

<b>Name of the student:</b>		<b>Roll No.</b>	
<b>Practical Number:2</b>		<b>Date of Practical:</b>	
<b>Relevant CO's: ITC802.2</b>	<b>At the end of the course students will be able to use tools like hadoop and NoSQL to solve big data related problems.</b>		
<b>Sign here to indicate that you have read all the relevant material provided before attempting this practical</b>			<b>Sign:</b>

**Practical grading using Rubrics**

Indicator	Very Poor	Poor	Average	Good	Excellent
<b>Timeline</b> (2)	More than a session late (0)	NA	NA	NA	Early or on time (2)
<b>Code de- sign</b> (2)	N/A	Very poor code design with no comments and indentation(0.5)	Poor code design with very comments and indentation (1)	Design with good coding standards (1.5)	Accurate design with better coding standards (2)
<b>Performance</b> (4)	Unable to perform the experiment (0)	Able to partially perform the experiment (1)	Able to perform the experiment for certain use cases (2)	Able to perform the experiment considering most of the use cases (3)	Able to perform the experiment considering all use cases (4)
<b>Postlab</b> (2)	No Execution(0)	N/A	Partially Executed (1)	N/A	Fully Executed (2)

<b>Total Marks (10)</b>	<b>Sign of instructor</b>

# Practical

COURSE TITLE: BIG DATA ANALYTICS

COURSE TERM: 2021-2022

INSTRUCTOR NAME: SAURABH KULKARNI

---

**Problem Statement: Counting number of words in given text file using map reduce.**

**Theory: Explain the working of word count using map reduce with small example and diagrams**

---



**Code:****code for mapper:**

```
import java.io.IOException;
import java.util.StringTokenizer;

import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Mapper.Context;

public class WCMapper extends Mapper<LongWritable, Text, Text,
    ↪ IntWritable>
{
    // Create object of type Text to hold strings created per word
    ↪ of given document
```

**Code for Reducer:****Code for Driver Class:**

**PostLab: Find inverted index**

In this assignment you have to implement a simple map reduce job that builds an inverted index on the set of input documents. An inverted index maps each word to a list of documents that contain the word, and additionally records the position of each occurrence of the word within the document. For the purpose of this assignment, the position will be based on counting words, not characters.

Ex: Assume below are the input Documents.

file1="data is good."

file2="data is not good?"

Output:

data (file1,1)(file2,1)

good (file1,3)(file2,4)

is (file1,2)(file2,2)

not (file2,3)

For more details on inverted indices, you can check out the Wikipedia page on inverted indices.

Now in this assignment you need to implement above map-reduce job.

Input: A set of documents

Output:

Map: word1 (filename, position)

word2 (filename, position)

word1 (filename, position)

and so on for each occurrence of each word.

Reduce: word1 (filename, position)(filename,position)

word2 (filename, position)

and so on for each word.

Code for getting file name in Hadoop, which can be used in the Map function:

```
String filename=null;  
filename = ((FileSplit) context.getInputSplit()).  
    ↪ getPath().getName();
```

**Code for postlab question**

# 8695 | YASH SANKPAL | BATCH-C

## THEORY:

In Hadoop, **MapReduce** is a computation that decomposes large manipulation jobs into individual tasks that can be executed in parallel across a cluster of servers. The results of tasks can be joined together to compute final results.

MapReduce consists of 2 steps:

- **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)

