0
18/04/04 17:36:28 INFO mapred.JobClient:     Total time spent by all maps waiting after reserving slots (ms)=0
18/04/04 17:36:28 INFO mapred.JobClient:     Launched map tasks=1
18/04/04 17:36:28 INFO mapred.JobClient:     Data-local map tasks=1
18/04/04 17:36:28 INFO mapred.JobClient:     SLOTS_MILLIS_REDUCES=0
18/04/04 17:36:28 INFO mapred.JobClient:   Lab7_q2$ct
18/04/04 17:36:28 INFO mapred.JobClient:     nt=476
18/04/04 17:36:28 INFO mapred.JobClient:     cnt=26
18/04/04 17:36:28 INFO mapred.JobClient:   File Output Format Counters
18/04/04 17:36:28 INFO mapred.JobClient:     Bytes Written=411
18/04/04 17:36:28 INFO mapred.JobClient:   FileSystemCounters
18/04/04 17:36:28 INFO mapred.JobClient:     HDFS_BYTES_READ=23851
18/04/04 17:36:28 INFO mapred.JobClient:     FILE_BYTES_WRITTEN=21100
18/04/04 17:36:28 INFO mapred.JobClient:     HDFS_BYTES_WRITTEN=411
18/04/04 17:36:28 INFO mapred.JobClient:   File Input Format Counters
18/04/04 17:36:28 INFO mapred.JobClient:     Bytes Read=23734
18/04/04 17:36:28 INFO mapred.JobClient:   Map-Reduce Framework
18/04/04 17:36:28 INFO mapred.JobClient:     Map input records=1147
18/04/04 17:36:28 INFO mapred.JobClient:     Physical memory (bytes) snapshot=41254912
18/04/04 17:36:28 INFO mapred.JobClient:     Spilled Records=0
18/04/04 17:36:28 INFO mapred.JobClient:     CPU time spent (ms)=90
18/04/04 17:36:28 INFO mapred.JobClient:     Total committed heap usage (bytes)=16252928
18/04/04 17:36:28 INFO mapred.JobClient:     Virtual memory (bytes) snapshot=382939136
18/04/04 17:36:28 INFO mapred.JobClient:     Map output records=26
18/04/04 17:36:28 INFO mapred.JobClient:     SPLIT_RAW_BYTES=117
```

Lab 8 :-

q1) Side data

Take data from console and use it in Mapreduce job.

Code:-

```
q1) // side data. taaking data from terminal as final argument.
```

```
import java.io.*;
```

```
import org.apache.hadoop.fs.Path;
```

```
import org.apache.hadoop.io.Text;
```

```
import org.apache.hadoop.conf.*;
```

```
import org.apache.hadoop.mapreduce.*;
```

```
import org.apache.hadoop.mapreduce.lib.input.*;
```

```
import org.apache.hadoop.mapreduce.lib.output.*;
```

```
public class student{

    public static class Map extends Mapper<Text,Text,Text,Text>

    {

        public void map(Text key, Text value, Context context) throws IOException, InterruptedException

        {

            String name = context.getConfiguration().get("name");

            String[] elements = key.toString().split(",");

            if(name.equalsIgnoreCase(elements[1]))

                context.write(key,value);

        }

    }

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

    {

        Configuration conf=new Configuration();

        conf.set("name",args[2]);

        Job job=new Job(conf,"student");

        job.setInputFormatClass(KeyValueTextInputFormat.class);

        job.setJarByClass(student.class);

        job.setOutputKeyClass(Text.class);

        job.setOutputValueClass(Text.class);

        job.setMapperClass(Map.class);

        job.setOutputFormatClass(TextOutputFormat.class);
```

```
FileOutputFormat.setOutputPath(job,new Path(args[1]));

job.waitForCompletion(true);

}
```

The screenshot shows a virtual machine environment with a Windows desktop. The desktop has a taskbar with various application icons. The active window is a Mozilla Firefox browser displaying the HDFS web interface. The browser's address bar shows the URL 'localhost:50075/browseBlock.jsp?blockid=6875316691922594645&blockSize=94&genstamp=2277&filename=%2Fdinesh%2Flab%2Finputs%2Flab8_q1.txt&datanodePort=50010&namenodein...'. The page title is 'File: /dinesh/lab/inputs/lab8_q1.txt'. The main content area shows a file list with columns for file name, size, and modification time. The file 'lab8_q1.txt' is listed with a size of 32768 bytes and a modification time of 2011-10-10 10:10:10. Below the file list, there are links for 'Download this file' and 'Tail this file'. The 'Download this file' link is highlighted. The 'Tail this file' link is also visible. The 'Chunk size to view (in bytes, up to file's DFS block size):' is set to 32768. The 'Total number of blocks: 1' is displayed. The 'Log directory' link is visible at the bottom of the page. The desktop taskbar shows the Start button, a search bar, and several application icons including Internet Explorer, Firefox, and various system utilities. The system clock in the bottom right corner shows the date '04-04-2013' and time '17:39'.

The screenshot shows a web browser window with the title "Hadoop_gromy - VMware Workstation 14 Player (Non-commercial use only)". The browser address bar shows the URL "HDFS/dinesh/lab/out/lab8_q1/part-r-00000 - Mozilla Firefox". The page content includes a navigation bar with tabs for "HDFS/dinesh/lab/out/lab8_q1/part-r-00000", "HDFS/dinesh/lab/inputs/la...", "HDFS/dinesh/lab/out/lab4...", and "HDFS/dinesh/lab/out/lab4...". The main content area displays the file path "File: /dinesh/lab/out/lab8_q1/part-r-00000" and a "Goto" button. Below the "Goto" button, there are links for "Go back to dir listing" and "Advanced view/download options". The file details section shows the file path "1. gourthaa.100.100.100.100" and a "Download this file" link. The "Local logs" section is also visible at the bottom.

