

Big Data 4

Tutorial 3: MapReduce

Dr. Nikos Ntarmos

Room M111, Lilybank Gardens
<nikos.ntarmos@glasgow.ac.uk>

School of Computing Science
University of Glasgow

17 October, 2013



Outline

- 1 JobTracker UI
 - General information
 - Job information
 - Task information

- 2 Hadoop programming
 - Hadoop data types
 - Basic Mapper/Reducer methods
 - Custom input/output formats
 - Word Count galore



JobTracker UI :: General information



bigdata-06 Hadoop Map/Reduce Administration

State: RUNNING

Started: Wed Oct 16 17:28:19 BST 2013

Version: 2.9.0-mr1-cdh4.4.0, Unknown

Compiled: Tue Sep 3 19:47:44 PDT 2013 by Jenkins from Unknown

Identifier: 201310161728

Cluster Summary (Heap Size is 81.06 MB/4.20 GB)

Running Map Tasks	Running Reduce Tasks	Total Submissions	Nodes	Occupied Map Slots	Occupied Reduce Slots	Reserved Map Slots	Reserved Reduce Slots	Map Task Capacity	Reduce Task Capacity	Avg. Tasks/Node	Blacklisted Nodes	Excluded Nodes
51	0	3	6	51	0	0	0	52	26	13.00	0	0

Scheduling Information

Queue Name	State	Scheduling Information
default	running	N/A

Filter (Jobid, Priority, User, Name)

Example: 'user:smith 3200' will filter by 'smith' only in the user field and '3200' in all fields

Running Jobs

Jobid	Priority	User	Name	Map % Complete	Map Total	Maps Completed	Reduce % Complete	Reduce Total	Reduces Completed	Job Scheduling Information	Diagnostic Info
job_201310161728_0004	NORMAL	nikos	MyWordCount	34.33% <div></div>	2329	777	0.00% <div></div>	13	0	NA	NA

Failed Jobs

Jobid	Priority	User	Name	Map % Complete	Map Total	Maps Completed	Reduce % Complete	Reduce Total	Reduces Completed	Job Scheduling Information	Diagnostic Info
job_201310161728_0002	NORMAL	nikos	MyWordCount	100.00% <div></div>	2329	97	100.00% <div></div>	13	0	NA	NA
job_201310161728_0003	NORMAL	nikos	MyWordCount	100.00% <div></div>	2329	0	100.00% <div></div>	13	0	NA	NA

Retired Jobs

none

Local Logs

[Log directory: Job Tracker History](#)

Hadoop, 2013.



bigdata-06 Hadoop Machine List

Active Task Trackers

Task Trackers															
Name	Host	# running tasks	Max Map Tasks	Max Reduce Tasks	Task Failures	Directory Failures	Node Health Status	Seconds Since Node Last Healthy	Total Tasks Since Start	Succeeded Tasks Since Start	Total Tasks Last Day	Succeeded Tasks Last Day	Total Tasks Last Hour	Succeeded Tasks Last Hour	Seconds since heartbeat
tracker_bigdata-06.dcs.gla.ac.uk:localhost@127.0.0.1:47880	bigdata-06.dcs.gla.ac.uk	32	32	16	0	0	N/A	0	669	605	0	0	654	590	0
tracker_bigdata-05.dcs.gla.ac.uk:localhost@127.0.0.1:44742	bigdata-05.dcs.gla.ac.uk	4	4	2	0	0	N/A	0	112	108	0	0	111	107	0
tracker_bigdata-02.dcs.gla.ac.uk:localhost@127.0.0.1:15813	bigdata-02.dcs.gla.ac.uk	4	4	2	0	0	N/A	0	102	98	0	0	101	97	0
tracker_bigdata-04.dcs.gla.ac.uk:localhost@127.0.0.1:43520	bigdata-04.dcs.gla.ac.uk	4	4	2	0	0	N/A	0	109	108	0	0	107	106	0
tracker_bigdata-01.dcs.gla.ac.uk:localhost@127.0.0.1:15834	bigdata-01.dcs.gla.ac.uk	4	4	2	0	0	N/A	0	112	108	0	0	111	107	0
tracker_bigdata-03.dcs.gla.ac.uk:localhost@127.0.0.1:44060	bigdata-03.dcs.gla.ac.uk	4	4	2	0	0	N/A	0	113	110	0	0	109	106	0



JobTracker UI :: Job information



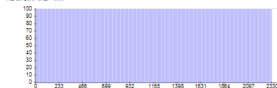
Hadoop job_201310161728_0004 on bigdata-06