## Work Flow of the Program

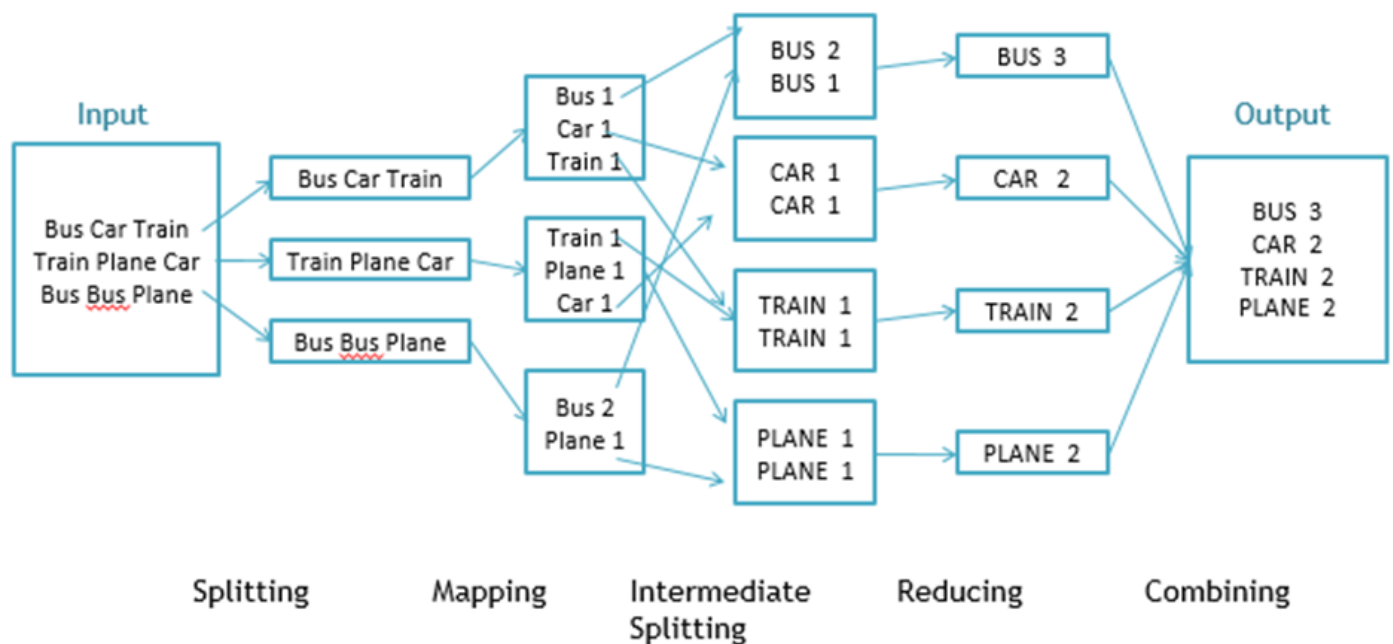


Fig. WorkFlow of MapReducing

Workflow of MapReduce consists of 5 steps:

1. **Splitting** – The splitting parameter can be anything, e.g. splitting by space, comma, semicolon, or even by a new line ('\n').
2. **Mapping** – as explained above.
3. **Intermediate splitting** – the entire process in parallel on different clusters. In order to group them in “Reduce Phase” the similar KEY data should be on the same cluster.
4. **Reduce** – it is nothing but mostly group by phase.
5. **Combining** – The last phase where all the data (individual result set from each cluster) is combined together to form a result.

## Code:

```
1 package exp2;
2
3 import java.io.IOException;
4 import java.util.StringTokenizer;
5
6 import org.apache.hadoop.io.IntWritable;
7 import org.apache.hadoop.io.LongWritable;
8 import org.apache.hadoop.io.Text;
9 import org.apache.hadoop.mapreduce.Mapper;
10 import org.apache.hadoop.mapreduce.Mapper.Context;
11
12 public class WCMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
13     // Create object of type Text to hold strings created per word of given document
14     Text word = new Text();
15     // Create final static variable of type IntWritable with value equal to 1 as according to algorithm map function puts 1 for each word encountered.
16     public static IntWritable one = new IntWritable(1);
17     //write a map function here
18     public void map(LongWritable ikey, Text ivalue, Context context)
19         throws IOException, InterruptedException {
20         //Define a String type variable and assign value which is equal to string equivalent of Text value passed to map function
21         String var = ivalue.toString();
22         //Convert this variable to tokens using StringTokenizer class as it separates out each word of the document
23         StringTokenizer st = new StringTokenizer(var);
24         //Until there are tokens in StringTokenizer,
25         while(st.hasMoreTokens()) {
26             // set the Text object created in first line of the code in WCMapper to next token in the StringTokenizer
27             word.set(st.nextToken());
28             context.write(word, one);
29         }
30     }
31     // Write this Text Variable and final static IntWritable Variable created to the context so that key-value pairs are generated.
32
33 }
34
35
36
```



