



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

Experiment No. 4
Experiment on Hadoop Map-Reduce
Date of Performance: 07/09/2023
Date of Submission: 14/09/2023



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

Aim: To write a program to implement a word count program using MapReduce.

Theory:

WordCount is a simple program which counts the number of occurrences of each word in a given text input data set. WordCount fits very well with the MapReduce programming model making it a great example to understand the Hadoop Map/Reduce programming style. The implementation consists of three main parts:

1. Mapper
2. Reducer
3. Driver

Step-1. Write a Mapper

A Mapper overrides the `map()` function from the Class "org.apache.hadoop.mapreduce.Mapper" which provides <key, value> pairs as the input. A Mapper implementation may output <key,value> pairs using the provided Context .

Input value of the WordCount Map task will be a line of text from the input data file and the key would be the line number <line_number, line_of_text> . Map task outputs <word, one> for each word in the line of text.

Pseudo-code

```
void Map (key,value){  
    for each word x in  
        value:  
            output.collect(x,1);  
}
```

Step-2. Write a Reducer



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

A Reducer collects the intermediate <key,value> output from multiple map tasks and assemble a single result. Here, the WordCount program will sum up the occurrence of each word to pairs as <word, occurrence>.

Pseudo-code

```
void Reduce (keyword, <list of value>){  
  for each x in <list of value>:  
    sum+=x;  
    final_output.collect(keyword, sum);  
}
```

Code:

```
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.Reducer  
; import
```



```
org.apache.hadoop.conf.Configuration;

import

org.apache.hadoop.mapreduce.Job;

import

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

import

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

import

org.apache.hadoop.mapreduce.lib.input.FileInputFormat;

import

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

import org.apache.hadoop.fs.Path;

public class WordCount

{

    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();

    StringTokenizer tokenizer = new

    StringTokenizer(line); while

    (tokenizer.hasMoreTokens()) {

    value.set(tokenizer.nextToken());

    context.write(value, new IntWritable(1));
```



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

```
}
```

```
}
```

```
}
```

```
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 x:
```

```
values)
```

```
{
```

```
sum+=x.get();
```

```
}
```

```
context.write(key, new IntWritable(sum));
```

```
}
```

```
}
```

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

```
Exception { Configuration conf= new
```

```
Configuration();
```

```
Job job = new Job(conf,"My Word Count
```

```
Program"); job.setJarByClass(WordCount.class);
```

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



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

```
job.setReducerClass(Reduce.class);
job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);
job.setInputFormatClass(TextInputFormat.class);
job.setOutputFormatClass(TextOutputFormat.class);
Path outputPath = new Path(args[1]);

//Configuring the input/output path from the filesystem into the job
FileInputFormat.addInputPath(job, new Path(args[0]));

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

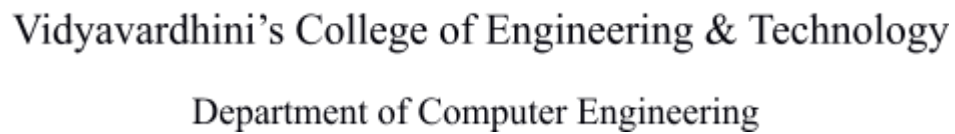
//deleting the output path automatically from hdfs so that we don't
have to delete it explicitly
outputPath.getFileSystem(conf).delete(outputPath);

//exiting the job only if the flag value becomes false
System.exit(job.waitForCompletion(true) ? 0 : 1);

}

}
```

Output:

[illegible]



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

```
1 package org.samarth;
2 import java.io.IOException;
3 import java.util.StringTokenizer;
4 import org.apache.hadoop.io.IntWritable;
5 import org.apache.hadoop.io.LongWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapred.MapReduceBase;
8 import org.apache.hadoop.mapred.Mapper;
9 import org.apache.hadoop.mapred.OutputCollector;
10 import org.apache.hadoop.mapred.Reporter;
11
12 public class WC_Mapper extends MapReduceBase implements Mapper<LongWritable,Text,Text,IntWritable> {
13     private final static IntWritable one = new IntWritable(1);
14     private Text word = new Text();
15
16     public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output,
17                     Reporter reporter) throws IOException {
18         String line = value.toString();
19         StringTokenizer tokenizer = new StringTokenizer(line);
20         while (tokenizer.hasMoreTokens()) {
21             word.set(tokenizer.nextToken());
22             output.collect(word, one);
23         }
24     }
25 }
```

```
1 package org.samarth;
2
3 import java.io.IOException;
4 import java.util.Iterator;
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapred.MapReduceBase;
8 import org.apache.hadoop.mapred.OutputCollector;
9 import org.apache.hadoop.mapred.Reducer;
10 import org.apache.hadoop.mapred.Reporter;
11
12 public class WC_Reducer extends MapReduceBase implements Reducer<Text,IntWritable,Text,IntWritable> {
13     public void reduce(Text key, Iterator<IntWritable> values,OutputCollector<Text,IntWritable> output,
14                       Reporter reporter) throws IOException {
15         int sum=0;
16         while (values.hasNext()) {
17             sum+=values.next().get();
18         }
19         output.collect(key,new IntWritable(sum));
20     }
21 }
```




Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

The screenshot shows an IDE with the following project structure on the left:

- Project: WordCount (C:\Users\admin\IdeaProjects\WordCount)
- src
 - main
 - java
 - org.samarth
 - WC_Mapper
 - WC_Reducer
 - WC_Runner
 - resources
 - test
- pom.xml
- External Libraries
- Scratches and Consoles

The main editor displays the `WC_Reducer.java` file with the following code:

```
1 package org.samarth;
2
3 import java.io.IOException;
4 import java.util.Iterator;
5 import org.apache.hadoop.io.IntWritable;
6 import org.apache.hadoop.io.Text;
7 import org.apache.hadoop.mapred.MapReduceBase;
8 import org.apache.hadoop.mapred.OutputCollector;
9 import org.apache.hadoop.mapred.Reducer;
10 import org.apache.hadoop.mapred.Reporter;
11
12 // usages
13 public class WC_Reducer extends MapReduceBase implements Reducer<Text, IntWritable, Text, IntWritable> {
14     public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text, IntWritable> output,
15         Reporter reporter) throws IOException {
16         int sum=0;
17         while (values.hasNext()) {
18             sum+=values.next().get();
19         }
20         output.collect(key, new IntWritable(sum));
21     }
22 }
```

The screenshot shows a Windows Command Prompt window with the following commands and output:

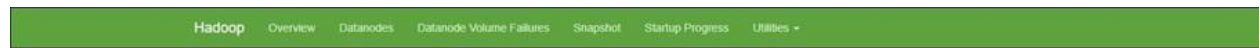
```
Microsoft Windows [Version 10.0.22000.2295]
(c) Microsoft Corporation. All rights reserved.

C:\Users\admin>cd Desktop
C:\Users\admin\Desktop>hadoop fs -mkdir /input
C:\Users\admin\Desktop>hadoop fs -put input.txt /input
C:\Users\admin\Desktop>
```



VidyaVardhini's College of Engineering & Technology

Department of Computer Engineering



Browse Directory

/input Go!   




Show 25 entries Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	rw-r--r--	admin	supergroup	36 B	Sep 13 04:53	1	128 MB	input.txt	


Showing 1 to 1 of 1 entries Previous 1 Next

Hadoop, 2022.

Browse Directory

/input Go!   

Show 25 entries Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	rw-r--r--	samarth	supergroup	36 B	Sep 13 04:53	1	128 MB	input.txt	

Showing 1 to 1 of 1 entries Previous 1 Next

Hadoop, 2022.

File information - input.txt

Download Head the file (first 32K) Tail the file (last 32K)

Block information — Block 0

Block ID: 1073741825
Block Pool ID: BP-1815554947-192.168.137.1-1695693903037
Generation Stamp: 1001
Size: 75
Availability:
• LAPTOP-K02APR2F.mshome.net

File contents

Hello World
Hello My name is Samarth Pandey I am Samarth
Welcome to World

Close



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering

```
Command Prompt
C:\Users\samar\Desktop>hadoop fs -mkdir /input
C:\Users\samar\Desktop>hadoop fs -put input.txt /input
C:\Users\samar\Desktop>hadoop jar C:\Users\samar\IdeaProjects\WordCount\target\hadoop-mapreduce-3.2.4.jar wordcount /input/input.txt /output
2023-09-29 18:57:09,119 INFO client.RetryProxy: Connecting to ResourceManager at /0.0.0.0:8032
2023-09-29 18:57:09,763 INFO mapreduce.JobResourceUploader: Disabling Erasure Coding for path: /tmp/hadoop-yarn/staging/samar/.staging/job_1695993949979_0001
2023-09-29 18:57:10,326 INFO input.FileInputFormat: Total input files to process : 1
2023-09-29 18:57:10,697 INFO mapreduce.JobSubmitter: number of splits:1
2023-09-29 18:57:11,007 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1695993949979_0001
2023-09-29 18:57:11,010 INFO mapreduce.JobSubmitter: Executing with tokens: []
2023-09-29 18:57:11,299 INFO conf.Configuration: resource-types.xml not found
2023-09-29 18:57:11,300 INFO resource.ResourceUtils: Unable to find 'resource-types.xml'.
2023-09-29 18:57:11,723 INFO impl.YarnClientImpl: Submitted application application_1695993949979_0001
2023-09-29 18:57:11,814 INFO mapreduce.Job: The url to track the job: http://LAPTOP-K02APR2F:8088/proxy/application_1695993949979_0001/
2023-09-29 18:57:11,816 INFO mapreduce.Job: Running job: job_1695993949979_0001
2023-09-29 18:57:22,135 INFO mapreduce.Job: Job job_1695993949979_0001 running in uber mode : false
2023-09-29 18:57:22,136 INFO mapreduce.Job: map 0% reduce 0%
2023-09-29 18:57:35,308 INFO mapreduce.Job: map 100% reduce 0%
2023-09-29 18:57:43,413 INFO mapreduce.Job: map 100% reduce 100%
2023-09-29 18:57:44,434 INFO mapreduce.Job: Job job_1695993949979_0001 completed successfully
2023-09-29 18:57:45,177 INFO mapreduce.Job: Counters: 54
File System Counters
  FILE: Number of bytes read=126
  FILE: Number of bytes written=478089
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=177
  HDFS: Number of bytes written=76
  HDFS: Number of read operations=8
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=2
  HDFS: Number of bytes read erasure-coded=0
Job Counters
  Launched map tasks=1
  Launched reduce tasks=1
  Data-local map tasks=1
  Total time spent by all maps in occupied slots (ms)=5488
  Total time spent by all reduces in occupied slots (ms)=5838
  Total time spent by all map tasks (ms)=5488
  Total time spent by all reduce tasks (ms)=5838
  Total vcore-milliseconds taken by all map tasks=5488
  Total vcore-milliseconds taken by all reduce tasks=5838
  Total megabyte-milliseconds taken by all map tasks=5619712
  Total megabyte-milliseconds taken by all reduce tasks=5978112
Map-Reduce Framework
  Map input records=1
```

