

MapReduce Tutorial

- [Purpose](#)
- [Prerequisites](#)
- [Overview](#)
- [Inputs and Outputs](#)
- [Example: WordCount v1.0](#)
 - [Source Code](#)
 - [Usage](#)
 - [Walk-through](#)
- [MapReduce - User Interfaces](#)
 - [Payload](#)
 - [Mapper](#)
 - [Reducer](#)
 - [Partitioner](#)
 - [Counter](#)
 - [Job Configuration](#)
 - [Task Execution & Environment](#)
 - [Memory Management](#)
 - [Map Parameters](#)
 - [Shuffle/Reduce Parameters](#)
 - [Configured Parameters](#)
 - [Task Logs](#)
 - [Distributing Libraries](#)
 - [Job Submission and Monitoring](#)
 - [Job Control](#)
 - [Job Input](#)
 - [InputSplit](#)
 - [RecordReader](#)
 - [Job Output](#)
 - [OutputCommitter](#)
 - [Task Side-Effect Files](#)
 - [RecordWriter](#)
 - [Other Useful Features](#)
 - [Submitting Jobs to Queues](#)
 - [Counters](#)
 - [DistributedCache](#)
 - [Profiling](#)
 - [Debugging](#)
 - [Data Compression](#)
 - [Skipping Bad Records](#)
 - [Example: WordCount v2.0](#)
 - [Source Code](#)
 - [Sample Runs](#)
 - [Highlights](#)

Purpose

This document comprehensively describes all user-facing facets of the Hadoop MapReduce framework and serves as a tutorial.

Prerequisites

Ensure that Hadoop is installed, configured and is running. More details:

- [Single Node Setup](#) for first-time users.
- [Cluster Setup](#) for large, distributed clusters.

Overview

Hadoop MapReduce is a software framework for easily writing applications which process vast amounts of data (multi-terabyte data-sets) in-parallel on large clusters (thousands of nodes) of commodity hardware in a reliable, fault-tolerant manner.

A MapReduce *job* usually splits the input data-set into independent chunks which are processed by the *map tasks* in a completely parallel manner. The framework sorts the outputs of the maps, which are then input to the *reduce tasks*. Typically both the input and the output of the job are stored in a file-system. The framework takes care of scheduling tasks, monitoring them and re-executes the failed tasks.


Typically the compute nodes and the storage nodes are the same, that is, the MapReduce framework and the Hadoop Distributed File System (see [HDFS Architecture Guide](#)) are running on the same set of nodes. This configuration allows the framework to effectively schedule tasks on the nodes where data is already present, resulting in very high aggregate bandwidth across the cluster.

The MapReduce framework consists of a single master ResourceManager, one worker NodeManager per cluster-node, and MRAppMaster per application (see [YARN Architecture Guide](#)).

Minimally, applications specify the input/output locations and supply *map* and *reduce* functions via implementations of appropriate interfaces and/or abstract-classes. These, and other job parameters, comprise the *job configuration*.

The Hadoop *job client* then submits the job (jar/executable etc.) and configuration to the ResourceManager which then assumes the responsibility of distributing the software/configuration to the workers, scheduling tasks and monitoring them, providing status and diagnostic information to the job-client.

Although the Hadoop framework is implemented in Java™, MapReduce applications need not be written in Java.

- [Hadoop Streaming](#) is a utility which allows users to create and run jobs with any executables (e.g. shell utilities) as the mapper and/or the reducer.
- [Hadoop Pipes](#) is a [SWIG](#) -compatible C++ API to implement MapReduce applications (non JNI™ based).

Inputs and Outputs

The MapReduce framework operates exclusively on <key, value> pairs, that is, the framework views the input to the job as a set of <key, value> pairs and produces a set of <key, value> pairs as the output of the job, conceivably of different types.

The key and value classes have to be serializable by the framework and hence need to implement the [Writable](#) interface. Additionally, the key classes have to implement the [WritableComparable](#) interface to facilitate sorting by the framework.

Input and Output types of a MapReduce job:

(input) <k1, v1> -> **map** -> <k2, v2> -> **combine** -> <k2, v2> -> **reduce** -> <k3, v3> (output)

Example: WordCount v1.0

Before we jump into the details, let's walk through an example MapReduce application to get a flavour for how they work.

WordCount is a simple application that counts the number of occurrences of each word in a given input set.

This works with a local-standalone, pseudo-distributed or fully-distributed Hadoop installation ([Single Node Setup](#)).