User info
 Job Name: /j/1617280004
 Job File Path: /hdfs://bigdata-06.gla.ac.uk:8020/user/hive/hive/job_201310161728_0004.yc
 Submit Host Address: 10.233.233.222
 Job-URL: No users are enrolled
 Job Setup: Success
 Status: Success
 Started at: Wed Oct 16 17:48:28 2013
 Finished at: Wed Oct 16 18:02:18 2013
 Finished in: 27mins, 5secs
 Job Cleanup: Success

Kind	% Complete	Num Tasks	Running	Complete	Failed	Failed/Retryable Attempts
map	100.000%	2229	0	2229	0	0/0
reduce	100.000%	12	0	0	12	0/0

	Counter	Map	Reduce	Total
File System Counters	File: Number of bytes read	6,359,039	20,249	6,379,288
	File: Number of bytes written	276,348,824	0/12,781	276,348,824
	File: Number of read operations	0	0	0
	File: Number of large read operations	0	0	0
	File: Number of write operations	0	0	0
	HDFS: Number of bytes read	0/0,871,047,744	0	0/0,871,047,744
	HDFS: Number of bytes written	0	229	229
	HDFS: Number of read operations	4,760	4	4,764
	HDFS: Number of large read operations	0	0	0
	HDFS: Number of write operations	0	12	12
JVM Counters	Uncompressed heap memory	0	0	0.000
	Uncompressed heap memory	0	0	0.000
	Uncompressed heap memory	0	0	0.000
	Uncompressed heap memory	0	0	0.000
	Uncompressed heap memory	0	0	0.000
	Uncompressed heap memory	0	0	0.000
	Uncompressed heap memory	0	0	0.000
Map Input Summary	Map Input Bytes	17,229,448,280	0	17,229,448,280
	Map Input Bytes	280/113	0	280/113
	Map Input Bytes	1,015,854,880	21,851	1,015,854,880
	Map Input Bytes	45,000	28	45,028
	Map Input Bytes	0	21	21
	Map Input Bytes	0	21,249	21,249
	Map Input Bytes	0	1,000	1,000
	Map Input Bytes	0	12	12
	Map Input Bytes	0	12	12
	Map Input Bytes	0	12	12
Map Output Summary	Map Output Bytes	1,015,854,880	21,851	1,015,854,880
	Map Output Bytes	45,000	28	45,028
	Map Output Bytes	0	21	21
	Map Output Bytes	0	21,249	21,249
	Map Output Bytes	0	1,000	1,000
	Map Output Bytes	0	12	12
	Map Output Bytes	0	12	12
	Map Output Bytes	0	12	12
	Map Output Bytes	0	12	12
	Map Output Bytes	0	12	12
Map Output Summary	Map Output Bytes	1,015,854,880	21,851	1,015,854,880
	Map Output Bytes	45,000	28	45,028
	Map Output Bytes	0	21	21
	Map Output Bytes	0	21,249	21,249
	Map Output Bytes	0	1,000	1,000
	Map Output Bytes	0	12	12
	Map Output Bytes	0	12	12
	Map Output Bytes	0	12	12
	Map Output Bytes	0	12	12
	Map Output Bytes	0	12	12

Map Output Summary



Reduce Output Summary



Hadoop job_201310161728_0004 on bigdata-06

User: nikos

Job Name: MyWordCount

Job File: hdfs://bigdata-06.dcs.gla.ac.uk:8020/user/nikos/.staging/job_201310161728_0004/job.xml

