# WordCount-Example-in-Hadoop

## Enabling technologies for Data Science

**2017msbda008**
**Samiksha Agarwal**

August 6, 2018
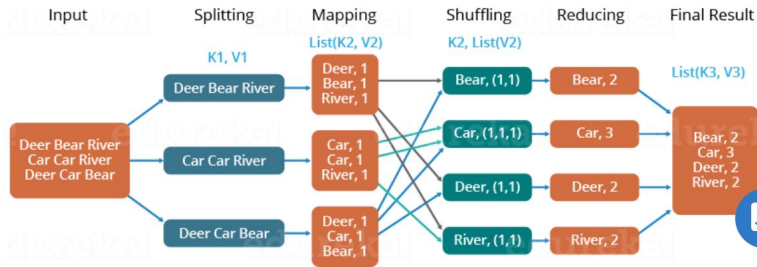
# Running a Map/Reduce job

# What is MapReduce?

MapReduce is a programming framework that allows us to perform distributed and parallel processing on large data sets in a distributed environment.

Let us understand, how a MapReduce works by taking an example where I have a text file whose contents are as follows:

**Dear, Bear, River, Car, Car, River, Deer, Car and Bear**

Now, suppose, we have to perform a word count using MapReduce. So, we will be finding the unique words and the number of occurrences of those unique words.
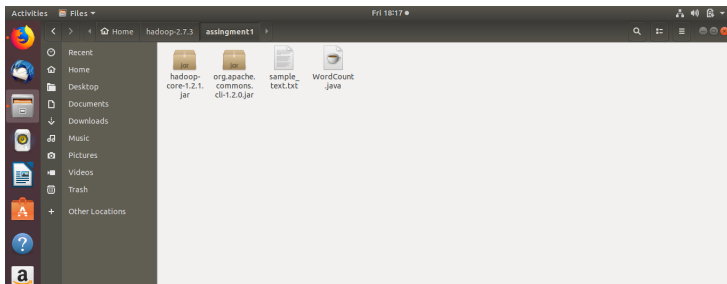
The Overall MapReduce Word Count Process

we will understand this process by program and for details of each step you can follow this link Click Here

# Steps for run the program :

**First, start all daemons**
Command: ./start-all.sh

- Step 1: Download the hadoop-1.2.1.jar.zip file Click Here
- Step 2: Download org.apache.commons.cli-1.2.0.jar.zip Click Here
- Step 3: Create a folder (assingment1) and extract both zip file here.
- Step 4: Make java file (WordCount.java) and for making this file you can prefer this link Click Here
- Step 5: Make a text file.

- Step 5: crtl+alt+t (open terminal).
- Step 6: Before you run the sample, you must create input and output locations in HDFS. Use the following commands to create the input directory /user/wordcount/input in HDFS:

  <span style="color:red">Command: hadoop fs -mkdir /home/samiksha/wordcount</span>
  <span style="color:red">Command: hadoop fs -chown wordcount /home/samiksha/wordcount</span>
  <span style="color:red">Command: hadoop fs -mkdir /home/samiksha/wordcount /input</span>
  <span style="color:red">Command: hadoop fs -mkdir /home/samiksha/wordcount /output</span>

```
samiksha@samiksha-HP-Pavilion-Notebook:~$ hadoop fs -mkdir /home/samiksha/wordcount
samiksha@samiksha-HP-Pavilion-Notebook:~$ hadoop fs -chown wordcount /home/samiksha/wordcount
samiksha@samiksha-HP-Pavilion-Notebook:~$ hadoop fs -mkdir /home/samiksha/wordcount/input
samiksha@samiksha-HP-Pavilion-Notebook:~$ hadoop fs -mkdir /home/samiksha/wordcount/output
```

  you can check your input and output file

  <span style="color:red">Command: hadoop fs -ls /home/samiksha/wordcount</span>

```
samiksha@samiksha-HP-Pavilion-Notebook:~$ hadoop fs -ls /home/samiksha/wordcountFound 2 items
drwxr-xr-x   - samiksha supergroup          0 2018-07-27 18:12 /home/samiksha/wordcount/input
drwxr-xr-x   - samiksha supergroup          0 2018-07-27 18:12 /home/samiksha/wordcount/output
```

