# Source Code – WordCount.java

```java
1.  package org.myorg;
2.
3.  import java.io.*;
4.  import java.util.*;
5.
6.  import org.apache.hadoop.fs.Path;
7.  import org.apache.hadoop.filecache.DistributedCache;
8.  import org.apache.hadoop.conf.*;
9.  import org.apache.hadoop.io.*;
10. import org.apache.hadoop.mapred.*;
11. import org.apache.hadoop.util.*;
12.
```

```
13.  public class WordCount extends Configured
     implements Tool {

14.

15.   public static class Map extends MapReduceBase implements
      Mapper<LongWritable, Text, Text, IntWritable> {

16.

17.     static enum Counters { INPUT_WORDS }

18.

19.     private final static IntWritable one = new
        IntWritable(1);

20.     private Text word = new Text();

21.

22.     private boolean caseSensitive = true;

23.     private Set<String> patternsToSkip = new
        HashSet<String>();
```

```
24.
25.    private long numRecords = 0;
26.    private String inputFile;
27.
28.    public void configure(JobConf job) {
29.      caseSensitive =
       job.getBoolean("wordcount.case.sensitive", true);
30.      inputFile = job.get("map.input.file");
31.
32.      if (job.getBoolean("wordcount.skip.patterns", false)) {
33.        Path[] patternsFiles = new Path[0];
34.        try {
35.          patternsFiles =
       DistributedCache.getLocalCacheFiles(job);
36.        } catch (IOException ioe) {
```

```java
37.        System.err.println("Caught exception while getting
   cached files: " + StringUtils.stringifyException(ioe));
38.      }
39.    for (Path patternsFile : patternsFiles) {
40.      parseSkipFile(patternsFile);
41.    }
42.     }
43.   }
44.
45.   private void parseSkipFile(Path patternsFile) {
46.    try {
47.     BufferedReader fis = new BufferedReader(new FileReader
   (patternsFile.toString()));
48.     String pattern = null;
49.     while ((pattern = fis.readLine()) != null) {
```

```java
50.      patternsToSkip.add(pattern);
51.     }
52.    } catch (IOException ioe) {
53.      System.err.println("Caught exception while parsing the
   cached file '" + patternsFile + "' : " +
   StringUtils.stringifyException(ioe));
54.    }
55.   }
56.
57.   public void map(LongWritable key, Text value,
   OutputCollector<Text, IntWritable> output, Reporter
   reporter) throws IOException {
58.     String line = (caseSensitive) ? value.toString() :
   value.toString().toLowerCase();
59.
```

```java
60.    for (String pattern : patternsToSkip) {
61.      line = line.replaceAll(pattern, "");
62.    }
63.
64.    StringTokenizer tokenizer = new StringTokenizer(line);
65.    while (tokenizer.hasMoreTokens()) {
66.      word.set(tokenizer.nextToken());
67.      output.collect(word, one);
68.      reporter.incrCounter(Counters.INPUT_WORDS, 1);
69.    }
70.
71.    if ((++numRecords % 100) == 0) {
72.      reporter.setStatus("Finished processing " + numRecords
+ " records " + "from the input file: " + inputFile);
```

```java
73.     }

74.    }

75.   }

76.

77.   public static class Reduce extends MapReduceBase
      implements Reducer<Text, IntWritable, Text, IntWritable> {

78.     public void reduce(Text key, Iterator<IntWritable>
        values, OutputCollector<Text, IntWritable> output,
        Reporter reporter)
        throws IOException {

79.       int sum = 0;

80.       while (values.hasNext()) {

81.         sum += values.next().get();

82.       }

83.       output.collect(key, new IntWritable(sum));
```

```
84.    }

85.    }

86.

87.  public int run(String[] args) throws Exception {

88.    JobConf conf = new JobConf(getConf(), WordCount.class);

89.    conf.setJobName("wordcount");

90.

91.    conf.setOutputKeyClass(Text.class);

92.    conf.setOutputValueClass(IntWritable.class);

93.

94.    conf.setMapperClass(Map.class);

95.    conf.setCombinerClass(Reduce.class);

96.    conf.setReducerClass(Reduce.class);

97.
```

```java
98.    conf.setInputFormat(TextInputFormat.class);
99.    conf.setOutputFormat(TextOutputFormat.class);
100.
101.   List<String> other_args = new ArrayList<String>();
102.   for (int i=0; i < args.length; ++i) {
103.    if ("-skip".equals(args[i])) {
104.     DistributedCache.addCacheFile(new
Path(args[++i]).toUri(), conf);
105.      conf.setBoolean("wordcount.skip.patterns", true);
106.    } else {
107.     other_args.add(args[i]);
108.    }
109.   }
110.
```

```java
111.    FileInputFormat.setInputPaths(conf, new Path
        (other_args.get(0)));
112.    FileOutputFormat.setOutputPath(conf, new Path
        (other_args.get(1)));
113.
114.    JobClient.runJob(conf);
115.    return 0;
116.  }
117.
118.  public static void main(String[] args) throws Exception {
119.    int res = ToolRunner.run(new Configuration(),
        new WordCount(), args);
120.    System.exit(res);
121.  }
122. }
```