Submit Host: bigdata-06.dcs.gla.ac.uk

Submit Host Address: 130.209.255.236

Job-ACLs: All users are allowed



Job Setup: Successful

Status: Running

Started at: Wed Oct 16 17:45:24 BST 2013

Running for: 8mins, 4sec

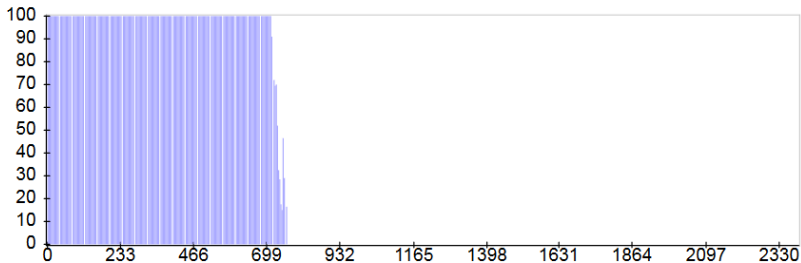
Job Cleanup: Pending

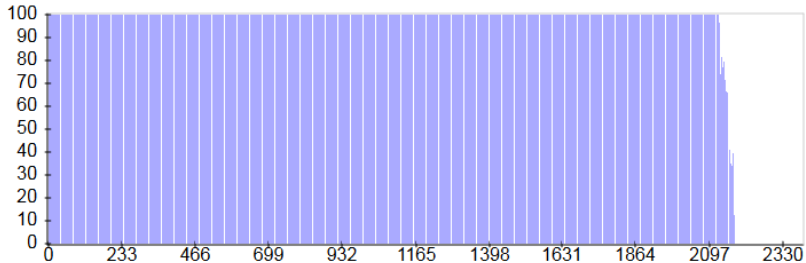
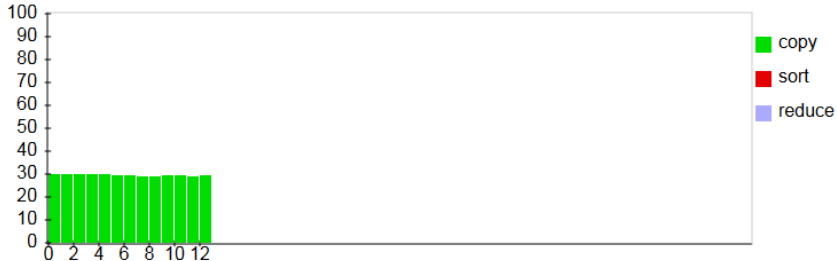
Kind	% Complete	Num Tasks	Pending	Running	Complete	Killed	Failed/Killed Task Attempts
map	25.87% 	2329	1695	52	582	0	0 / 0
reduce	0.00% 	13	13	0	0	0	0 / 0

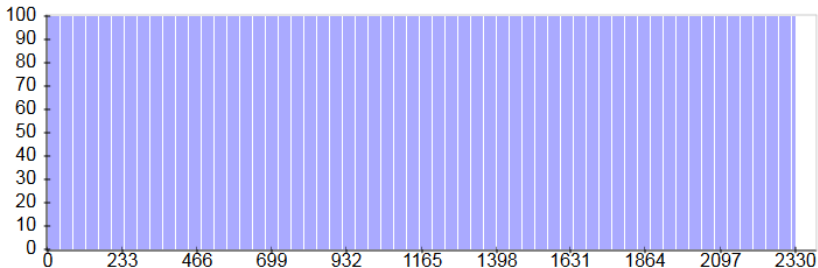


	Counter	Map	Reduce	Total
File System Counters	FILE: Number of bytes read	55,946	0	55,946
	FILE: Number of bytes written	113,678,798	0	113,678,798
	FILE: Number of read operations	0	0	0
	FILE: Number of large read operations	0	0	0
	FILE: Number of write operations	0	0	0
	HDFS: Number of bytes read	92,903,754,407	0	92,903,754,407
	HDFS: Number of bytes written	0	0	0
	HDFS: Number of read operations	1,460	0	1,460
	HDFS: Number of large read operations	0	0	0
	HDFS: Number of write operations	0	0	0
Job Counters	Launched map tasks	0	0	741
	Data-local map tasks	0	0	407
	Rack-local map tasks	0	0	334
	Total time spent by all maps in occupied slots (ms)	0	0	28,198,158
Map-Reduce Framework	Map input records	18,582,761	0	18,582,761
	Map output records	241,575,893	0	241,575,893
	Map output bytes	2,750,248,628	0	2,750,248,628
	Input split bytes	89,856	0	89,856
	Combine input records	240,381,830	0	240,381,830
	Combine output records	9,100	0	9,100
	Spilled Records	9,386	0	9,386
	CPU time spent (ms)	18,670,340	0	18,670,340
	Physical memory (bytes) snapshot	458,600,370,176	0	458,600,370,176
	Virtual memory (bytes) snapshot	1,236,150,968,320	0	1,236,150,968,320
	Total committed heap usage (bytes)	594,153,635,840	0	594,153,635,840
uk.ac.gla.dcs.bd4.WordCount\$MyMapper\$Counters	NUM_BYTES	92,900,417,853	0	92,900,417,853
	NUM_LINES	241,575,893	0	241,575,893
	NUM_RECORDS	18,582,761	0	18,582,761



Map Completion Graph - [close](#)Reduce Completion Graph - [close](#)

Map Completion Graph - [close](#)Reduce Completion Graph - [close](#)

Map Completion Graph - [close](#)Reduce Completion Graph - [close](#)