Source Code

```
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);
        }
    }
}

public static class IntSumReducer
    extends Reducer<Text, IntWritable, Text, IntWritable> {
    private IntWritable result = new IntWritable();

    public void reduce(Text key, Iterable<IntWritable> values,
        Context context
        ) throws IOException, InterruptedException {

        int sum = 0;
        for (IntWritable val : values) {
            sum += val.get();
        }
        result.set(sum);
        context.write(key, result);
    }
}

public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "word count");
    job.setJarByClass(WordCount.class);
    job.setMapperClass(TokenizerMapper.class);
    job.setCombinerClass(IntSumReducer.class);
    job.setReducerClass(IntSumReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    System.exit(job.waitForCompletion(true) ? 0 : 1);
}
}

```

Usage

Assuming environment variables are set as follows:

```

export JAVA_HOME=/usr/java/default
export PATH=${JAVA_HOME}/bin:${PATH}
export HADOOP_CLASSPATH=${JAVA_HOME}/lib/tools.jar

```

Compile WordCount.java and create a jar:

```
$ bin/hadoop com.sun.tools.javac.Main WordCount.java
$ jar cf wc.jar WordCount*.class
```

Assuming that:

- /user/joe/wordcount/input - input directory in HDFS
- /user/joe/wordcount/output - output directory in HDFS

Sample text-files as input:

```
$ bin/hadoop fs -ls /user/joe/wordcount/input/
/user/joe/wordcount/input/file01
/user/joe/wordcount/input/file02

$ bin/hadoop fs -cat /user/joe/wordcount/input/file01
Hello World Bye World

$ bin/hadoop fs -cat /user/joe/wordcount/input/file02
Hello Hadoop Goodbye Hadoop
```

Run the application:

```
$ bin/hadoop jar wc.jar WordCount /user/joe/wordcount/input /user/joe/wordcount/output
```

Output:

```
$ bin/hadoop fs -cat /user/joe/wordcount/output/part-r-00000
Bye 1
Goodbye 1
Hadoop 2
Hello 2
World 2
```

Applications can specify a comma separated list of paths which would be present in the current working directory of the task using the option `-files`. The `-libjars` option allows applications to add jars to the classpaths of the maps and reduces. The option `-archives` allows them to pass comma separated list of archives as arguments. These archives are unarchived and a link with name of the archive is created in the current working directory of tasks. More details about the command line options are available at [Commands Guide](https://hadoop.apache.org/docs/stable/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html).

Running wordcount example with `-libjars`, `-files` and `-archives`:

```
bin/hadoop jar hadoop-mapreduce-examples-<ver>.jar wordcount -files cachefile.txt -
```

Here, myarchive.zip will be placed and unzipped into a directory by the name "myarchive.zip".

Users can specify a different symbolic name for files and archives passed through `-files` and `-archives` option, using `#`.

For example,

```
bin/hadoop jar hadoop-mapreduce-examples-<ver>.jar wordcount -files dir1/dict.txt#d
```

Here, the files `dir1/dict.txt` and `dir2/dict.txt` can be accessed by tasks using the symbolic names `dict1` and `dict2` respectively. The archive `mytar.tgz` will be placed and unarchived into a directory by the name "tgzdir".

Applications can specify environment variables for mapper, reducer, and application master tasks by specifying them on the command line using the options `-Dmapreduce.map.env`, `-Dmapreduce.reduce.env`, and `-Dyarn.app.mapreduce.am.env`, respectively.

For example the following sets environment variables `FOO_VAR=bar` and `LIST_VAR=a,b,c` for the mappers and reducers,

```
bin/hadoop jar hadoop-mapreduce-examples-<ver>.jar wordcount -Dmapreduce.map.env.FOO
```

➔ Walk-through

The WordCount application is quite straight-forward.

```
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);
    }
}
```

The Mapper implementation, via the `map` method, processes one line at a time, as provided by the specified `TextInputFormat`. It then splits the line into tokens separated by whitespaces, via the `StringTokenizer`, and emits a key-value pair of `< <word>, 1>`.

For the given sample input the first map emits:

```
< Hello, 1>
< World, 1>
< Bye, 1>
< World, 1>
```

The second map emits:

```
< Hello, 1>
< Hadoop, 1>
< Goodbye, 1>
< Hadoop, 1>
```

We'll learn more about the number of maps spawned for a given job, and how to control them in a fine-grained manner, a bit later in the tutorial.

```
job.setCombinerClass(IntSumReducer.class);
```

WordCount also specifies a combiner. Hence, the output of each map is passed through the local combiner (which is same as the Reducer as per the job configuration) for local aggregation, after being sorted on the *keys*.

