

# Breaking down MapReduce Concepts

## 1: Datatypes

Hadoop uses a wrapper around java datatypes. For eg: IntWritable is used for declaring an integer instead of int.

Similarly, Conversion needed from Java to Hadoop:

- \* byte -> ByteWritable
- \* short -> ShortWritable
- \* long -> LongWritable
- \* float -> FloatWritable
- \* double -> DoubleWritable
- \* boolean -> BooleanWritable
- \* string -> Text

Explore more dataTypes : <https://hadoop.apache.org/docs/r2.6.2/api/org/apache/hadoop/io/>

Also have a look at : <http://stackoverflow.com/questions/19441055/why-does-hadoop-need-classes-like-text-or-intwritable-instead-of-string-or-integ>

Using these hadoop datatypes are easy. Let us look at an example.

1\* To create a IntWritable from a int.

Method: public IntWritable(int value)

Eg: IntWritable var1 = new IntWritable(5);

2\* To set value to a declared IntWritable variable.

Method: public void set (int value)

Eg: IntWritable var1 = new IntWritable();  
var1.set(5);

3\* To get the value stored in the IntWritable variable.

Method: public int get()

Eg: var1.get(); // output will be 5;

Explore more methods at:

<https://hadoop.apache.org/docs/r2.6.2/api/org/apache/hadoop/io/IntWritable.html>

## 2: Understanding the mapreduce components

Note: Start reading from the main method for better understanding.

```
1 import java.io.IOException;
2 import java.util.StringTokenizer;
3
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.Text;
8 import org.apache.hadoop.mapreduce.Job;
9 import org.apache.hadoop.mapreduce.Mapper;
10 import org.apache.hadoop.mapreduce.Reducer;
11 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
12 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
13
14 /* The following class has the implementation of Mapper,Combiner,Reducer and the main method. For simplicity,
15 we will use the same class for Reducer and Combiner.
16 */
17 public class WordCount {
18
19     /* Extend the custom class with Mapper<IpKey, IpValue, OpKey, OpValue >
20     IpKey ->InputKey to mapper class (Object in most cases)
21     IpValue ->InputValue to mapper class (Text in most cases)
22     OpKey ->OutputKey from mapper class (Will be the InputKey for Reducer)
23     OpValue ->OutputValue from mapper class(Will be InputValue for Reducer)
24     */
25     public static class TokenizerMapper
26         extends Mapper<Object, Text, Text, IntWritable>{
27
28         //Create new IntWritable object
29         private final static IntWritable one = new IntWritable(1);
30
31         //Create new Text object
32         private Text word = new Text();
33
34         /*User logic is placed inside map function.
35         Output (Key,value) pair is written to context in every stage,
36         context facilitates the data movement from stage to another stage.
37         Eg: from mapper class to reducer class
38         */
39         public void map(Object key, Text value, Context context)
40             throws IOException, InterruptedException {
41
42             /*Splits strings based on a token. Words separated by spaces in a sentence are divided into an array of words
43             here. Iterate through the array of words and assign value '1'to every word. Here word is the key and 1 is the value
44             for the corresponding key.
45             */
46             StringTokenizer itr = new StringTokenizer(value.toString());
47             while (itr.hasMoreTokens()) {
48                 word.set(itr.nextToken());
49
50                 // context.write(key,value)
51                 context.write(word, one);
52             }
53         }
54     }
55 }
56
57
58
59
60
61
```

```

62
63 /* Extend the custom class with Reducer<IpKey, IpValue, OpKey, OpValue >
64 IpKey ->InputKey to reducer class (OutputKey from mapper class)
65 IpValue ->InputValue to reducer class (OuputValue from mapper class)
66 OpKey ->OutputKey from reducer class (Will be written to file)
67 OpValue ->OutputValue from reducer class(Will be written to file)
68 */
69 public static class IntSumReducer
70     extends Reducer<Text,IntWritable,Text,IntWritable>{
71     private IntWritable result = new IntWritable();
72
73     /*User logic is placed inside reducefunction.
74     Output (Key,value) pair is written to context.
75     DataType of key in reduce is Text because outputKey from mapper is a Text. So the dataType must be
76 changed according to the mapper output.
77     DataType of Value is Iterable<IntWritable>because outputValue from mapper is a IntWritable. As a particular
78 key can have multiple values,
79 we adopt Iterable<InputValue DataType>as the standard.
80 */
81     public void reduce(Text key, Iterable<IntWritable>values,
82         Context context
83         ) throws IOException, InterruptedException {
84         int sum = 0;
85
86         //For every key, sum the value
87         for (IntWritable val : values) {
88             sum += val.get();
89         }
90         result.set(sum);
91         context.write(key, result);
92     }
93 }
94
95
96
97 public static void main(String[] args) throws Exception {
98     //Create an object of org.apache.hadoop.conf.Configuration
99     Configuration conf = new Configuration();
100
101     /* Create an object of org.apache.hadoop.mapreduce.Job.
102     Creates a new Job with a given jobName. The Job makes a copy of the Configuration so that any necessary
103 internal modifications do not reflect on the incoming parameter.
104 */
105     Job job = Job.getInstance(conf, "word count");
106
107     // Set the class name in which main method resides.
108     job.setJarByClass(WordCount.class);
109     // Set the mapper class name
110     job.setMapperClass(TokenizerMapper.class);
111
112     //Set the combiner class name
113     job.setCombinerClass(IntSumReducer.class);
114
115     //Set the reducer class name
116     job.setReducerClass(IntSumReducer.class);
117
118     /*Output from reducer is the final result and Output is Key,Value pair.
119     Inform Hadoop about the dataType of Key and Value from reducer
120 */
121     job.setOutputKeyClass(Text.class);
122     job.setOutputValueClass(IntWritable.class);
123
124
125

```

```

126  /* The following command is used to run a mapreduce program.
127      hadoop jar wc.jar WordCount input_dir output_dir
128      hadoop jar <jar name><pgm class name><input path><output path>
129      The input and output path to the dataset are set as below
130  */
131  FileInputFormat.addInputPath(job, new Path(args[0]));
132  FileOutputFormat.setOutputPath(job, new Path(args[1]));
133
134  /* Wait for the program execution to complete.
135      Exit the program after job execution.
136  */
137  System.exit(job.waitForCompletion(true) ? 0 : 1);
  }
}

```

### 3: Understanding the mapreduce workflow with an example