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

Browse Directory

Show 25 entries

Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	<input type="checkbox"/>
<input type="checkbox"/>	drwxr-xr-x	samar	supergroup	0 B	Sep 29 18:56	0	0 B	input	<input type="checkbox"/>
<input type="checkbox"/>	drwxr-xr-x	samar	supergroup	0 B	Sep 29 18:57	0	0 B	output	<input type="checkbox"/>
<input type="checkbox"/>	drwx-----	samar	supergroup	0 B	Sep 29 18:57	0	0 B	tmp	<input type="checkbox"/>

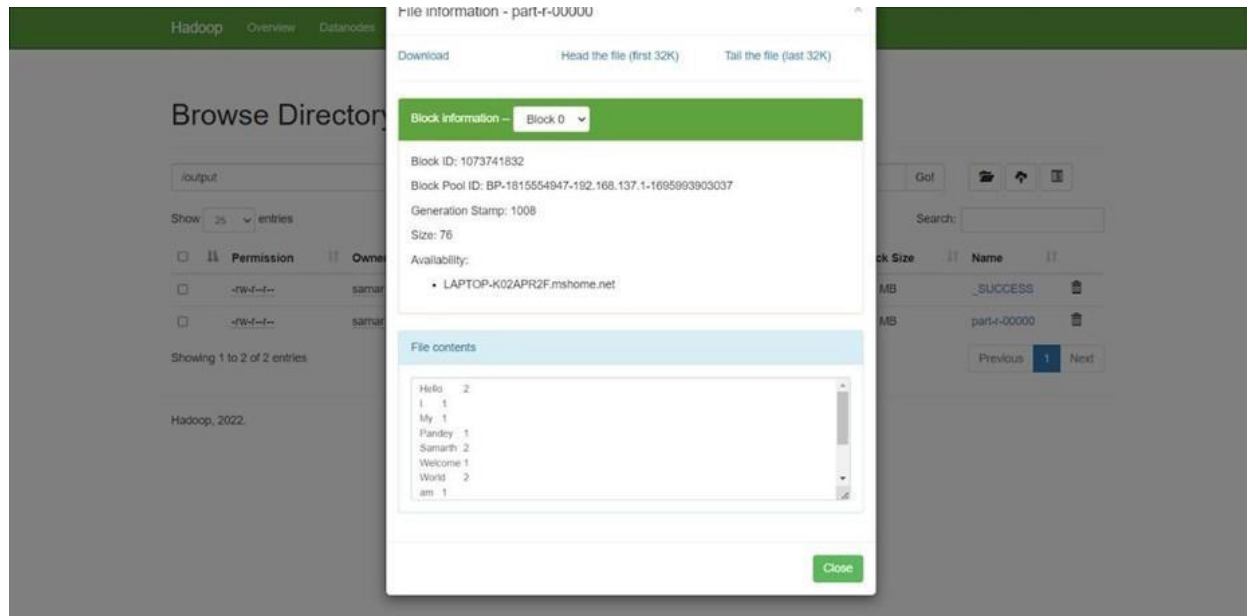
Showing 1 to 3 of 3 entries

2022



Vidyavardhini's College of Engineering & Technology

Department of Computer Engineering



Conclusion:

In this experiment, the primary emphasis was on the implementation of a Word Count program using the MapReduce paradigm, serving as a fundamental demonstration of Hadoop's data processing prowess. The Word Count program was deconstructed into its three crucial components: Mapper, Reducer, and Driver. The Mapper is responsible for handling the input data, breaking it into individual words, and emitting pairs in the form of <word, 1>. The Reducer plays a pivotal role in aggregating and summing these pairs to yield the ultimate word count. Meanwhile, the Driver assumes the role of orchestrating the MapReduce operation by configuring input/output paths and executing the tasks..