(/)

**REFCARDS (/REFCARDZ) TREND REPORTS (/TRENDREPORTS)** (/users/login.html)

Culture and Methodologies
(/culture-and-methodologies)

Data Engineering
(/data-engineering)

Software Design and Architecture
(/software-design-and-architecture)

Coding
(/coding)

Testing, Deployment, and Maintenance
(/testing-deployment-and-maintenance)

# RELATED

Redis-Based Tomcat Session Management (/articles/redis-based-tomcat-session-management)

Optimizing Java Applications for AWS Lambda (/articles/java-apps-aws-lambda)

Buildpacks: An Open-Source Alternative to Chainguard (/articles/buildpacks-open-source-alternative-to-chainguard)

Efficient Asynchronous Processing Using CyclicBarrier and CompletableFuture in Java (/articles/efficient-asynchronous-processing-using-cyclicbarr)

# Partner Resources

![DZone](/) (/) **REFCARDS (/REFCARDZ) REPORTS (/TREND-REPORTS) EVENTS (/EVENTS)** (/users/login.html)

Culture and Methodologies (/culture-and-methodologies)

Data Engineering (/data-engineering)

Software Design and Architecture (/software-design-and-architecture)

Coding (/coding)

Testing, Deployment, and Maintenance (/testing-deployment-and-maintenance)

Culture and Methodologies          Data Engineering          Software Design and Architecture          Coding          Testing, Deployment, and Maintenance
(/culture-and-          (/data-          (/software-design-and-          (/coding)          (/testing-deployment-and-
methodologies)          engineering)          architecture)          maintenance)

# Word Count Program With MapReduce and Java

In this post, we provide an introduction to the basics of MapReduce, along with a tutorial to create a word count app using Hadoop and Java.

By   Shital Kat (/users/2752223/shitalkatkar.html) · Mar. 03, 16 · Tutorial

In Hadoop, MapReduce (https://dzone.com/articles/mapreduce-design-patterns-1) 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) |
  |---|---|---|
  |  |  | (BUS,7), |

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.

# Now Let's See the Word Count Program in Java

Fortunately, we don't have to write all of the above steps, we only need to write the splitting parameter, Map function logic, and Reduce function logic. The rest of the remaining steps will execute automatically.

Make sure that Hadoop is installed on your system with the Java SDK.

DZone® (/) REFCARDS (/REFCARDZ) REPORTS (/TREND-REPORTS) ZONES (/ZONES) (/users/login.html)

Culture and Methodologies (/culture-and-methodologies) | Data Engineering (/data-engineering) | Software Design and Architecture (/software-design-and-architecture) | Coding (/coding) | Testing, Deployment, and Maintenance (/testing-deployment-and-maintenance)