```
code:-q2) //basic reduce side join
```

53

```
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;

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

import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;


public class multi {

    public static class Map1 extends Mapper<LongWritable, Text, Text, Text>{

        public void map(LongWritable key, Text value, Context context) throws IOException,
        InterruptedException

        {

            String[] line=value.toString().split(",");

            context.write(new Text(line[0]),new Text(line[1] + " " + line[2]));

        }

    }

    public static class Map2 extends Mapper<LongWritable, Text, Text, Text>{

        public void map(LongWritable key, Text value, Context context) throws IOException,
        InterruptedException

        {

            String[] line=value.toString().split(",");

            context.write(new Text(line[0]),new Text(line[1]));

        }

    }

    public static class Reduce extends Reducer<Text, Text, Text, Text>{

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

```
        String line = null;

        for(Text val:values)

        {

            line = val.toString();

        }

        context.write(key, new Text(line));

    }

}

public static void main(String[] args) throws IOException, ClassNotFoundException,
InterruptedException {

    Configuration conf = new Configuration();

    Job job = new Job(conf, "multi");

    job.setJarByClass(multi.class);

    job.setOutputKeyClass(Text.class);

    job.setOutputValueClass(Text.class);

    job.setMapperClass(Map1.class);

    job.setMapperClass(Map2.class);

    job.setInputFormatClass(TextInputFormat.class);

    job.setOutputFormatClass(TextOutputFormat.class);

    MultipleInputs.addInputPath(job,new Path(args[0]), TextInputFormat.class,
Map1.class);

    MultipleInputs.addInputPath(job,new Path(args[1]), TextInputFormat.class,
Map2.class);

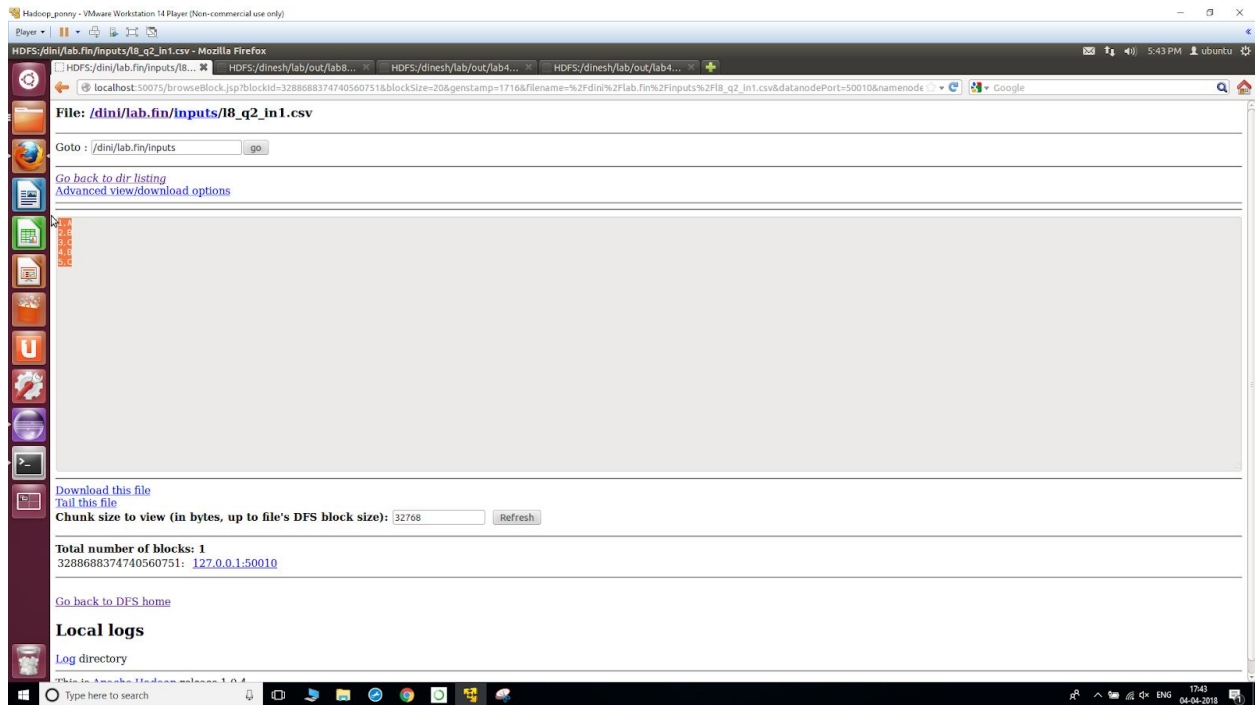
    FileOutputFormat.setOutputPath(job, new Path(args[2]));
```