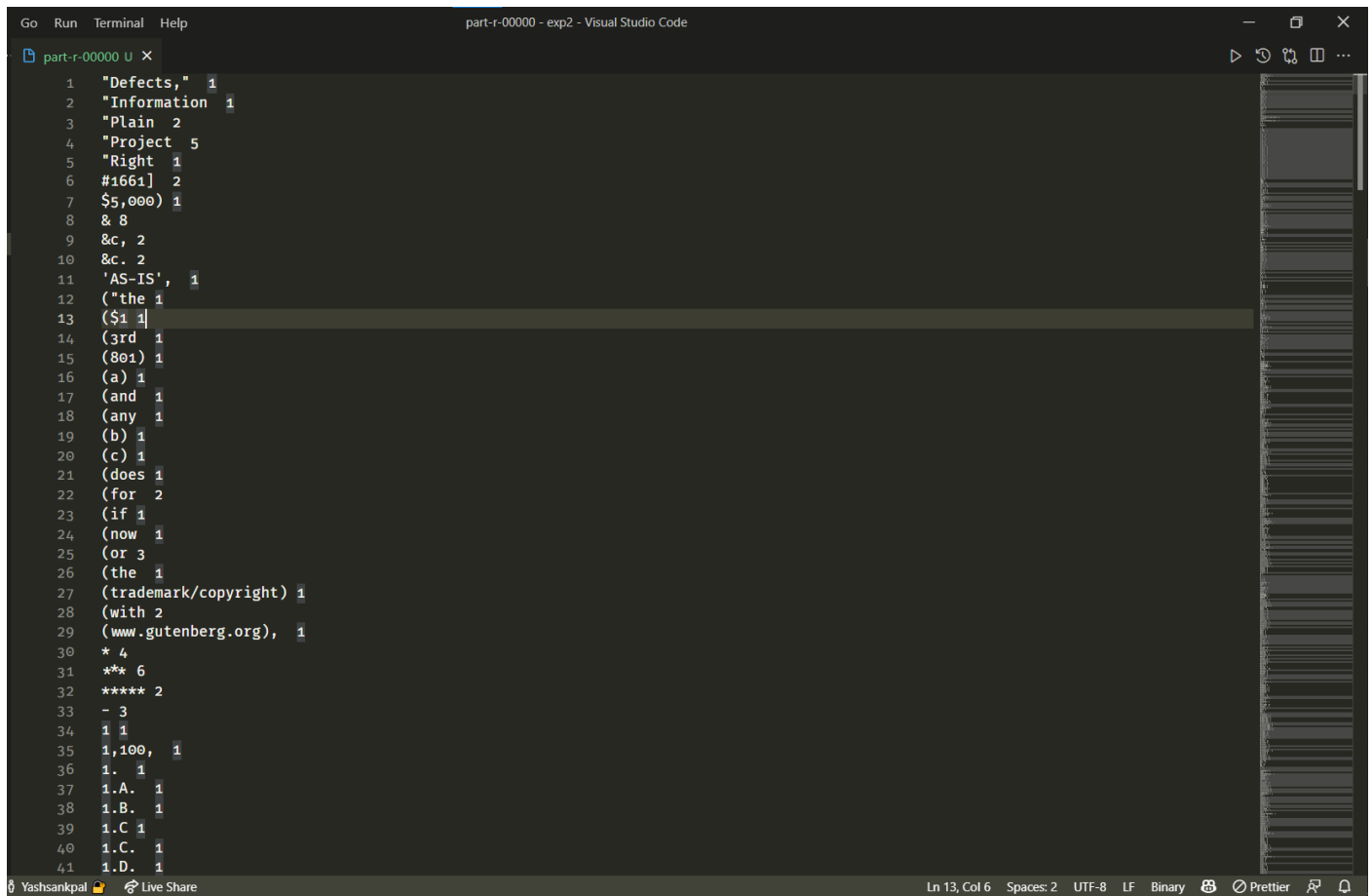
```
1 package exp2;
2
3 import java.io.IOException;
4
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapreduce.Reducer;
8 import org.apache.hadoop.mapreduce.Reducer.Context;
9
10
11 public class WCReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
12
13     public void reduce(Text key, Iterable<IntWritable> values, Context context)
14         throws IOException, InterruptedException {
15         // process values
16         //initialize sum=0
17         int sum = 0;
18         //Iterate over the collection of values to get count of each word i.e. key
19         for(IntWritable i: values) {
20             sum += i.get();
21         }
22         //Write this count to the context.
23         context.write(key, new IntWritable(sum));
24     }
25 }
26
27
28
```

```

1 package exp2;
2
3
4
5 import java.util.Date;
6 import java.util.Formatter;
7 import org.apache.hadoop.conf.Configuration;
8 import org.apache.hadoop.fs.Path;
9 import org.apache.hadoop.io.IntWritable;
10 import org.apache.hadoop.io.Text;
11 import org.apache.hadoop.mapreduce.Job;
12 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
13 import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
14 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
15 import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
16 import org.apache.hadoop.util.GenericOptionsParser;
17 public class WCDriver {
18
19     public static void main(String[] args) throws Exception {
20         Configuration conf = new Configuration();
21         GenericOptionsParser parser = new GenericOptionsParser(conf, args);
22         args = parser.getRemainingArgs();
23
24         //Job job = new Job(conf, "wordcount");
25         Job job=new Job(conf,"wordcount");
26         job.setJarByClass(WCDriver.class);
27
28         //job.setOutputKeyClass(Text.class);
29         // job.setOutputValueClass(IntWritable.class);
30         job.setMapperClass(WCMapper.class);
31         job.setReducerClass(WCReducer.class);
32
33         job.setMapOutputKeyClass(Text.class);
34         job.setMapOutputValueClass(IntWritable.class);
35
36         job.setInputFormatClass(TextInputFormat.class);
37         job.setOutputFormatClass(TextOutputFormat.class);
38
39         Formatter formatter = new Formatter();
40         String outpath = "Out"
41         + formatter.format("%1$tM%1$td%1$tH%1$tM%1$tS", new Date());
42         FileInputFormat.setInputPaths(job, new Path("hdfs://localhost:9000/input/wordcountdata.txt"));
43         FileOutputFormat.setOutputPath(job, new Path("hdfs://localhost:9000/output/exp2"));
44
45         System.out.println(job.waitForCompletion(true));
46     }
47
48 }
49