JobTracker UI :: Task information



Task	Complete	Status	Start Time	Finish Time	Errors	Counters
task_201310161728_0004_m_001070	96.87% 		16-Oct-2013 18:01:38			24
task_201310161728_0004_m_001074	94.63% 		16-Oct-2013 18:01:39			24
task_201310161728_0004_m_001076	85.89% 		16-Oct-2013 18:01:40			24
task_201310161728_0004_m_001081	87.50% 		16-Oct-2013 18:01:47			24
task_201310161728_0004_m_001083	84.84% 		16-Oct-2013 18:01:48			24
task_201310161728_0004_m_001085	84.52% 		16-Oct-2013 18:01:49			24
task_201310161728_0004_m_001087	74.89% 		16-Oct-2013 18:01:54			24
task_201310161728_0004_m_001089	67.75% 		16-Oct-2013 18:02:03			24
task_201310161728_0004_m_001091	66.82% 		16-Oct-2013 18:02:04			24
task_201310161728_0004_m_001094	79.34% 		16-Oct-2013 18:02:07			24
task_201310161728_0004_m_001096	60.58% 		16-Oct-2013 18:02:07			24
task_201310161728_0004_m_001098	70.95% 		16-Oct-2013 18:02:08			24
task_201310161728_0004_m_001100	62.50% 		16-Oct-2013 18:02:10			24
task_201310161728_0004_m_001102	64.82% 		16-Oct-2013 18:02:14			24
task_201310161728_0004_m_001104	45.31% 		16-Oct-2013 18:02:19			24
task_201310161728_0004_m_001107	56.77% 		16-Oct-2013 18:02:21			24
task_201310161728_0004_m_001108	96.17% 		16-Oct-2013 18:02:24			24
task_201310161728_0004_m_001109	47.66% 		16-Oct-2013 18:02:25			24
task_201310161728_0004_m_001111	54.49% 		16-Oct-2013 18:02:25			24
task_201310161728_0004_m_001113	50.00% 		16-Oct-2013 18:02:25			24
task_201310161728_0004_m_001114	89.32% 		16-Oct-2013 18:02:26			24
task_201310161728_0004_m_001115	37.17% 		16-Oct-2013 18:02:26			24
task_201310161728_0004_m_001116	68.67% 		16-Oct-2013 18:02:32			24



Job [job_201310161728_0004](#)

All Task Attempts

Task Attempts	Machine	Status	Progress	Start Time	Finish Time	Errors	Task Logs	Counters	Actions
attempt_201310161728_0004_m_001210_0	/default/bigdata-06.dcs.gla.ac.uk	RUNNING	32.47% <div><div></div></div>	16-Oct-2013 18:04:02			Last 4KB Last 8KB All	21	

Input Split Locations

/default/bigdata-06.dcs.gla.ac.uk



Counters for task_201310161728_0004_m_001165

File System Counters		
	FILE: Number of bytes read	0
	FILE: Number of bytes written	164,566
	FILE: Number of read operations	0
	FILE: Number of large read operations	0
	FILE: Number of write operations	0
	HDFS: Number of bytes read	134,221,626
	HDFS: Number of bytes written	0
	HDFS: Number of read operations	2
	HDFS: Number of large read operations	0
	HDFS: Number of write operations	0
Map-Reduce Framework		
	Map input records	31,337
	Map output records	407,381
	Map output bytes	4,637,876
	Input split bytes	128
	Combine input records	407,381
	Combine output records	13
	Spilled Records	13
	CPU time spent (ms)	33,510
	Physical memory (bytes) snapshot	660,152,320
	Virtual memory (bytes) snapshot	1,778,925,568
	Total committed heap usage (bytes)	977,600,512
uk.ac.gla.dcs.bd4.WordCount\$MyMapper\$Counters		
	NUM_BYTES	134,216,895
	NUM_LINES	407,381
	NUM_RECORDS	31,337



Hadoop programming :: Hadoop data types



Hadoop data types

```
public interface Writable {  
    void readFields(DataInput in);  
    void write(DataOutput out);  
}
```

```
1 public class MyWritable implements Writable {  
2     private int value;  
3     private long timestamp;  
4  
5     public void write(DataOutput out) throws IOException {  
6         out.writeInt(value);  
7         out.writeLong(timestamp);  
8     }  
9  
10    public void readFields(DataInput in) throws IOException {  
11        value = in.readInt();  
12        timestamp = in.readLong();  
13    }  
14  
15    public static MyWritable read(DataInput in) throws IOException {  
16        MyWritable w = new MyWritable();  
17        w.readFields(in);  
18        return w;  
19    }  
20 }
```