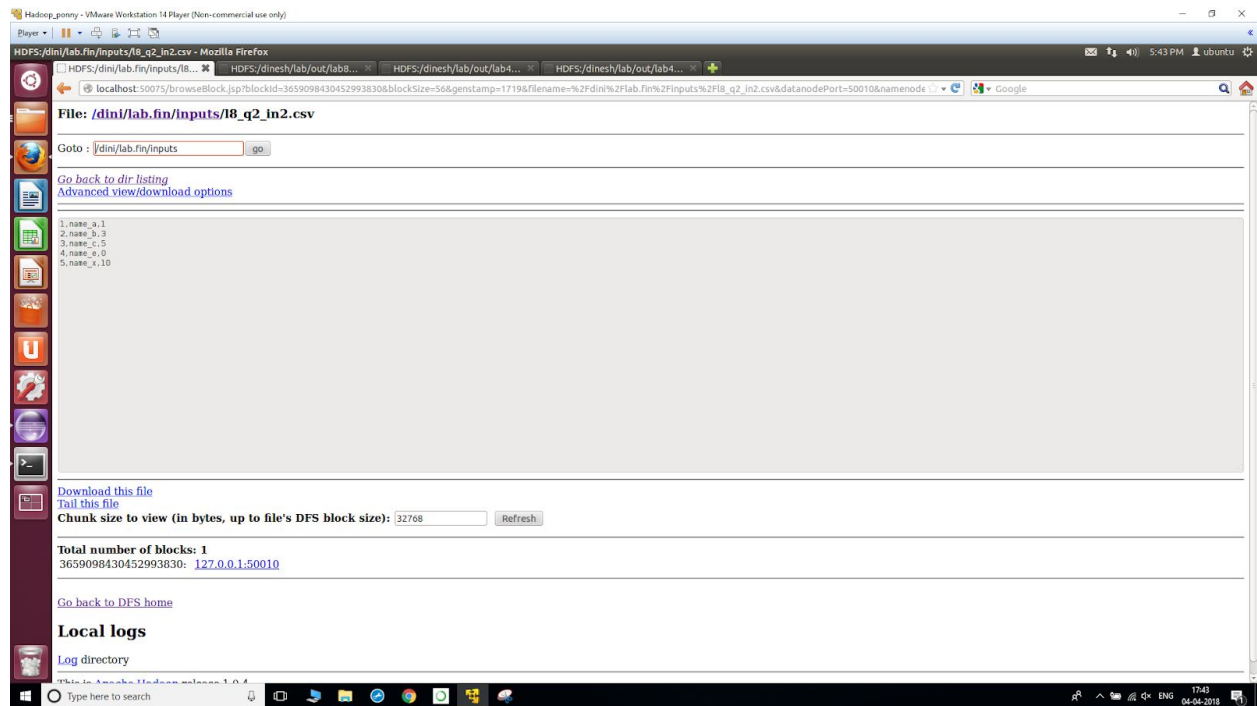
```
job.waitForCompletion(true);
```

```
}
```