The output of the first map:

```
< Bye, 1>
< Hello, 1>
< World, 2>
```

The output of the second map:

```
< Goodbye, 1>
< Hadoop, 2>
< Hello, 1>
```

```
public void reduce(Text key, Iterable<IntWritable> values,
                  Context context
                  ) throws IOException, InterruptedException {
    int sum = 0;
    for (IntWritable val : values) {
        sum += val.get();
    }
    result.set(sum);
}
```

```
context.write(key, result);  
}
```

The Reducer implementation, via the reduce method just sums up the values, which are the occurrence counts for each key (i.e. words in this example).

Thus the output of the job is:

```
< Bye, 1>  
< Goodbye, 1>  
< Hadoop, 2>  
< Hello, 2>  
< World, 2>
```

The main method specifies various facets of the job, such as the input/output paths (passed via the command line), key/value types, input/output formats etc., in the Job. It then calls the `job.waitForCompletion` to submit the job and monitor its progress.

We'll learn more about Job, InputFormat, OutputFormat and other interfaces and classes a bit later in the tutorial.

MapReduce - User Interfaces

This section provides a reasonable amount of detail on every user-facing aspect of the MapReduce framework. This should help users implement, configure and tune their jobs in a fine-grained manner. However, please note that the javadoc for each class/interface remains the most comprehensive documentation available; this is only meant to be a tutorial.

Let us first take the Mapper and Reducer interfaces. Applications typically implement them to provide the map and reduce methods.

We will then discuss other core interfaces including Job, Partitioner, InputFormat, OutputFormat, and others.

Finally, we will wrap up by discussing some useful features of the framework such as the DistributedCache, IsolationRunner etc.

Payload

Applications typically implement the Mapper and Reducer interfaces to provide the map and reduce methods. These form the core of the job.

Mapper

Mapper maps input key/value pairs to a set of intermediate key/value pairs.

Maps are the individual tasks that transform input records into intermediate records. The transformed intermediate records do not need to be of the same type as the input records. A given input pair may map to zero or many output pairs.

The Hadoop MapReduce framework spawns one map task for each InputSplit generated by the InputFormat for the job.

Overall, mapper implementations are passed to the job via `Job.setMapperClass(Class)` method. The framework then calls `map(WritableComparable, Writable, Context)` for each key/value pair in the InputSplit for that task. Applications can then override the `cleanup(Context)` method to perform any required cleanup.

Output pairs do not need to be of the same types as input pairs. A given input pair may map to zero or many output pairs. Output pairs are collected with calls to `context.write(WritableComparable, Writable)`.

Applications can use the Counter to report its statistics.

All intermediate values associated with a given output key are subsequently grouped by the framework, and passed to the Reducer(s) to determine the final output. Users can control the grouping by specifying a Comparator via [Job.setGroupingComparatorClass\(Class\)](#).

The Mapper outputs are sorted and then partitioned per Reducer. The total number of partitions is the same as the number of reduce tasks for the job. Users can control which keys (and hence records) go to which Reducer by implementing a custom Partitioner.

Users can optionally specify a combiner, via [Job.setCombinerClass\(Class\)](#), to perform local aggregation of the intermediate outputs, which helps to cut down the amount of data transferred from the Mapper to the Reducer.

The intermediate, sorted outputs are always stored in a simple (key-len, key, value-len, value) format. Applications can control if, and how, the intermediate outputs are to be compressed and the [CompressionCodec](#) to be used via the Configuration.

How Many Maps?

The number of maps is usually driven by the total size of the inputs, that is, the total number of blocks of the input files.

The right level of parallelism for maps seems to be around 10-100 maps per-node, although it has been set up to 300 maps for very cpu-light map tasks. Task setup takes a while, so it is best if the maps take at least a minute to execute.

Thus, if you expect 10TB of input data and have a blocksize of 128MB, you'll end up with 82,000 maps, unless `Configuration.set(MRJobConfig.NUM_MAPS, int)` (which only provides a hint to the framework) is used to set it even higher.

Reducer

[Reducer](#) reduces a set of intermediate values which share a key to a smaller set of values.

The number of reduces for the job is set by the user via [Job.setNumReduceTasks\(int\)](#).

Overall, Reducer implementations are passed the Job for the job via the [Job.setReducerClass\(Class\)](#) method and can override it to initialize themselves. The framework then calls [reduce\(WritableComparable, Iterable<Writable>, Context\)](#) method for each <key, (list of values)> pair in the grouped inputs. Applications can then override the `cleanup(Context)` method to perform any required cleanup.

Reducer has 3 primary phases: shuffle, sort and reduce.

