


DZone (/) > Big Data Zone (/big-data-analytics-tutorials-tools-news) > Word Count Program With MapReduce and Java

Word Count Program With MapReduce and Java

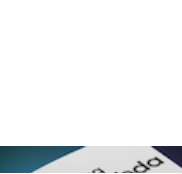
 (/users/2752223/shitalkatkar.html) by **Shital Kat** (/users/2752223/shitalkatkar.html) · Mar. 03, 16 · Big Data Zone (/big-data-analytics-tutorials-tools-news) · Tutorial

 Like (31)

 Comment (5)

 Save

 Tweet

 O'Reilly | Build and Deploy Server Applications with Java

Learn to build

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	Set of Tuples	(Bus,1), (Car,1), (bus,1), (car,1), (train,1), (car,1), (bus,1), (car,1), (train,1), (bus,1),
-------	---------------	--

 Output of Map function) REFCARDZ (/refcardz) TREND REPORTS (/trendreports) WEBINARS (/webinars) ZONES		(TRAIN,1), (BUS,1), (bus,1), (car,1), (CAR,1), (/search) (/users/login.html) (car,1), (BUS,1), (TRAIN,1)
Output	Converts into smaller set of tuples	(BUS,7), (CAR,7), (TRAIN,4)

Work Flow of the Program

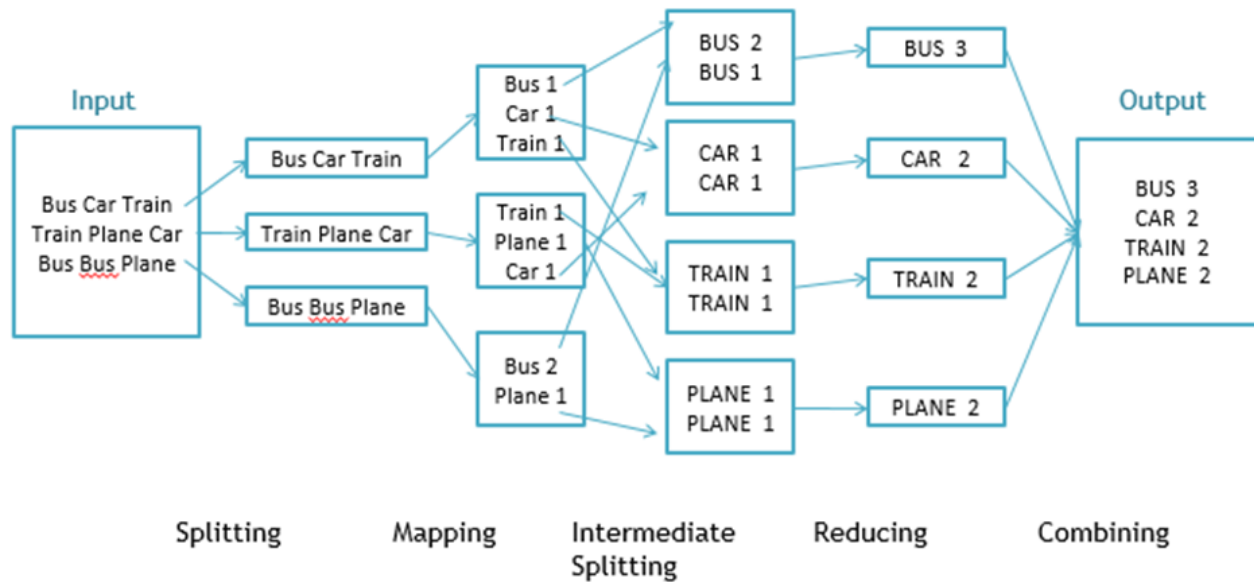


Fig. WorkFlow of MapReducing

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.

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.



(/users/login.html)



(/search)

REFCARDZ (/refcardz)

TREND REPORTS (/trendreports)

WEBINARS (/webinars)

ZONES ▾

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         {
39             String line = value.toString();
40             String[] words=line.split(",");

```


[DZone \(/\)](#)


[\(/users/login.html\)](/users/login.html)


[\(/search\)](/search)

[REFCARDZ \(/refcardz\)](#)
[TREND REPORTS \(/trendreports\)](#)
[WEBINARS \(/webinars\)](#)
[ZONES ▾](#)

```

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, InterruptedException
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 }