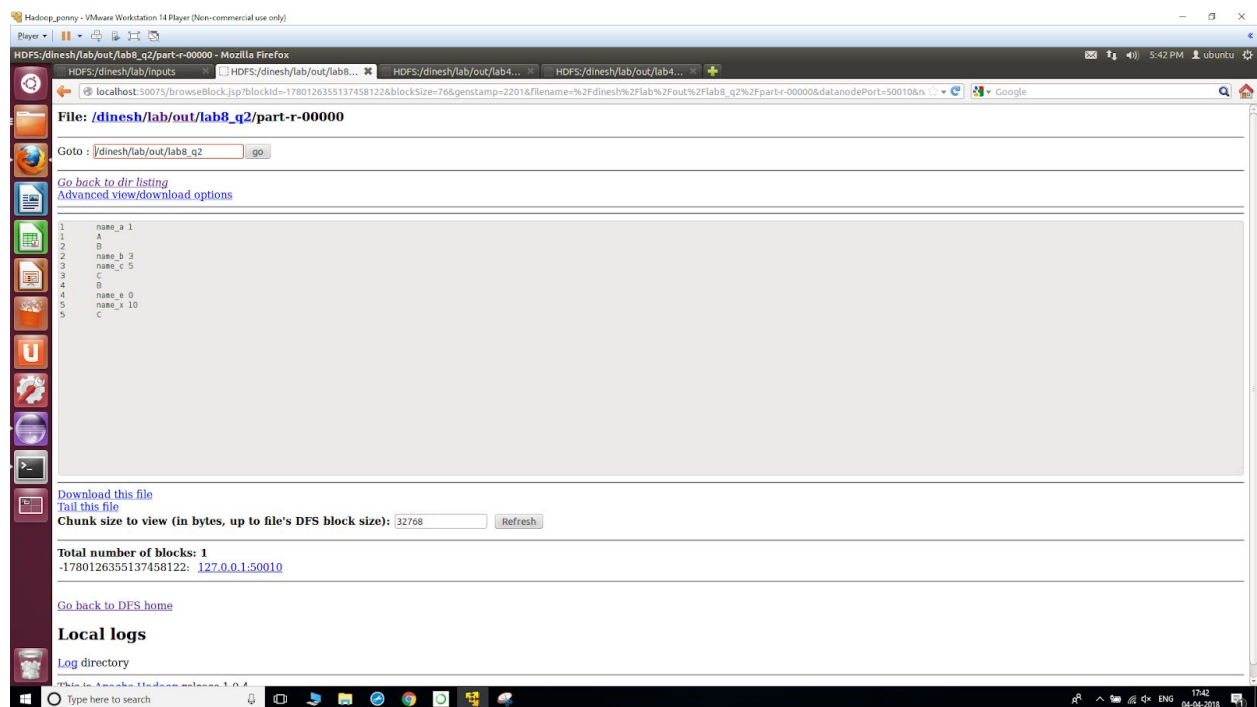
```
}
```

Input:- 2 files





Output:- writing the data after joining 2 files.



l9)

q1) distributed cache.

One of the 2 files is cached for faster access. The file that is more frequently used is cached.

Code:-

```
import java.io.*;

import java.util.*;

import org.apache.hadoop.conf.Configuration;

import org.apache.hadoop.fs.Path;

import org.apache.hadoop.io.*;

import org.apache.hadoop.mapreduce.*;

import org.apache.hadoop.filecache.*;

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 emp{

    public static class Map extends Mapper<LongWritable,Text,Text,Text> {

        Path[] cfile=new Path[0];

        ArrayList<Text> empl=new ArrayList<Text>();

        public void setup(Context context)

        {
```

```
Configuration conf=context.getConfiguration();

try
{
    cfile = DistributedCache.getLocalCacheFiles(conf);

    BufferedReader reader=new BufferedReader(new FileReader(cfile[0].toString()));

    String line;

    while ((line=reader.readLine())!=null)
    {
        Text tt=new Text(line);

        empl.add(tt);
    }
}

catch(IOException e)
{
    e.printStackTrace();
}

}

public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException {

    String[] elements = value.toString().split(",");

    for(Text e:empl)
    {
        String[] line1 = e.toString().split(",");

        if(elements[0].equals(line1[0]))
```

```

        {

            context.write(new Text(elements[0]),new Text(elements[1]+"."+elements[2]+"."+line1[1]));

        }

    }

}

```

```

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

    Configuration conf = new Configuration();

    Job job = new Job(conf,"emp");

    job.setJarByClass(emp.class);

    DistributedCache.addCacheFile(new Path(args[0]).toUri(),job.getConfiguration());

    job.setMapperClass(Map.class);

    job.setOutputKeyClass(Text.class);

    job.setOutputValueClass(Text.class);

    //job.setInputFormatClass(TextInputFormat.class);

    //job.setOutputFormatClass(TextOutputFormat.class);

    FileInputFormat.addInputPath(job, new Path(args[1]));

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

    job.waitForCompletion(true);

}
}

```

Input and output will be same as normal join but the execution will be faster.

q2) reduce side join customer and transaction details.

Find total amount of transaction of a customer.

code:-

```
import java.io.IOException;

import org.apache.hadoop.conf.Configuration;

import org.apache.hadoop.fs.Path;

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.MultipleInputs;

import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;

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


public class ReduceJoin {

    public static class CustsMapper extends Mapper<Object, Text, Text, Text>

    {

        public void map(Object key, Text value, Context context)

        throws IOException, InterruptedException

        {

            String record = value.toString();

            String[] parts = record.split(",");

            context.write(new Text(parts[0]), new Text("cust\t" + parts[1]));

        }

    }

}
```

```
}
```

```
}
```

```
public static class TxnsMapper extends Mapper<Object, Text, Text, Text>
```

```
{
```

```
public void map(Object key, Text value, Context context)
```

```
throws IOException, InterruptedException
```

```
{
```

```
String record = value.toString();
```

```
String[] parts = record.split(",");
```

```
context.write(new Text(parts[2]), new Text("txn\t" + parts[3]));
```

```
}
```

```
}
```

```
public static class ReduceJoinReducer extends Reducer<Text, Text, Text, Text>
```

```
{
```

```
public void reduce(Text key, Iterable<Text> values, Context context)
```

```
throws IOException, InterruptedException
```

```
{
```

```
String name = "";
```

```
double total = 0.0;
```

```
int count = 0;
```

```
for (Text t : values)
```

```
{

String parts[] = t.toString().split("\\t");

if (parts[0].equals("tnxn"))

{

count++;

total += Float.parseFloat(parts[1]);

}

else if (parts[0].equals("cust"))

{

name = parts[1];

}

}

String str = String.format("%d\\t%f", count, total);

context.write(new Text(name), new Text(str));

}

}

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

Configuration conf = new Configuration();

Job job = new Job(conf, "Reduce-side join");

job.setJarByClass(ReduceJoin.class);

job.setReducerClass(ReduceJoinReducer.class);

job.setOutputKeyClass(Text.class);
```

```
job.setOutputValueClass(Text.class);
```

```
MultipleInputs.addInputPath(job, new Path(args[0]),TextInputFormat.class, CustsMapper.class);
```

```
MultipleInputs.addInputPath(job, new Path(args[1]),TextInputFormat.class, TxnsMapper.class);
```

```
Path outputPath = new Path(args[2]);
```

```
FileOutputFormat.setOutputPath(job, outputPath);
```

```
outputPath.getFileSystem(conf).delete(outputPath);
```