- Step 7: move your sample_text.txt to /home/samiksha/wordcount /input directory in HDFS.
  Command: hadoop fs -put /home/samiksha/hadoop-2.7.3/assingment1/sample_text.txt /home/samiksha/wordcount /input

  `samiksha@samiksha-HP-Pavilion-Notebook:~$ hadoop fs -put /home/samiksha/hadoop-2.7.3/assingment1/sample_text.txt /home/samiksha/wordcount/inp ut`

- Step 8: open WordCount.java and You can use an appropriate package for your domain, or keep the generic version.
  package org.myorg.WordCount;



```
package org.myorg.WordCount;

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;
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);
      }
```

- step 9: now, go to the folder assingment1 and open terminal.
- Step 10: Compile the WordCount class.
  Command: mkdir -p build
  Command: javac -cp
  hadoop-core-1.2.1.jar:org.apache.commons.cli-1.2.0.jar
  WordCount.java -d build -Xlint



- step 11: Create a JAR file for the WordCount application.
  Command: jar -cvf wordcount.jar -C build/ .

you can check jar file into the assingment1 folder :



- step 12: Run the WordCount application from the JAR file, passing the paths to the input and output directories in HDFS.
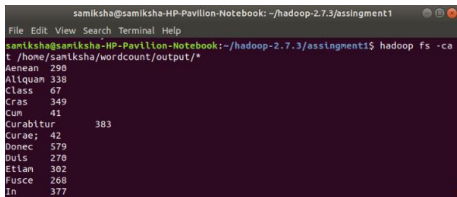  Command: hadoop jar wordcount.jar
  org.myorg.WordCount.WordCount /home/samiksha/wordcount/input
  /home/samiksha/wordcount/output

- step 12: When you look at the output, The number of occurrences from all input files has been reduced to a single sum for each word.
  Command: hadoop fs -cat /home/samiksha/wordcount/output/*



- Step 13: If you want to run the sample again, you first need to remove the output directory. Use the following command.
  Command: hadoop fs -rm -r /home/samiksha/wordcount/output
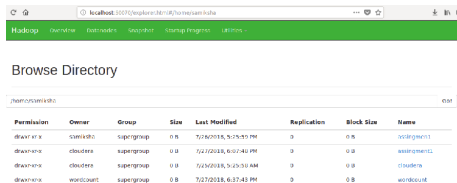  **Don't forget to stop all daemons**
  Command: ./stop-all.sh

**if you want check the output file**
go to Utilities → Browse the file system



now, go to home → samiksha → wordcount

here is your input and output file



Browse Directory

/home/samiksha/wordcount     Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| drwxr-xr-x | samiksha | supergroup | 0 B | 7/27/2018, 6:15:38 PM | 0 | 0 B | input |
| drwxr-xr-x | samiksha | supergroup | 0 B | 7/27/2018, 6:37:53 PM | 0 | 0 B | output |

now, go to output $\rightarrow$



Browse Directory

/home/samiksha/wordcount/output     Go!

| Permission | Owner | Group | Size | Last Modified | Replication | Block Size | Name |
|---|---|---|---|---|---|---|---|
| -rw-r--r-- | samiksha | supergroup | 0 B | 7/27/2018, 6:37:53 PM | 1 | 128 MB | _SUCCESS |
| -rw-r--r-- | samiksha | supergroup | 5.02 KB | 7/27/2018, 6:37:53 PM | 1 | 128 MB | part-r-00000 |

open part-r-0000 $\rightarrow$ download the file

# **YARN Web UI**

(for more details you can follow this link Click Here )

# Thank You