## Steps

1. Open Eclipse> File > New > Java Project >( Name it - MRProgramsDemo) > Finish.

2. Right Click > New > Package ( Name it - PackageDemo) > Finish.

3. Right Click on Package > New > Class (Name it - WordCount).

4. Add Following Reference Libraries:

    1. Right Click on Project > Build Path> Add External

        1. */usr/lib/hadoop-0.20/**hadoop-core.jar***

        2. *Usr/lib/hadoop-0.20/lib/**Commons-cli-1.2.jar***

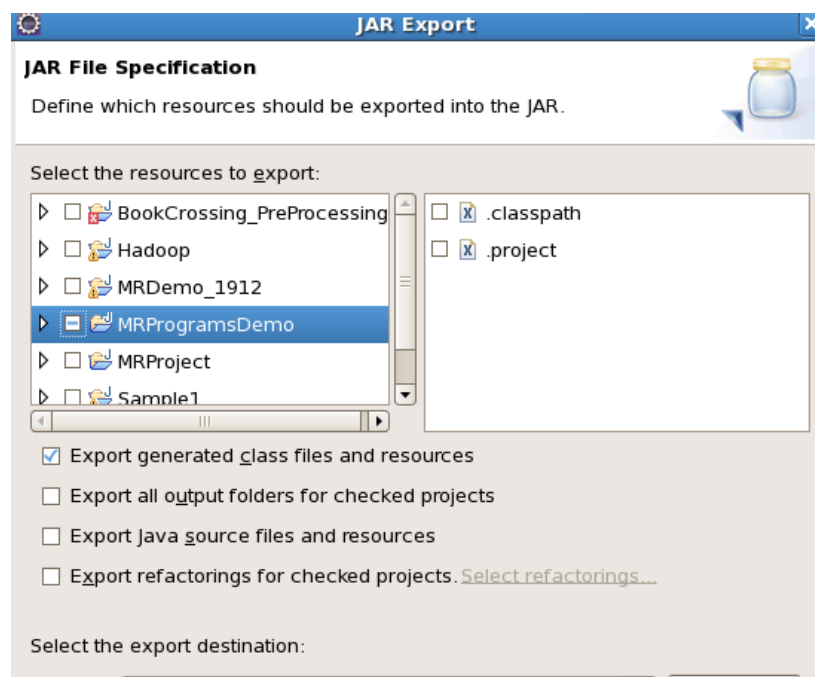5. Type the following code:

```
1  package PackageDemo;
2
3  import java.io.IOException;
4  import org.apache.hadoop.conf.Configuration;
5  import org.apache.hadoop.fs.Path;
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.Job;
10 import org.apache.hadoop.mapreduce.Mapper;
11 import org.apache.hadoop.mapreduce.Reducer;
12 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
13 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
14 import org.apache.hadoop.util.GenericOptionsParser;
15
16
17
18
19 public class WordCount {
20 public static void main(String [] args) throws Exception
21 {
22 Configuration c=new Configuration();
23 String[] files=new GenericOptionsParser(c,args).getRemainingArgs();
24 Path input=new Path(files[0]);
25 Path output=new Path(files[1]);
26 Job j=new Job(c,"wordcount");
27 j.setJarByClass(WordCount.class);
28 j.setMapperClass(MapForWordCount.class);
29 j.setReducerClass(ReduceForWordCount.class);
30 j.setOutputKeyClass(Text.class);
31 j.setOutputValueClass(IntWritable.class);
32 FileInputFormat.addInputPath(j, input);
33 FileOutputFormat.setOutputPath(j, output);
34 System.exit(j.waitForCompletion(true)?0:1);
35 }
36 public static class MapForWordCount extends Mapper<LongWritable, Text, Text, IntWritable>{
37 public void map(LongWritable key, Text value, Context con) throws IOException, InterruptedException
38 {
```
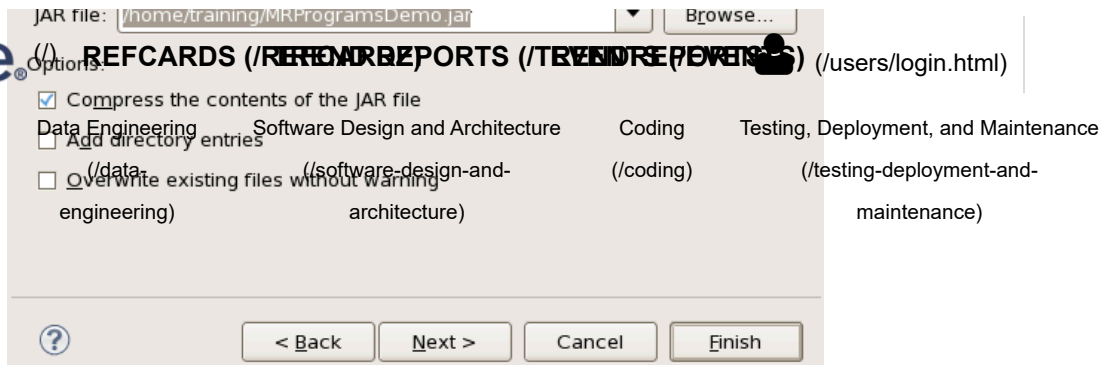
REFCARDS (/REFCARDZ) REPORTS (/TRENDREPORTS) (/users/login.html)