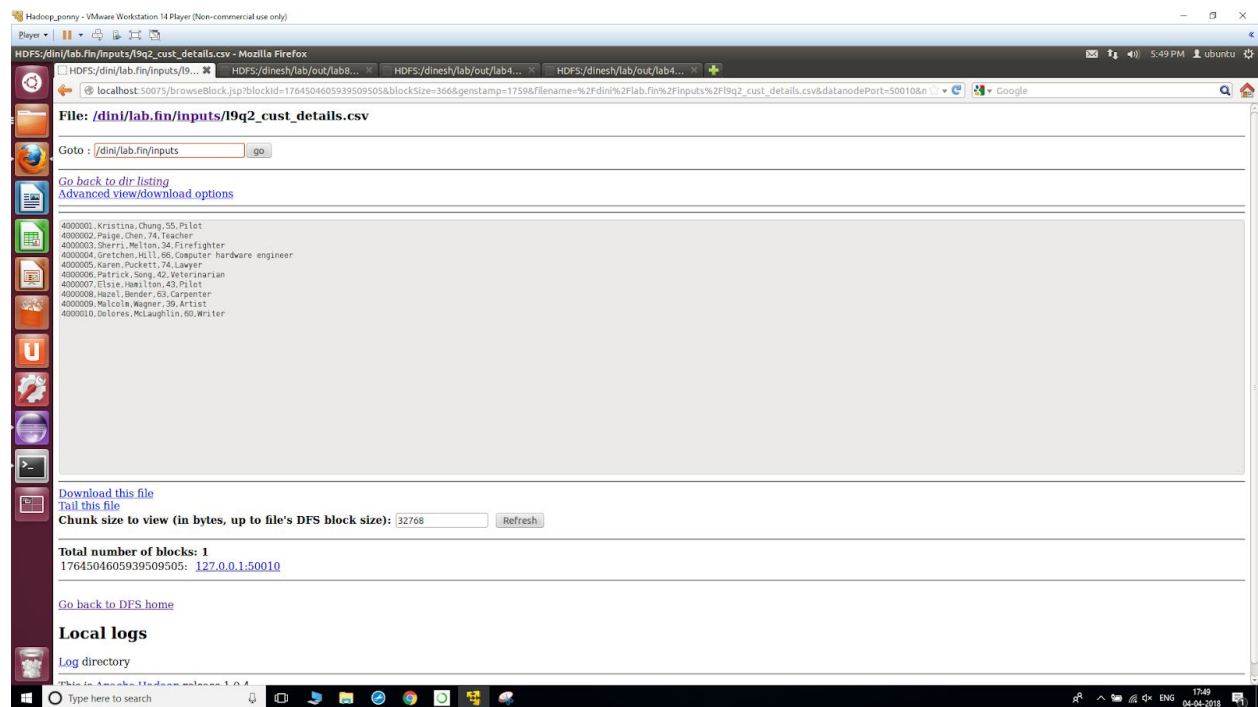
```
System.exit(job.waitForCompletion(true) ? 0 : 1);
```

```
}
```

```
}
```

Input:-

Cust_details



transaction_details

The screenshot shows a Hadoop VM workspace with a web browser displaying the file `File: /dini/lab.fin/inputs/I9q2_trans_details.csv`. The browser interface includes a sidebar with icons for various applications and a main content area showing a list of transactions. The transactions are listed in a table with columns for date, time, and description. The total number of blocks is 1, and the total size is 1531604329078060047. The browser also shows a 'Local logs' section with a 'Log directory' link.