Shuffle

Input to the Reducer is the sorted output of the mappers. In this phase the framework fetches the relevant partition of the output of all the mappers, via HTTP.

Sort

The framework groups Reducer inputs by keys (since different mappers may have output the same key) in this stage.

The shuffle and sort phases occur simultaneously; while map-outputs are being fetched they are merged.

Secondary Sort

If equivalence rules for grouping the intermediate keys are required to be different from those for grouping keys before reduction, then one may specify a Comparator via [Job.setSortComparatorClass\(Class\)](#). Since [Job.setGroupingComparatorClass\(Class\)](#) can be used to control how intermediate keys are grouped, these can be used in conjunction to simulate *secondary sort on values*.

Reduce

In this phase the `reduce(WritableComparable, Iterable<Writable>, Context)` method is called for each <key, (list of values)> pair in the grouped inputs.

The output of the reduce task is typically written to the [FileSystem](#) via `Context.write(WritableComparable, Writable)`.

Applications can use the Counter to report its statistics.

The output of the Reducer is *not sorted*.

How Many Reduces?

The right number of reduces seems to be 0.95 or 1.75 multiplied by (*<no. of nodes> * <no. of maximum containers per node>*).

With 0.95 all of the reduces can launch immediately and start transferring map outputs as the maps finish. With 1.75 the faster nodes will finish their first round of reduces and launch a second wave of reduces doing a much better job of load balancing.

Increasing the number of reduces increases the framework overhead, but increases load balancing and lowers the cost of failures.

The scaling factors above are slightly less than whole numbers to reserve a few reduce slots in the framework for speculative-tasks and failed tasks.

Reducer NONE

It is legal to set the number of reduce-tasks to *zero* if no reduction is desired.

In this case the outputs of the map-tasks go directly to the `FileSystem`, into the output path set by `FileOutputFormat.setOutputPath(Job, Path)`. The framework does not sort the map-outputs before writing them out to the `FileSystem`.

Partitioner

`Partitioner` partitions the key space.

`Partitioner` controls the partitioning of the keys of the intermediate map-outputs. The key (or a subset of the key) is used to derive the partition, typically by a *hash function*. The total number of partitions is the same as the number of reduce tasks for the job. Hence this controls which of the m reduce tasks the intermediate key (and hence the record) is sent to for reduction.

`HashPartitioner` is the default `Partitioner`.

Counter

`Counter` is a facility for MapReduce applications to report its statistics.

Mapper and Reducer implementations can use the Counter to report statistics.

Hadoop MapReduce comes bundled with a `library` of generally useful mappers, reducers, and partitioners.

Job Configuration

`Job` represents a MapReduce job configuration.

`Job` is the primary interface for a user to describe a MapReduce job to the Hadoop framework for execution. The framework tries to faithfully execute the job as described by `Job`, however:

- Some configuration parameters may have been marked as final by administrators (see [Final Parameters](#)) and hence cannot be altered.
- While some job parameters are straight-forward to set (e.g. `Job.setNumReduceTasks(int)`), other parameters interact subtly with the rest of the framework and/or job configuration and are more complex to set (e.g. `Configuration.set(JobContext.NUM_MAPS, int)`).

`Job` is typically used to specify the Mapper, combiner (if any), `Partitioner`, Reducer, `InputFormat`, `OutputFormat` implementations. `FileInputFormat` indicates the set of input files (`FileInputFormat.setInputPaths(Job, Path...)/FileInputFormat.addInputPath(Job, Path)`) and (`FileInputFormat.setInputPaths(Job, String...)/FileInputFormat.addInputPaths(Job, String)`) and where the output files should be written (`FileOutputFormat.setOutputPath(Path)`).

Optionally, `Job` is used to specify other advanced facets of the job such as the `Comparator` to be used, files to be put in the `DistributedCache`, whether intermediate and/or job outputs are to be compressed (and how), whether job tasks can be executed in a *speculative* manner (`setMapSpeculativeExecution(boolean)`)/

`setReduceSpeculativeExecution(boolean))`, maximum number of attempts per task (`setMaxMapAttempts(int)/setMaxReduceAttempts(int)`) etc.

Of course, users can use `Configuration.set(String, String)/ Configuration.get(String)` to set/get arbitrary parameters needed by applications. However, use the DistributedCache for large amounts of (read-only) data.

Task Execution & Environment

The MRAppMaster executes the Mapper/Reducer *task* as a child process in a separate jvm.

The child-task inherits the environment of the parent MRAppMaster. The user can specify additional options to the child-jvm via the `mapreduce.{map|reduce}.java.opts` and configuration parameter in the Job such as non-standard paths for the run-time linker to search shared libraries via `-Djava.library.path=<>` etc. If the `mapreduce.{map|reduce}.java.opts` parameters contains the symbol `@taskid@` it is interpolated with value of taskid of the MapReduce task.

