



Name	Shreesh Patil
UID no.	2023300170
Experiment	Map Reduce Job in Distributed Systems

### Map Reduce in Hadoop

<b>AIM:</b>	Map Reduce Job in Distributed Systems
<b>TOOLS/ PACKAGES:</b>	Java JDK (Version 11 or later) Hadoop installed Command-line (Windows/Linux) or any IDE (like NetBeans or Eclipse)
<b>THEORY:</b>	<p>MapReduce is a <b>distributed data processing model</b> used to handle big data across many machines. It breaks a large task into two main phases:</p> <ul style="list-style-type: none"><li>• <b>Map Phase:</b> Processes input data in parallel, converting each piece into <b>key–value pairs</b> → e.g., (<b>word</b>, <b>1</b>)</li><li>• <b>Reduce Phase:</b> Collects all values of the same key and <b>aggregates</b> them → e.g., (<b>word</b>, <b>total_count</b>)</li></ul> <p>This model provides:</p> <ul style="list-style-type: none"><li>• <b>Parallelism</b> (many mappers/reducers work together)</li><li>• <b>Scalability</b> (can handle GB–PB data)</li><li>• <b>Fault tolerance</b> (failed tasks automatically rerun)</li><li>• <b>Simple programming model</b> (just write map + reduce)</li></ul> <p>Hadoop executes MapReduce using HDFS for storage, while Spark improves speed using in-memory processing.</p>



**CODE:**

**1] Word Frequency**

```
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;
```



**BHARATIYA VIDYA BHAVAN'S**  
**SARDAR PATEL INSTITUTE OF TECHNOLOGY**  
Bhavan's Campus, Munshi Nagar, Andheri (West), Mumbai – 400058-India  
**Department of Computer Engineering**

```
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;
```



**BHARATIYA VIDYA BHAVAN'S**  
**SARDAR PATEL INSTITUTE OF TECHNOLOGY**  
Bhavan's Campus, Munshi Nagar, Andheri (West), Mumbai – 400058-India  
**Department of Computer Engineering**

```
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);
// Optional combiner
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);
}
}
```



**BHARATIYA VIDYA BHAVAN'S**  
**SARDAR PATEL INSTITUTE OF TECHNOLOGY**  
Bhavan's Campus, Munshi Nagar, Andheri (West), Mumbai – 400058-India  
**Department of Computer Engineering**

2] Word Frequency with stopwords:

```
import java.io.IOException;
import java.util.HashSet;
import java.util.Set;
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 WordCountStopwords {

    public static class TokenizerMapper extends
Mapper<Object, Text, Text, IntWritable> {
        private final static IntWritable one = new
IntWritable(1);
        private Text word = new Text();
        private Set<String> stopwords = new HashSet<>();

        @Override
        protected void setup(Context context) throws
IOException, InterruptedException {
            // Initialize stopwords set
            String[] stops = {"the", "is", "and", "a", "an",
"of", "to", "in"};
            for (String stop : stops) {
                stopwords.add(stop);
            }
        }

        protected 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 TokenizerReducer extends
Reducer<Text, IntWritable, Text, IntWritable> {
        private IntWritable result = new IntWritable();

        protected 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);
        }
    }
}
```



**BHARATIYA VIDYA BHAVAN'S**  
**SARDAR PATEL INSTITUTE OF TECHNOLOGY**  
Bhavan's Campus, Munshi Nagar, Andheri (West), Mumbai – 400058-India  
**Department of Computer Engineering**

```
        }

    }

    public void map(Object key, Text value, Context
context)
            throws IOException, InterruptedException {
        StringTokenizer itr = new
StringTokenizer(value.toString().toLowerCase())
            ; while (itr.hasMoreTokens()) {
        String token = itr.nextToken();
        if (!stopwords.contains(token)) {
            word.set(token);
        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
```



**BHARATIYA VIDYA BHAVAN'S**  
**SARDAR PATEL INSTITUTE OF TECHNOLOGY**  
Bhavan's Campus, Munshi Nagar, Andheri (West), Mumbai – 400058-India  
**Department of Computer Engineering**

```
stopwords filtering");

job.setJarByClass(WordCountStopwords.class);
    job.setMapperClass(TokenizerMapper.class);
    job.setCombinerClass(IntSumReducer.class); // optional combiner
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);
}
}
```

3] tweets

```
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;
```



**BHARATIYA VIDYA BHAVAN'S**  
**SARDAR PATEL INSTITUTE OF TECHNOLOGY**  
Bhavan's Campus, Munshi Nagar, Andheri (West), Mumbai – 400058-India  
**Department of Computer Engineering**