Date	Time	Description
00000000	06-26-2011	4000001,040.33,Exercise & Fitness,Cardio Machine Accessories,Clarksville,Tennessee,credit
00000001	05-26-2011	4000002,198.44,Exercise & Fitness,Weightlifting Gloves,Long Beach,California,credit
00000002	06-01-2011	4000002,005.56,Exercise & Fitness,Weightlifting Machine Accessories,Anaheim,California,credit
00000003	06-05-2011	4000003,198.10,Gymnastics,Gymnastics Rings,Milwaukee,Wisconsin,credit
00000004	12-17-2011	4000002,098.81,Team Sports,Field Hockey,Nashville,Tennessee,credit
00000005	02-14-2011	4000004,199.63,Outdoor Recreation,Camping & Backpacking & Hiking,Chicago,Illinois,credit
00000006	10-28-2011	4000005,027.89,Puzzles,Jigsaw Puzzles,Charleston,South Carolina,credit
00000007	07-14-2011	4000006,096.01,Outdoor Play Equipment,Sandboxes,Columbus,Ohio,credit
00000008	01-17-2011	4000006,010.44,Winter Sports,Snowboarding,Des Moines,Iowa,credit
00000009	05-17-2011	4000006,152.46,Jumping,Bungee Jumping,St. Petersburg,Florida,credit
00000010	05-29-2011	4000007,180.26,Outdoor Recreation,Archery,Reno,Nevada,credit
00000011	06-18-2011	4000009,121.39,Outdoor Play Equipment,Swing Sets,Columbus,Ohio,credit
00000012	02-08-2011	4000009,041.52,Indoor Games,Bowling,San Francisco,California,credit
00000013	03-19-2011	4000010,107.80,Team Sports,Field Hockey,Honolulu,Hawaii,credit
00000014	02-25-2011	4000010,036.81,Gymnastics,Vaulting Horses,Los Angeles,California,credit
00000015	10-20-2011	4000001,137.64,Combat Sports,Fencing,Honolulu,Hawaii,credit
00000016	05-28-2011	4000010,035.56,Exercise & Fitness,Free Weight Bars,Columbia,South Carolina,credit
00000017	10-18-2011	4000008,075.55,Water Sports,Scuba Diving & Snorkeling,Oakland,Nebraska,credit
00000018	11-18-2011	4000008,088.65,Team Sports,Basball,Salt Lake City,Utah,credit
00000019	08-28-2011	4000008,051.81,Water Sports,Life Jackets,Newark,New Jersey,credit
00000020	06-29-2011	4000005,041.55,Exercise & Fitness,Weightlifting Belts,New Orleans,Louisiana,credit
00000021	02-14-2011	4000005,045.79,Air Sports,Parachutes,New York,New York,credit
00000022	10-10-2011	4000009,019.64,Water Sports,Kitesurfing,Saint Paul,Minnesota,credit
00000023	05-02-2011	4000009,099.50,Gymnastics,Gymnastics Rings,Springfield,Illinois,credit
00000024	06-10-2011	4000003,151.20,Water Sports,Surfing,Plano,Texas,credit
00000025	10-14-2011	4000009,144.20,Indoor Games,Darts,Phoenix,Arizona,credit

Output:-

The screenshot shows a Hadoop VM workspace with a web browser displaying the file `File: /dini/lab.fin/outputs/I9/out_I9q2f/part-r-000000`. The browser interface includes a sidebar with icons for various applications and a main content area showing a list of output records. The records are listed in a table with columns for name, age, and address. The total number of blocks is 1, and the total size is 4947296396033794800. The browser also shows a 'Local logs' section with a 'Log directory' link.

Name	Age	Address
Kristina	8	651.049999
Paige	6	706.970007
Sherris	3	527.589996
Gretchen	5	337.060005
Karen	5	325.150005
Patrick	5	539.980019
Elise	6	699.550003
Hazel	10	859.419999
Melissa	6	457.520006
Dolores	6	447.990004

Lab10:- Secondary sorting

Hadoop automatically sorts by key. If we want to sort by value also, or a part of value, that process is called secondary sorting. With a slight manipulation to the format of the key object, secondary sorting gives us the ability to take the value into account during the sort phase by creating a composite key by adding a part of, or the entire value to the natural key to achieve your sorting objectives.

Input contains country name and state name our aim to sort by country and inside the country sort by state.