## Sample Runs

Sample text-files as input:

```
$ bin/hadoop dfs -ls /usr/joe/wordcount/input/
/usr/joe/wordcount/input/file01
/usr/joe/wordcount/input/file02

$ bin/hadoop dfs -cat /usr/joe/wordcount/input/file01
Hello World, Bye World!

$ bin/hadoop dfs -cat /usr/joe/wordcount/input/file02
Hello Hadoop, Goodbye to hadoop.
```

Run the application:

```
$ bin/hadoop jar /usr/joe/wordcount.jar org.myorg.WordCount
/usr/joe/wordcount/input /usr/joe/wordcount/output
```

**Output:**

```
$ bin/hadoop dfs -cat /usr/joe/wordcount/output/part-00000
Bye 1
Goodbye 1
Hadoop, 1
Hello 2
World! 1
World, 1
hadoop. 1
to 1
```

Notice that the inputs differ from the first version we looked at, and how they affect the outputs.

Now, lets plug-in a pattern-file which lists the word-patterns to be ignored, via the `DistributedCache`.

```
$ hadoop dfs -cat /user/joe/wordcount/patterns.txt
\.
\,
\!
to
```

## Run it again, this time with more options:

```
$ bin/hadoop jar /usr/joe/wordcount.jar org.myorg.WordCount
-Dwordcount.case.sensitive=true /usr/joe/wordcount/input
/usr/joe/wordcount/output -skip
/user/joe/wordcount/patterns.txt
```

## As expected, the output:

```
$ bin/hadoop dfs -cat /usr/joe/wordcount/output/part-00000
Bye 1
Goodbye 1
Hadoop 1
Hello 2
World 2
hadoop 1
```

## Run it once more, this time switch-off case-sensitivity:

```
$ bin/hadoop jar /usr/joe/wordcount.jar org.myorg.WordCount
-Dwordcount.case.sensitive=false /usr/joe/wordcount/input
/usr/joe/wordcount/output -skip
/user/joe/wordcount/patterns.txt
```

Sure enough, the output:

```
$ bin/hadoop dfs -cat /usr/joe/wordcount/output/part-00000
bye 1
goodbye 1
hadoop 2
hello 2
world 2
```

**Highlights**

The second version of `WordCount` improves upon the previous one by using some features offered by the MapReduce framework:

- Demonstrates how applications can access configuration parameters in the `configure` method of the `Mapper` (and `Reducer`) implementations (lines 28-43).

- Demonstrates how the `DistributedCache` can be used to distribute read-only data needed by the jobs. Here it allows the user to specify word-patterns to skip while counting (line 104).

- Demonstrates the utility of the `Tool` interface and the `GenericOptionsParser` to handle generic Hadoop command-line options (lines 87-116, 119).

- Demonstrates how applications can use `Counters` (line 68) and how they can set application-specific status information via the `Reporter` instance passed to the `map` (and `reduce`) method (line 72).

*Java and JNI are trademarks or registered trademarks of Sun Microsystems, Inc. in the United States and other countries.*