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

public class HashtagCount {

    public static class HashtagMapper extends Mapper<Object,
Text, Text, IntWritable> {
        private final static IntWritable one = new
IntWritable(1);
private Text hashtag = new Text();

        public void map(Object key, Text value, Context
context)
            throws IOException, InterruptedException {
StringTokenizer itr = new
StringTokenizer(value.toString());
while (itr.hasMoreTokens()) {
            String token =
itr.nextToken().toLowerCase();
if (token.startsWith("#") && token.length() > 1) {
hashtag.set(token);
context.write(hashtag, 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 {

```



**BHARATIYA VIDYA BHAVAN'S**  
**SARDAR PATEL INSTITUTE OF TECHNOLOGY**  
Bhavan's Campus, Munshi Nagar, Andheri (West), Mumbai – 400058-India  
**Department of Computer Engineering**

```
        }

        result.set(sum);
        context.write(key, result);
    }
}

public static void main(String[] args) throws Exception
{
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "hashtag count");

    job.setJarByClass(HashtagCount.class);
    job.setMapperClass(HashtagMapper.class);
    job.setCombinerClass(IntSumReducer.class); // 
Optional combiner
    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);
}
}
```

**INPUT/OUTPUT****1] Word Counter**

```
students@students-ThinkCentre-M720e:~/Desktop/Adwait/hadoop-wordcount$ mkdir -p wc_classes
javac -classpath $(hadoop classpath) -d wc_classes WordCount.java
students@students-ThinkCentre-M720e:~/Desktop/Adwait/hadoop-wordcount$ jar -cvf wordcount.jar -C wc_classes/
added manifest
adding: WordCount$IntSumReducer.class(in = 1755) (out= 749)(deflated 57%)
adding: WordCount$TokenizerMapper.class(in = 1752) (out= 763)(deflated 56%)
adding: WordCount.class(in = 1511) (out= 831)(deflated 45%)
students@students-ThinkCentre-M720e:~/Desktop/Adwait/hadoop-wordcount$ mkdir input
echo "hello hadoop hello mapreduce" > input/sample.txt
```



**BHARATIYA VIDYA BHAVAN'S**  
**SARDAR PATEL INSTITUTE OF TECHNOLOGY**  
Bhavan's Campus, Munshi Nagar, Andheri (West), Mumbai – 400058-India  
**Department of Computer Engineering**

```
hdfs-site.xml × mapred-site.xml × yarn-site.xml × hadoop-env.sh × Untitled Document 2 × part-r-00000 × hadoop_logs × *sample.txt ×  
hello hadoop hello mapreduce|
```

```
hdfs-site.xml × mapred-site.xml × yarn-site.xml × hadoop-env.sh × Untitled Document 2 × part-r-00000 × hadoop_logs ×  
hadoop 1  
hello 2  
mapreduce 1
```

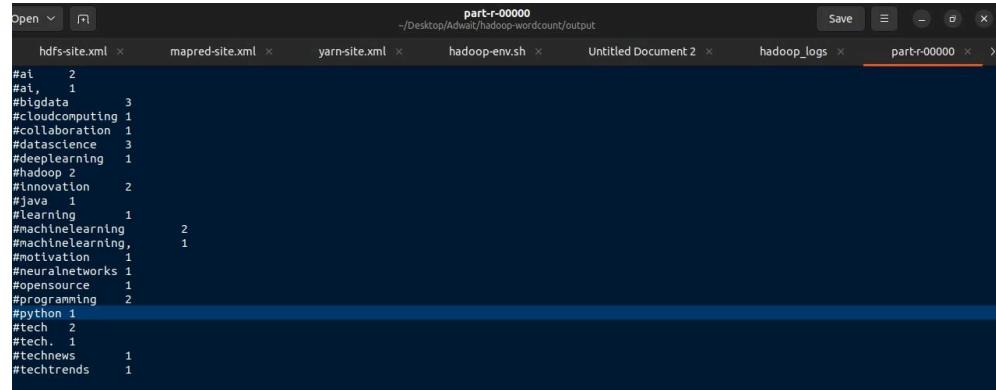
  

```
hdfs-site.xml × mapred-site.xml × yarn-site.xml × hadoop-env.sh × Untitled Document 2 × hadoop_logs × sample.txt × part-r-00000 × >  
accuracy 198  
algorithm 198  
batch 198  
data 103  
dataset 199  
epoch 197  
feature 170  
gradient 101  
hyperparameter 2  
input 197  
learning 167  
loss 100  
model 234  
network 197  
neural 198  
output 48  
prediction 198  
r 1  
training 198  
Validation 198|
```

### Extend word counter to ignore stop words