Hadoop data types

```
public interface WritableComparable<T> extends Writable, Comparable<T> {  
    // Writable -> void readFields(DataInput in), void write(DataOutput out);  
    // Comparable<T> -> int compareTo(T o);  
}
```

```
1 public class MyWritableComparable implements WritableComparable {  
2     private int value;  
3     private long timestamp;  
4  
5     public void write(DataOutput out) throws IOException {  
6         out.writeInt(value); out.writeLong(timestamp);  
7     }  
8  
9     public void readFields(DataInput in) throws IOException {  
10         value = in.readInt(); timestamp = in.readLong();  
11     }  
12  
13     public int compareTo(MyWritableComparable o) {  
14         return (this.value < o.value ? -1 : (this.value == o.value ? 0 : 1));  
15     }  
16  
17     public int hashCode() {  
18         final int prime = 31;  
19         return prime * (prime + value) + (int) (timestamp ^ (timestamp >>> 32));  
20     }  
21 }
```



Hadoop data types

```
public abstract class BinaryComparable implements Comparable<BinaryComparable> {  
    abstract byte[] getBytes();  
    abstract int getLength();  
  
    int compareTo(BinaryComparable other);  
    int compareTo(byte[] other, int off, int len);  
    boolean equals(Object other);  
    int hashCode();  
}
```



Hadoop data types

- ObjectWritable
- GenericWritable
- NullWritable

- BooleanWritable
- ByteWritable
- ShortWritable, IntWritable, LongWritable
- FloatWritable, DoubleWritable
- VIntWritable, VLongWritable
- Text

- BytesWritable
- ArrayPrimitiveWritable, ArrayWritable, TwoDArrayWritable

- MapWritable, SortedMapWritable
- EnumSetWritable



Hadoop programming :: Basic Mapper/Reducer methods



Basic Mapper/Reducer methods

```
public class Mapper<KEYIN,VALUEIN,KEYOUT,VALUEOUT> {
    static class Context { ... }
    protected void setup(Context context) { ... }
    protected void map(KEYIN key, VALUEIN value, Context context) { ... }
    protected void cleanup(Context context) { ... }
    void run(Context context) {
        setup(context);
        while (context.nextKeyValue())
            map(context.getCurrentKey(), context.getCurrentValue(), context);
        cleanup(context);
    }
}

public class Reducer<KEYIN,VALUEIN,KEYOUT,VALUEOUT> {
    static class Context { ... }
    protected void setup(Context context) { ... }
    protected void reduce(KEYIN key, Iterable<VALUEIN> values, Context context) {
        ... }
    protected void cleanup(Context context) { ... }
    void run(Context context) {
        setup(context);
        while (context.nextKey())
            reduce(context.getCurrentKey(), context.getValues(), context);
        cleanup(context);
    }
}
```



Hadoop programming :: Custom input/output formats



Custom InputFormat

```
public abstract class InputFormat<KEY,VALUE> {  
    abstract List<InputSplit> getSplits(JobContext context);  
    abstract RecordReader<KEY,VALUE> createRecordReader(InputSplit split,  
        TaskAttemptContext context);  
}  
  
public abstract class RecordReader<KEY,VALUE> implements Closeable {  
    abstract void initialize(InputSplit split, TaskAttemptContext context);  
    abstract void close();  
    abstract boolean nextKeyValue();  
    abstract KEY getCurrentKey();  
    abstract VALUE getCurrentValue();  
    abstract float getProgress();  
}  
  
public abstract class InputSplit {  
    abstract long getLength();  
    abstract String[] getLocations();  
}
```



Custom OutputFormat

```
public abstract class OutputFormat<KEY,VALUE> {
    abstract void checkOutputSpecs(JobContext context);
    abstract OutputCommitter getOutputCommitter(TaskAttemptContext context);
    abstract RecordWriter<KEY,VALUE> getRecordWriter(TaskAttemptContext context);
}

public abstract class RecordWriter<KEY,VALUE> {
    abstract void close(TaskAttemptContext context);
    abstract void write(KEY key, VALUE value);
}

public abstract class OutputCommitter {
    abstract void setupTask(TaskAttemptContext taskContext);
    abstract boolean needsTaskCommit(TaskAttemptContext taskContext);
    abstract void commitTask(TaskAttemptContext taskContext);
    abstract void abortTask(TaskAttemptContext taskContext);
    boolean isRecoverySupported();
    void recoverTask(TaskAttemptContext taskContext);

    abstract void setupJob(JobContext jobContext);
    void commitJob(JobContext jobContext);
    void abortJob(JobContext jobContext, JobStatus.State state);
}
```