Code:-

```
//secondary sorting

import java.io.BufferedReader;

import java.io.DataInput;

import java.io.DataOutput;

import java.io.FileReader;

import java.io.IOException;

import java.nio.file.FileSystem;

import java.util.*;

import org.apache.hadoop.filecache.DistributedCache;

import org.apache.hadoop.fs.Path;

import org.apache.hadoop.conf.*;

import org.apache.hadoop.io.*;

import org.apache.hadoop.mapreduce.*;

import org.apache.hadoop.mapreduce.Reducer.Context;

import org.apache.hadoop.mapreduce.lib.output.*;
```

```
import org.apache.hadoop.mapreduce.lib.input.*;

public class lab10 {

    public static class CompositeKeyWritable implements Writable, WritableComparable<CompositeKeyWritable> {

        private String deptNo;

        private String emp;

        public CompositeKeyWritable() {

        }

        public CompositeKeyWritable(String deptNo, String emp) {

            this.deptNo = deptNo;

            this.emp = emp;

        }

        public String toString() {

            return (new StringBuilder().append(deptNo).append("\t").append(emp)).toString();

        }

        public void readFields(DataInput dataInput) throws IOException {

            deptNo = WritableUtils.readString(dataInput);

            emp = WritableUtils.readString(dataInput);

        }

        public void write(DataOutput dataOutput) throws IOException {

            WritableUtils.writeString(dataOutput,deptNo);

            WritableUtils.writeString(dataOutput,emp);

        }

    }

}
```

```

    }

    public int compareTo(CompositeKeyWritable objKeyPair) {

        int result = deptNo.compareTo(objKeyPair.deptNo);

        if (0 == result) {

            result = emp.compareTo(objKeyPair.emp);

        }

        return result;

    }

}

public static class mapper1 extends Mapper<LongWritable, Text,
CompositeKeyWritable, NullWritable> {

    public void map(LongWritable key, Text value, Context context) throws
IOException, InterruptedException {

        if (value.toString().length() > 0) {

            String arrEmpAttributes[] = value.toString().split(",");

            context.write(new

                CompositeKeyWritable(arrEmpAttributes[0].toString(),(arrEmpAttributes[1].toString()),NullWritable.get());

        }

    }

}

public static class SecondarySortBasicPartitioner extends Partitioner<CompositeKeyWritable, NullWritable> {

    public int getPartition(CompositeKeyWritable key, NullWritable

        value, int numReduceTasks) {

```

```

        return (key.deptNo.hashCode() % numReduceTasks);
    }
}

public static class SecondarySortBasicCompKeySortComparator extends WritableComparator {

    protected SecondarySortBasicCompKeySortComparator()

    {
        super(CompositeKeyWritable.class, true);
    }

    public int compare(WritableComparable w1, WritableComparable w2) {

        CompositeKeyWritable key1 = (CompositeKeyWritable) w1;
        CompositeKeyWritable key2 = (CompositeKeyWritable) w2;

        int cmpResult = key1.deptNo.compareTo(key2.deptNo);

        if (cmpResult == 0)
        {
            return -key1.emp.compareTo(key2.emp);
        }

        return cmpResult;
    }
}

public static class SecondarySortBasicGroupingComparator extends WritableComparator {

    protected SecondarySortBasicGroupingComparator() {

        super(CompositeKeyWritable.class, true);
    }
}