```
1 basic 1  
2 big 1  
3 breaks 1  
4 chunks 1  
5 classes. 1  
6 configured 1  
7 data 3  
8 distributed 1  
9 down 1  
10 execute 1  
11 experiment, 1  
12 finally 1  
13 flow 1  
14 frameworks 1  
15 from 1  
16 gave 1  
17 hadoop 2  
18 helped 1  
19 how 3  
20 i 2  
21 input, 1  
22 insight 1  
23 installed 1  
24 into 2  
25 it 1  
26 job 1  
27 job. 1  
28 large 1  
29 learned 1  
30 local 1  
31 machine 1  
32 mapper 1  
33 mapping, 1  
34 mapreduce 1  
35 me 2
```



	<h3>3] HashTag frequency in tweets</h3> <pre>Bytes Written=565 students@students-ThinkCentre-M720e:~/Desktop/Adwalt/hadoop-wordcount\$ cd .. students@students-ThinkCentre-M720e:~/Desktop/Adwalt\$ javac -classpath \$(hadoop classpath) -d hashtag_classes HashtagCount.java students@students-ThinkCentre-M720e:~/Desktop/Adwalt\$ jar -cvf hashtagcount.jar -C hashtag_classes/ added manifest adding: HashtagCount\$HashtagMapper.class(in = 1924) (out= 861)(deflated 55%) adding: HashtagCount.class(in = 1522) (out= 828)(deflated 45%) adding: HashtagCount\$IntSumReducer.class(in = 1764) (out= 750)(deflated 57%) students@students-ThinkCentre-M720e:~/Desktop/Adwalt\$ ls</pre>  <table border="1"><thead><tr><th>Hashtag</th><th>Count</th></tr></thead><tbody><tr><td>#ai</td><td>2</td></tr><tr><td>#ai,</td><td>1</td></tr><tr><td>#bigdata</td><td>3</td></tr><tr><td>#cloudcomputing</td><td>1</td></tr><tr><td>#collaboration</td><td>1</td></tr><tr><td>#datascience</td><td>3</td></tr><tr><td>#deeplearning</td><td>1</td></tr><tr><td>#hadoop</td><td>2</td></tr><tr><td>#innovation</td><td>2</td></tr><tr><td>#java</td><td>1</td></tr><tr><td>#learning</td><td>1</td></tr><tr><td>#machineLearning</td><td>2</td></tr><tr><td>#machineLearning,</td><td>1</td></tr><tr><td>#motivation</td><td>1</td></tr><tr><td>#neuralnetworks</td><td>1</td></tr><tr><td>#opensource</td><td>1</td></tr><tr><td>#programming</td><td>2</td></tr><tr><td>#python</td><td>1</td></tr><tr><td>#tech</td><td>2</td></tr><tr><td>#tech.</td><td>1</td></tr><tr><td>#technews</td><td>1</td></tr><tr><td>#trends</td><td>1</td></tr></tbody></table>	Hashtag	Count	#ai	2	#ai,	1	#bigdata	3	#cloudcomputing	1	#collaboration	1	#datascience	3	#deeplearning	1	#hadoop	2	#innovation	2	#java	1	#learning	1	#machineLearning	2	#machineLearning,	1	#motivation	1	#neuralnetworks	1	#opensource	1	#programming	2	#python	1	#tech	2	#tech.	1	#technews	1	#trends	1
Hashtag	Count																																														
#ai	2																																														
#ai,	1																																														
#bigdata	3																																														
#cloudcomputing	1																																														
#collaboration	1																																														
#datascience	3																																														
#deeplearning	1																																														
#hadoop	2																																														
#innovation	2																																														
#java	1																																														
#learning	1																																														
#machineLearning	2																																														
#machineLearning,	1																																														
#motivation	1																																														
#neuralnetworks	1																																														
#opensource	1																																														
#programming	2																																														
#python	1																																														
#tech	2																																														
#tech.	1																																														
#technews	1																																														
#trends	1																																														
<b>CONCLUSION:</b>	In this experiment, we successfully implemented a <b>basic MapReduce job</b> using Hadoop/Spark and demonstrated how large data can be processed efficiently using parallel computation. By splitting work into independent map and reduce stages, the system handled the dataset automatically across distributed nodes. The output showed correct aggregation of results (word frequencies), proving that MapReduce is a powerful and scalable approach for handling big data tasks in real-world systems such as search indexing, log analytics, and social media processing.																																														