Hadoop programming :: Word Count galore



Word Count v0 :: Built-in mappers/reducers



Word Count v0

```
1 public class WordCount extends Configured implements Tool {
2     public int run(String[] args) throws Exception {
3         Job job = new Job();
4         job.setJobName("WordCount-v0");
5         job.setJarByClass(WordCount.class);
6
7         job.setMapperClass(TokenCounterMapper.class);
8         job.setReducerClass(IntSumReducer.class);
9
10        job.setOutputKeyClass(Text.class);
11        job.setOutputValueClass(IntWritable.class);
12
13        FileInputFormat.addInputPath(job, new Path(args[0]));
14        FileOutputFormat.setOutputPath(job, new Path(args[1]));
15
16        job.submit();
17        return (job.waitForCompletion(true) ? 0 : 1);
18    }
19
20    public static void main(String[] args) throws Exception {
21        System.exit(ToolRunner.run(new Configuration(), new WordCount(), args));
22    }
23 }
```



Word Count v1 :: User-defined mappers/reducers



Word Count v1

```
1 public class WordCount extends Configured implements Tool {
2     static class Map extends org.apache.hadoop.mapreduce.Mapper<LongWritable,
3         Text, Text, IntWritable> {
4         private final static IntWritable one = new IntWritable(1);
5         private Text word = new Text();
6
7         public void map(LongWritable key, Text value, Context context) throws
8             IOException, InterruptedException {
9             String line = value.toString();
10            StringTokenizer tokenizer = new StringTokenizer(line);
11            while (tokenizer.hasMoreTokens()) {
12                word.set(tokenizer.nextToken());
13                context.write(word, one);
14            }
15        }
16
17        public static class Reduce extends Reducer<Text, IntWritable, Text,
18            IntWritable> {
19            public void reduce(Text key, Iterable<IntWritable> values, Context
20                context) throws IOException, InterruptedException {
21                int sum = 0;
22                for (IntWritable value: values)
23                    sum += value.get();
24                context.write(key, new IntWritable(sum));
25            }
26        }
27    }
```



Word Count v1 (cont.)

```
24 public int run(String[] args) throws Exception {
25     Job job = new Job();
26     job.setJobName("WordCount-v1");
27     job.setJarByClass(WordCount.class);
28
29     job.setMapperClass(Map.class);
30     job.setCombinerClass(Reduce.class);
31     job.setReducerClass(Reduce.class);
32
33     job.setInputFormatClass(TextInputFormat.class);
34     FileInputFormat.setInputPaths(job, new Path(args[0]));
35
36     job.setOutputKeyClass(Text.class);
37     job.setOutputValueClass(IntWritable.class);
38     job.setOutputFormatClass(TextOutputFormat.class);
39     FileOutputFormat.setOutputPath(job, new Path(args[1]));
40
41     job.submit();
42     return (job.waitForCompletion(true) ? 0 : 1);
43 }
44
45 public static void main(String[] args) throws Exception {
46     System.exit(ToolRunner.run(new Configuration(), new WordCount(), args));
47 }
48 }
```



Word Count v2 ::

Distributed cache + configuration + counters +
status messages + progress report



Word Count v2

- Count occurrences of words in the input stream, but also:
 - Allow user to define patterns/words to be skipped
 - Allow user to turn case-sensitivity on/off
 - Count total number of words processed
 - Report progress and update status messages as we go



Word Count v2

WordCount.java

```
1 public class WordCount extends Configured implements Tool {
2     public int run(String[] args) throws Exception {
3         Job job = new Job();
4         job.setJobName("WordCount-v3");
5         job.setJarByClass(WordCount.class);
6         job.setMapperClass(Map.class);
7         job.setCombinerClass(Reduce.class);
8         job.setReducerClass(Reduce.class);
9         job.setInputFormatClass(TextInputFormat.class);
10        FileInputFormat.setInputPaths(job, new Path(other_args.get(0)));
11        job.setOutputKeyClass(Text.class);
12        job.setOutputValueClass(IntWritable.class);
13        job.setOutputFormatClass(TextOutputFormat.class);
14        FileOutputFormat.setOutputPath(job, new Path(other_args.get(1)));
15
16        List<String> other_args = new ArrayList<String>();
17        for (int i = 0; i < args.length; ++i) {
18            if ("-skip".equals(args[i])) {
19                DistributedCache.addCacheFile(new Path(args[++i]).toUri(),
20                    job.getConfiguration());
21                job.getConfiguration().setBoolean("wordcount.skip.patterns", true);
22            } else
23                other_args.add(args[i]);
24        }
25    }
26 }
```

Word Count v2 (cont.)

WordCount.java (cont.)

```
24     job.submit();  
25     return job.waitForCompletion(true) ? 0 : 1;  
26 }  
27  
28 public static void main(String[] args) throws Exception {  
29     System.exit(ToolRunner.run(new Configuration(), new WordCount(), args));  
30 }  
31 }
```



Word Count v2 (cont.)