Culture and Methodologies     Data Engineering     Software Design and Architecture     Coding     Testing, Deployment, and Maintenance

(/culture-and-               (/data-                (/software-design-and-                (/coding)               (/testing-deployment-and-
methodologies)               engineering)           architecture)                                                 maintenance)

```
39  String line = value.toString();
40  String[] words = line.split(" ");
41  for(String word: words )
42  {
43      Text outputKey = new Text(word.toUpperCase().trim());
44      IntWritable outputValue = new IntWritable(1);
45      con.write(outputKey, outputValue);
46  }
47  }
48  }
49
50  public static class ReduceForWordCount extends Reducer<Text, IntWritable, Text, IntWritable>
51  {
52  public void reduce(Text word, Iterable<IntWritable> values, Context con) throws IOException, InterruptedEx
53  {
54  int sum = 0;
55     for(IntWritable value : values)
56     {
57     sum += value.get();
58     }
59     con.write(word, new IntWritable(sum));
60  }
61
62  }
63
64  }
```

The above program consists of three classes:

- Driver class (Public, void, static, or main; this is the entry point).

- The `Map` class which **extends** the public class Mapper<KEYIN,VALUEIN,KEYOUT,VALUEOUT> and implements the `Map` function.

- The `Reduce` class which extends the public class Reducer<KEYIN,VALUEIN,KEYOUT,VALUEOUT> and implements the `Reduce` function.

6. Make  a jar file

Right Click on Project> Export> Select export destination as **Jar File** > next> Finish.

DZone
REFCARDS (/REFCARDZ) REPORTS (/TRENDREPORTS) (/users/login.html)

Culture and Methodologies       Data Engineering       Software Design and Architecture       Coding       Testing, Deployment, and Maintenance
(/culture-and-                   (/data-                (/software-design-and-                   (/coding)     (/testing-deployment-and-
methodologies)                   engineering)           architecture)                                          maintenance)

7. Take a text file and move it into HDFS format:

To move this into Hadoop directly, open the terminal and enter the following commands:

```
1 [training@localhost ~]$ hadoop fs -put wordcountFile wordCountFile
```

8. Run the jar file:

*(Hadoop jar jarfilename.jar packageName.ClassName  PathToInputTextFile PathToOutputDirectry)*

```
1 [training@localhost ~]$ hadoop jar MRProgramsDemo.jar PackageDemo.WordCount wordCountFile MRDir1
```

9. Open the result:

```
1 [training@localhost ~]$ hadoop fs -ls MRDir1
2
3 Found 3 items
4
5 -rw-r--r--   1 training supergroup          0 2016-02-23 03:36 /user/training/MRDir1/_SUCCESS
6 drwxr-xr-x   - training supergroup          0 2016-02-23 03:36 /user/training/MRDir1/_logs
7 -rw-r--r--   1 training supergroup         20 2016-02-23 03:36 /user/training/MRDir1/part-r-00000
```

```
1 [training@localhost ~]$ hadoop fs -cat MRDir1/part-r-00000
2 BUS     7
3 CAR     4
4 TRAIN   6
```

Hadoop       MapReduce       Java (Programming Language)

Opinions expressed by DZone contributors are their own.

---

# RELATED

DZone.

Redis-Based Tomcat Session Management (/refcardz/...) REFCARDS (/refcardz) REFCARDZ TREND REPORTS (/trendreports) TREND REPORTS EVENTS (/events) (/users/login.html)

Optimizing Java Applications for AWS Lambda
(/culture-and-methodologies)
Culture and Methodologies Data Engineering (/data-engineering) Software Design and Architecture (/software-design-and-architecture) Coding (/coding) Testing, Deployment, and Maintenance (/testing-deployment-and-maintenance)

Buildpacks: An Open-Source Alternative to Chainguard

Efficient Asynchronous Processing Using CyclicBarrier and CompletableFuture in Java