```

## OUTPUT:



```
1 "Defects," 1
2 "Information" 1
3 "Plain" 2
4 "Project" 5
5 "Right" 1
6 "#1661" 2
7 "$5,000" 1
8 "&" 8
9 "&c," 2
10 "&c." 2
11 'AS-IS', 1
12 ("the" 1
13 ("the" 1
14 (3rd 1
15 (801) 1
16 (a) 1
17 (and 1
18 (any 1
19 (b) 1
20 (c) 1
21 (does 1
22 (for 2
23 (if 1
24 (now 1
25 (or 3
26 (the 1
27 (trademark/copyright) 1
28 (with 2
29 (www.gutenberg.org), 1
30 * 4
31 *** 6
32 ***** 2
33 - 3
34 1 1
35 1,100, 1
36 1. 1
37 1.A. 1
38 1.B. 1
39 1.C 1
40 1.C. 1
41 1.D. 1
```

## POSTLAB:

### Inverted Index

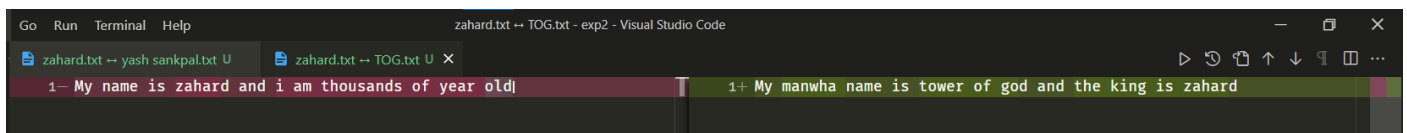
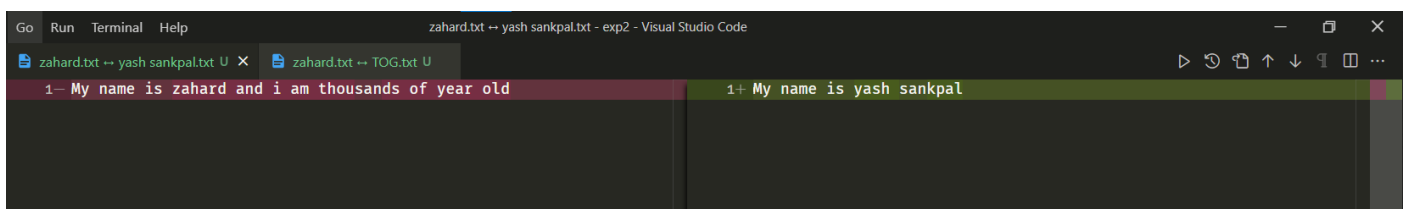
```
1 package exp2_postlab;
2
3 import java.io.IOException;
4 import java.util.StringTokenizer;
5
6 import org.apache.hadoop.io.LongWritable;
7 import org.apache.hadoop.io.Text;
8 import org.apache.hadoop.mapreduce.Mapper;
9 import org.apache.hadoop.mapreduce.lib.input.FileSplit;
10
11 public class WCMapper extends Mapper<LongWritable, Text, Text, Text> {
12     private Text nameKey = new Text();
13     private Text fileNameValue = new Text();
14
15     public void map(LongWritable ikey, Text ivalue, Context context) throws IOException, InterruptedException {
16         String str = ivalue.toString();
17         String[] st = str.split(" ");
18         String filename = null;
19
20         for(int i=0;i<st.length;i++) {
21             filename = ((FileSplit) context.getInputSplit()).getPath().getName();
22             nameKey.set(st[i]);
23             fileNameValue.set(" "+filename+", "+i+" ");
24             context.write(nameKey, fileNameValue);
25         }
26     }
27
28 }
29
```

```
1 package exp2_postlab;
2
3 import org.apache.hadoop.conf.Configuration;
4 import org.apache.hadoop.fs.Path;
5 import org.apache.hadoop.io.Text;
6 import org.apache.hadoop.mapreduce.Job;
7 import org.apache.hadoop.mapreduce.Mapper;
8 import org.apache.hadoop.mapreduce.Reducer;
9 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
10 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
11
12 public class WCDriver {
13
14     public static void main(String[] args) throws Exception {
15         Configuration conf = new Configuration();
16         Job job = Job.getInstance(conf, "JobName");
17         job.setJarByClass(WCDriver.class);