Map.java

```
1 public class Map extends Mapper<LongWritable, Text, Text, IntWritable>{
2     static enum Counters { INPUT_WORDS }
3     private final static IntWritable one = new IntWritable(1);
4     private Text word = new Text();
5     private boolean caseSensitive = true;
6     private Set<String> patternsToSkip = new HashSet<String>();
7     private long numRecords = 0;
8     private String inputFile;
9
10    private void parseSkipFile(Path patternsFile) {
11        try {
12            BufferedReader fis = new BufferedReader(new
13                FileReader(patternsFile.toString()));
14            String pattern = null;
15            while ((pattern = fis.readLine()) != null)
16                patternsToSkip.add(pattern);
17            fis.close();
18        } catch (IOException ioe) {
19            System.err.println("Caught exception while parsing the cached file '" +
20                patternsFile + "' : " + StringUtils.stringifyException(ioe));
21        }
22    }
```

Word Count v2 (cont.)

Map.java (cont.)

```
21 public void setup(Context context) {
22     Configuration conf = context.getConfiguration();
23     caseSensitive = conf.getBoolean("wordcount.case.sensitive", true);
24     inputFile = conf.get("map.input.file");
25
26     if (conf.getBoolean("wordcount.skip.patterns", false)) {
27         Path[] patternsFiles = new Path[0];
28         try {
29             patternsFiles = DistributedCache.getLocalCacheFiles(conf);
30         } catch (IOException ioe) {
31             System.err.println("Caught exception while getting cached files: " +
32                 StringUtils.stringifyException(ioe));
33         }
34         for (Path patternsFile : patternsFiles)
35             parseSkipFile(patternsFile);
36     }
37
38     public void cleanup(Context context) {
39         patternsToSkip.clear();
40     }
```

Word Count v2 (cont.)

Map.java (cont.)

```
1 public void map(LongWritable key, Text value, Context context) throws
2     IOException, InterruptedException {
3     String line = caseSensitive ? value.toString() :
4         value.toString().toLowerCase();
5
6     for (String pattern : patternsToSkip)
7         line = line.replaceAll(pattern, "");
8
9     StringTokenizer tokenizer = new StringTokenizer(line);
10    while (tokenizer.hasMoreTokens()) {
11        word.set(tokenizer.nextToken());
12        context.write(word, one);
13        context.getCounter(Counters.INPUT_WORDS).increment(1);
14    }
15
16    if ((++numRecords % 100) == 0)
17        context.setStatus("Finished processing " + numRecords + " records " +
18            "from the input file: " + inputFile);
19    }
20 }
```

Word Count v2 (cont.)

Reduce.java

```
1 public static class Reduce extends Reducer<Text, IntWritable, Text, IntWritable>
2 {
3     public void reduce(Text key, Iterable<IntWritable> values, Context context)
4         throws IOException, InterruptedException {
5         int sum = 0;
6         for (IntWritable value: values)
7             sum += value.get();
8         context.write(key, new IntWritable(sum));
9     }
10 }
```



Word Count v3 :: Custom InputFormat + partitioner + more counters



Word Count v3

- Count occurrences of first word in each “record” in the input stream, but this time:
 - Input records are spread across multiple lines ...
 - ... with a special sequence (`\t\t\t`) separating them ...
 - Also make a custom partitioner so there are 27 reducers (for `[a-z]` + `<everything else>`)...
 - Also count total number of bytes, lines, and records processed



Word Count v3

MyPartitioner.java

```
1 public class MyPartitioner extends Partitioner<Text, IntWritable> {  
2     int getPartition(Text key, IntWritable value, int numPartitions) {  
3         int c = Character.toLowerCase(key.toString().charAt(0));  
4         if (c < 'a' || c > 'z')  
5             return numPartitions - 1;  
6         return (int) Math.floor((float) (numPartitions - 2) * (c - 'a') / ('z' - 'a'));  
7     }  
8 }
```

MyInputFormat.java

```
1 public class MyInputFormat extends FileInputFormat<LongWritable, Text> {  
2     public RecordReader<LongWritable, Text> createRecordReader(InputSplit split,  
3         TaskAttemptContext context) {  
4         return new MyRecordReader();  
5     }  
6 }
```

Word Count v3 (cont.)