```

```
public int compare(WritableComparable w1,WritableComparable w2) {

    CompositeKeyWritable key1 = (CompositeKeyWritable) w1;

    CompositeKeyWritable key2 = (CompositeKeyWritable) w2;

    return key1.deptNo.compareTo(key2.deptNo);

}

}

public static class SecondarySortBasicReducer extends Reducer<CompositeKeyWritable, NullWritable, CompositeKeyWritable,

NullWritable>{

    public void reduce(CompositeKeyWritable key, Iterable<NullWritable>

        values, Context context) throws IOException, InterruptedException

    {

        for(NullWritable val:values)

        {

            context.write(key,NullWritable.get());

        }

    }

}

}

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

    Configuration conf=new Configuration();

    Job job=new Job(conf,"country");

    job.setJarByClass(lab10.class);
```

```
job.setOutputKeyClass(Text.class);

job.setOutputValueClass(Text.class);

job.setMapperClass(mapper1.class);

job.setMapOutputKeyClass(CompositeKeyWritable.class);

job.setMapOutputValueClass(NullWritable.class);

job.setPartitionerClass(SecondarySortBasicPartitioner.class);

job.setSortComparatorClass(SecondarySortBasicCompKeySortComparator.class);

job.setGroupingComparatorClass(SecondarySortBasicGroupingComparator.class);

job.setReducerClass(SecondarySortBasicReducer.class);

job.setInputFormatClass(TextInputFormat.class);

job.setOutputFormatClass(TextOutputFormat.class);

FileInputFormat.addInputPath(job, new Path(args[0]));

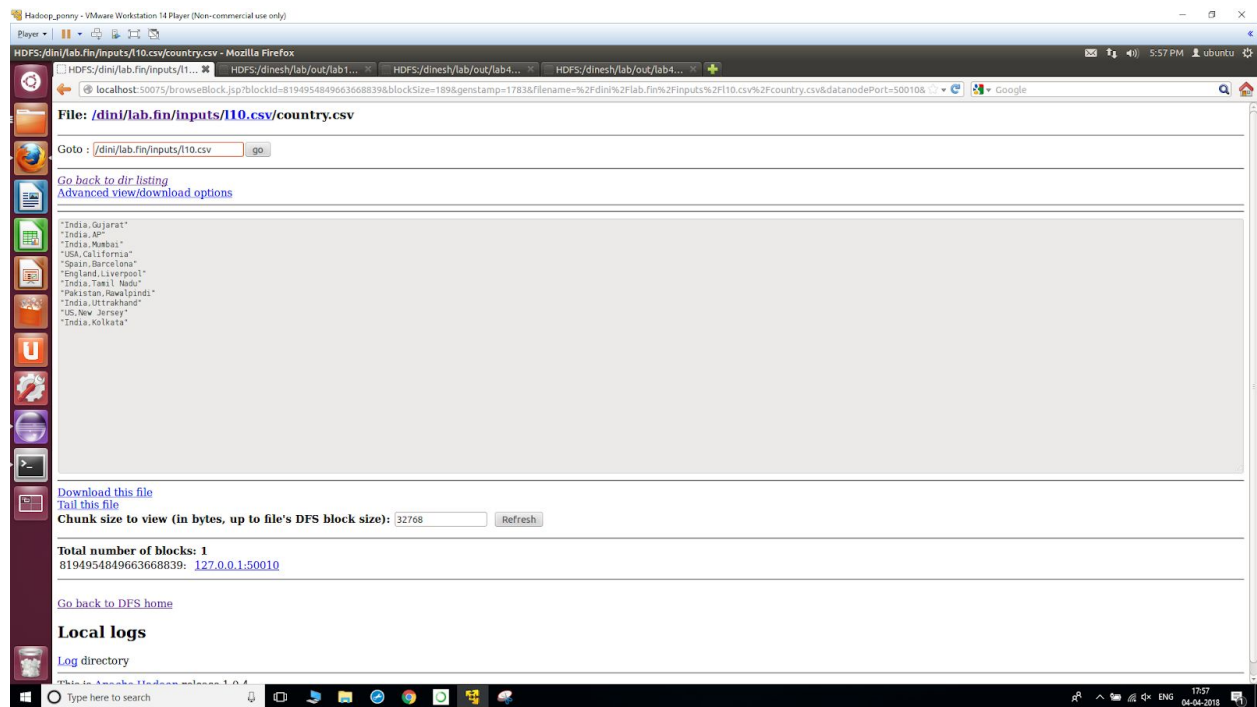
FileOutputFormat.setOutputPath(job, new Path(args[1]));

job.waitForCompletion(true);

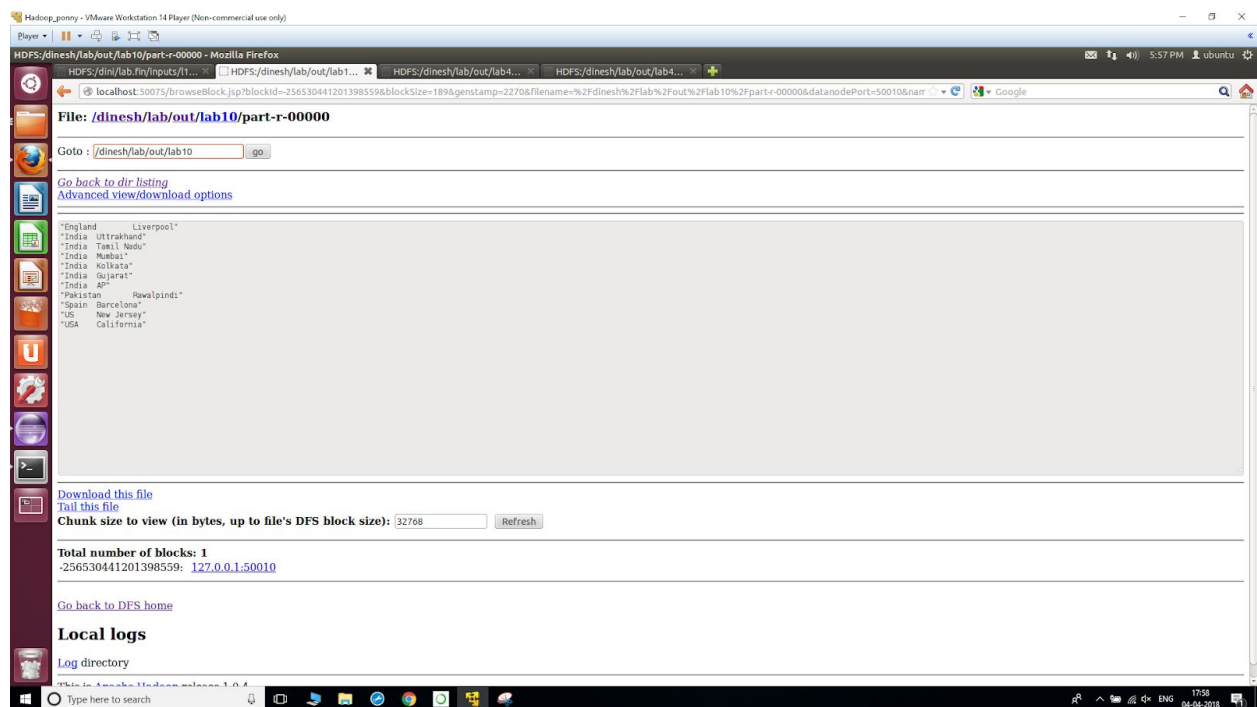
}

}
```

Input:-



Output:-



reference:-<https://dzone.com/articles/mapreduce-algorithms-%E2%80%93>