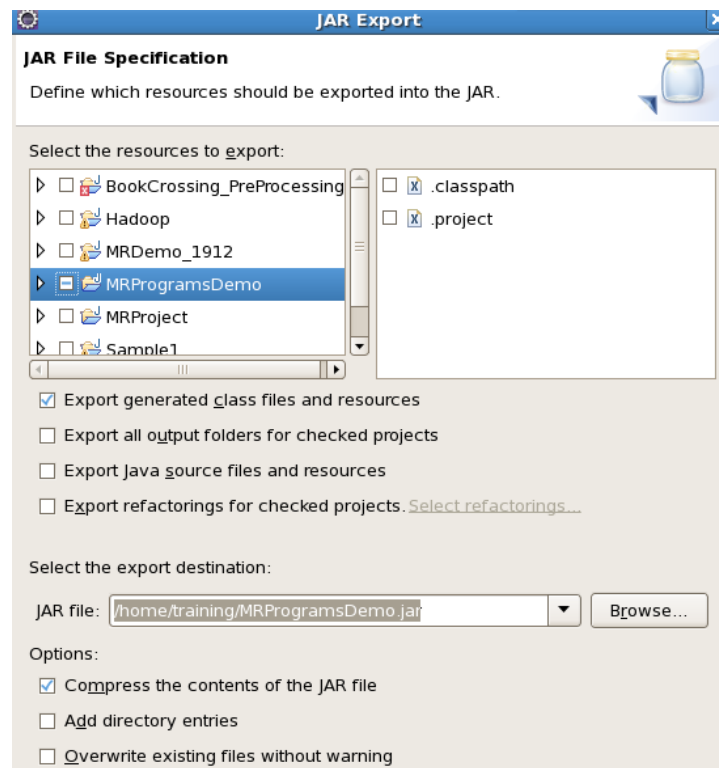
```

The above program consists of three classes:

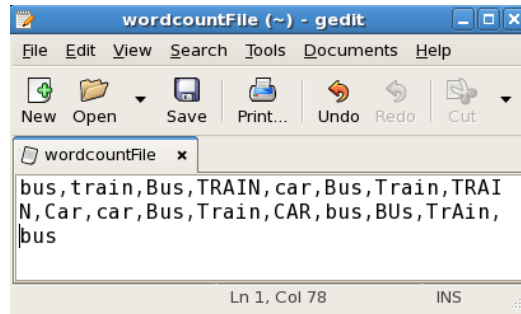
- Driver class (Public, void, static, or main; this is the entry point).
- The Map class which **extends** the public class Mapper<KEYIN,VALUEIN,KEYOUT,VALUEOUT> and implements the Map function.
- The Reduce class which extends the public class Reducer<KEYIN,VALUEIN,KEYOUT,VALUEOUT> and implements the Reduce function.

6. Make a jar file

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



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 filename.jar packageName.ClassName PathToInputTextFile PathToOutputDirectory)

```
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
```

O'Reilly | What is Distributed SQL?





Traditional SQL struggles to scale. NoSQL struggles with consistency. Meet the r of the database: Distributed SQL. [Free Book](#) ▶

Topics: MAPREDUCE, JAVA, HADOOP, BIG DATA, TUTORIAL, WORDCOUNT

Opinions expressed by DZone contributors are their own.

Popular on DZone


- [The Most Comprehensive Guide on WebRTC \(/articles/webrtc-comprehensive-guide?fromrel=true\)](/articles/webrtc-comprehensive-guide?fromrel=true)

 **DZone** (/) [Method Builder With Lombok @Builder \(/articles/method-builder-with-lombok-builder-from-rel=true\)](#)  [\(/search\)](#)

RECOMMENDED (/refcardz) TREND REPORTS (/trendreports) WEBINARS (/webinars) ZONES

- [Spring Boot Delete User Details API Test Client Using Rest Assured \(/articles/spring-boot-delete-user-details-api-test-client-us?fromrel=true\)](#)
- [MongoDB to Couchbase for Developers, Part 1: Architecture \(/articles/mongodb-to-couchbase-for-developers-part-1-archite-1?fromrel=true\)](#)

Big Data Partner Resources

 **wrike**

Deliver more of your best work, faster, with Wrike.

Wrike is an award-winning work management software that enables teams to plan and track projects, collaborate in real-time, and automate reports. [Start for free today!](#) ►

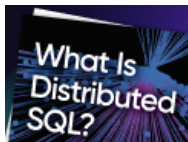
Presented by **Wrike**



Learn to build Serverless Apps with Google Cloud Run

Never babysit infrastructure again. Learn apps. [Free O'Reilly Book](#) ►

Presented by **Cockroach Labs**



O'Reilly | What is Distributed SQL?

Traditional SQL struggles to scale. NoSQL struggles with consistency. Meet the next evolution of the database: Distributed SQL. [Free O'Reilly Book](#) ►

Presented by **Cockroach Labs**



Monday.com helps teams work more efficiently to execute projects that deliver results on time. [Try for free](#) ►

Presented by **Monday.com**



Learn to build & deploy serverless Java apps with AWS Lambda

Completely free: All 10 chapters of O'Reil AWS Lambda [Free O'Reilly Book](#) ►

Presented by **Cockroach Labs**



All Your Work In One Place

Asana helps teams orchestrate their work, from daily tasks to strategic initiatives. With Asana, teams are more confident, move faster, and accomplish more with less, no matter where they are located. [Try for free](#) ►

Presented by **Asana**

ABOUT US

About DZone (/pages/about)

Send feedback (mailto:support@dzone.com)

Careers (https://devada.com/careers/)

Sitemap (/sitemap)

ADVERTISE

Advertise with DZone (https://advertise.dzone.com)

CONTRIBUTE ON DZONE

 [Article Submission Guidelines \(/articles/dzones-article-submission-guidelines\)](#)

[MVB Program \(/pages/mvb\)](#)



[\(/users/login.html\)](#)



[\(/search\)](#)

REFCARDZ [\(/refcardz\)](#) **TRENDREPORTS** [\(/trendreports\)](#) **WEBINARS** [\(/webinars\)](#) **ZONES** ▾

[Visit the Writers' Zone \(/writers-zone\)](#)

LEGAL

[Terms of Service \(/pages/tos\)](#)

[Privacy Policy \(/pages/privacy\)](#)

CONTACT US




600 Park Offices Drive

Suite 300

Durham, NC 27709

support@dzone.com (<mailto:support@dzone.com>)

+1 (919) 678-0300 (<tel:+19196780300>)

Let's be friends:    

[\(/pages/help/https://www.dzone.com/dzonecompany/dzone/\)](#)

DZone.com is powered by



[\(https://devada.com/answerhub/\)](https://devada.com/answerhub/)