MyRecordReader.java

```
1 public class MyRecordReader extends RecordReader<LongWritable, Text> {
2     private static final byte[] recordSeparator = "\t\t\t".getBytes();
3     private FSDataInputStream fsin;
4     private long start, end;
5     private boolean stillInChunk = true;
6     private DataOutputBuffer buffer = new DataOutputBuffer();
7     private LongWritable key = new LongWritable();
8     private Text value = new Text();
9
10    public void initialize(InputSplit inputSplit, TaskAttemptContext context)
11        throws IOException {
12        FileSplit split = (FileSplit) inputSplit;
13        Configuration conf = context.getConfiguration();
14        Path path = split.getPath();
15        FileSystem fs = path.getFileSystem(conf);
16
17        fsin = fs.open(path);
18        start = split.getStart();
19        end = split.getStart() + split.getLength();
20        fsin.seek(start);
21
22        if (start != 0)
23            readRecord(false);
24    }
```

Word Count v3 (cont.)

MyRecordReader.java (cont.)

```
24 private boolean readRecord(boolean withinBlock) throws IOException {  
25     int i = 0, b;  
26     while (true) {  
27         if ((b = fsin.read()) == -1)  
28             return false;  
29         if (withinBlock)  
30             buffer.write(b);  
31         if (b == recordSeparator[i]) {  
32             if (++i == recordSeparator.length)  
33                 return fsin.getPos() < end;  
34         } else  
35             i = 0;  
36     }  
37 }
```



Word Count v3 (cont.)

MyRecordReader.java (cont.)

```
38 public boolean nextKeyValue() throws IOException {
39     if (!stillInChunk)
40         return false;
41     boolean status = readRecord(true);
42     value = new Text();
43     value.set(buffer.getData(), 0, buffer.getLength());
44     key.set(fsин.getPos());
45     buffer.reset();
46     if (!status)
47         stillInChunk = false;
48     return true;
49 }
50
51 public LongWritable getCurrentKey() { return key; }
52
53 public Text getCurrentValue() { return value; }
54
55 public float getProgress() throws IOException {
56     return (float) (fsин.getPos() - start) / (end - start);
57 }
58
59 public void close() throws IOException { fsин.close(); }
60 }
```

Word Count v3 (cont.)

MyMapper.java

```
1 public class MyMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
2     static enum Counters { NUM_RECORDS, NUM_LINES, NUM_BYTES }
3     private Text _key = new Text();
4     private IntWritable _value = new IntWritable();
5
6     protected void map(LongWritable key, Text value, Context context) throws
7         IOException, InterruptedException {
8         StringTokenizer tokenizer = new StringTokenizer(value.toString(), "\n");
9         while (tokenizer.hasMoreTokens()) {
10             String line = tokenizer.nextToken();
11             int sep = line.indexOf(' ');
12             _key.set((sep == -1) ? line : line.substring(0, line.indexOf(' ')));
13             _value.set(1);
14             context.write(_key, _value);
15             context.getCounter(Counters.NUM_LINES).increment(1);
16         }
17         context.getCounter(Counters.NUM_BYTES).increment(value.getLength());
18         context.getCounter(Counters.NUM_RECORDS).increment(1);
19     }
```

Word Count v3 (cont.)

MyReducer.java

```
1 public class MyReducer extends Reducer<Text, IntWritable, Text, IntWritable> {  
2     private IntWritable _value = new IntWritable();  
3     protected void reduce(Text key, Iterable<IntWritable> values, Context  
4         context) throws IOException, InterruptedException {  
5         int sum = 0;  
6         for (Iterator<IntWritable> it = values.iterator(); it.hasNext();) {  
7             sum += it.next().get();  
8             _value.set(sum);  
9             context.write(key, _value);  
10        }  
11    }
```



Word Count v3 (cont.)

WordCount.java

```
1 public class WordCount extends Configured implements Tool {
2     public int run(String[] args) throws Exception {
3         Job job = new Job();
4         job.setJobName("MyWordCount(" + args[0] + ")");
5         job.setJarByClass(WordCount.class);
6         job.setInputFormatClass(MyInputFormat.class);
7         job.setOutputFormatClass(TextOutputFormat.class);
8         job.setMapperClass(MyMapper.class);
9         job.setPartitionerClass(MyPartitioner.class);
10        job.setMapOutputKeyClass(Text.class);
11        job.setMapOutputValueClass(IntWritable.class);
12        job.setReducerClass(MyReducer.class);
13        job.setCombinerClass(MyReducer.class);
14        FileInputFormat.setInputPaths(job, new Path(args[0]));
15        FileOutputFormat.setOutputPath(job, new Path(job.getJobName() +
16            "_output"));
17        job.submit();
18        return job.waitForCompletion(true) ? 0 : 1;
19    }
20    public static void main(String[] args) throws Exception {
21        System.exit(ToolRunner.run(new WordCount(), args));
22    }
23 }
```