18         // TODO: specify a mapper
19         job.setMapperClass(WCMapper.class);
20         // TODO: specify a reducer
21         job.setReducerClass(WCReducer.class);
22
23         // TODO: specify output types
24         job.setOutputKeyClass(Text.class);
25         job.setOutputValueClass(Text.class);
26
27         // TODO: specify input and output DIRECTORIES (not files)
28         FileInputFormat.setInputPaths(job, new Path("hdfs://localhost:9000/input/exp2-postlab"));
29         FileOutputFormat.setOutputPath(job, new Path("hdfs://localhost:9000/exp2-postlab"));
30
31         if (!job.waitForCompletion(true))
32             return;
33         System.out.println(job.waitForCompletion(true));
34     }
35 }
36
37 }
```

```

1  package exp2_postlab;
2
3  import java.io.IOException;
4
5  import org.apache.hadoop.io.IntWritable;
6  import org.apache.hadoop.io.Text;
7  import org.apache.hadoop.mapreduce.Reducer;
8  import org.apache.hadoop.mapreduce.Reducer.Context;
9
10
11 public class WCReducer extends Reducer<Text, Text, Text, Text> {
12
13     public void reduce(Text key, Iterable<Text> values, Context context)
14         throws IOException, InterruptedException {
15         // process values
16         //initialize sum=0
17         //Iterate over the collection of values to get count of each word i.e. key
18         StringBuilder sb = new StringBuilder();
19
20         for(Text i: values) {
21             sb.append(i+" ");
22         }
23
24
25         //Write this count to the context.
26         context.write(key,new Text(sb.toString()));
27     }
28 }
29

```



## OUTPUT:



```
1 My ( yash sankpal.txt,0 ) ( TOG.txt,0 ) ( zahard.txt,0 )
2 am ( zahard.txt,6 )
3 and ( zahard.txt,4 ) ( TOG.txt,7 )
4 god ( TOG.txt,6 )
5 i ( zahard.txt,5 )
6 is ( zahard.txt,2 ) ( yash sankpal.txt,2 ) ( TOG.txt,10 ) ( TOG.txt,3 )
7 king ( TOG.txt,9 )
8 manwha ( TOG.txt,1 )
9 name ( TOG.txt,2 ) ( zahard.txt,1 ) ( yash sankpal.txt,1 )
10 of ( TOG.txt,5 ) ( zahard.txt,8 )
11 old ( zahard.txt,10 )
12 sankpal ( yash sankpal.txt,4 )
13 the ( TOG.txt,8 )
14 thousands ( zahard.txt,7 )
15 tower ( TOG.txt,4 )
16 yash ( yash sankpal.txt,3 )
17 year ( zahard.txt,9 )
18 zahard ( zahard.txt,3 ) ( TOG.txt,11 )
19
```