Here is an example with multiple arguments and substitutions, showing jvm GC logging, and start of a passwordless JVM JMX agent so that it can connect with jconsole and the likes to watch child memory, threads and get thread dumps. It also sets the maximum heap-size of the map and reduce child jvm to 512MB & 1024MB respectively. It also adds an additional path to the `java.library.path` of the child-jvm.

```
<property>
  <name>mapreduce.map.java.opts</name>
  <value>
    -Xmx512M -Djava.library.path=/home/mycompany/lib -verbose:gc -Xloggc:/tmp/@taskid
    -Dcom.sun.management.jmxremote.authenticate=false -Dcom.sun.management.jmxremote.
  </value>
</property>

<property>
  <name>mapreduce.reduce.java.opts</name>
  <value>
    -Xmx1024M -Djava.library.path=/home/mycompany/lib -verbose:gc -Xloggc:/tmp/@taskid
    -Dcom.sun.management.jmxremote.authenticate=false -Dcom.sun.management.jmxremote.
  </value>
</property>
```

Memory Management

Users/admins can also specify the maximum virtual memory of the launched child-task, and any sub-process it launches recursively, using `mapreduce.{map|reduce}.memory.mb`. Note that the value set here is a per process limit. The value for `mapreduce.{map|reduce}.memory.mb` should be specified in mega bytes (MB). And also the value must be greater than or equal to the `-Xmx` passed to JavaVM, else the VM might not start.

Note: `mapreduce.{map|reduce}.java.opts` are used only for configuring the launched child tasks from MRAppMaster. Configuring the memory options for daemons is documented in [Configuring the Environment of the Hadoop Daemons](#).

The memory available to some parts of the framework is also configurable. In map and reduce tasks, performance may be influenced by adjusting parameters influencing the concurrency of operations and the frequency with which data will hit disk. Monitoring the filesystem counters for a job- particularly relative to byte counts from the map and into the reduce- is invaluable to the tuning of these parameters.

Map Parameters

A record emitted from a map will be serialized into a buffer and metadata will be stored into accounting buffers. As described in the following options, when either the serialization buffer or the metadata exceed a threshold, the contents of the buffers will be sorted and written to disk in the background while the map continues to output records. If either buffer fills completely while the spill is in progress, the map thread will block. When the map is finished, any remaining records are written to disk and all on-disk segments are merged into a single file. Minimizing the number of spills to disk can decrease map time, but a larger buffer also decreases the memory available to the mapper.

Name	Type	Description
mapreduce.task.io.sort.mb	int	The cumulative size of the serialization and accounting buffers storing records emitted from the map, in megabytes.
mapreduce.map.sort.spill.percent	float	The soft limit in the serialization buffer. Once reached, a thread will begin to spill the contents to disk in the background.

Other notes

- If either spill threshold is exceeded while a spill is in progress, collection will continue until the spill is finished. For example, if `mapreduce.map.sort.spill.percent` is set to 0.33, and the remainder of the buffer is filled while the spill runs, the next spill will include all the collected records, or 0.66 of the buffer, and will not generate additional spills. In other words, the thresholds are defining triggers, not blocking.
- A record larger than the serialization buffer will first trigger a spill, then be spilled to a separate file. It is undefined whether or not this record will first pass through the combiner.

Shuffle/Reduce Parameters

As described previously, each reduce fetches the output assigned to it by the Partitioner via HTTP into memory and periodically merges these outputs to disk. If intermediate compression of map outputs is turned on, each output is decompressed into memory. The following options affect the frequency of these merges to disk prior to the reduce and the memory allocated to map output during the reduce.

Name	Type	Description
mapreduce.task.io.soft.factor	int	Specifies the number of segments on disk to be merged at the same time. It limits the number of open files and compression codecs during merge. If the number of files exceeds this limit, the merge will proceed in several passes. Though this limit also applies to the map, most jobs should be configured so that hitting this limit is unlikely there.
mapreduce.reduce.merge.inmem.thresholds	int	The number of sorted map outputs fetched into memory before being merged to disk. Like the spill thresholds in the preceding note, this is not defining a unit of partition, but a trigger. In practice, this is usually set very high (1000) or disabled (0), since merging in-memory segments is often less expensive than merging from disk (see notes following this table). This threshold influences only the frequency of in-memory merges during the shuffle.
mapreduce.reduce.shuffle.merge.percent	float	The memory threshold for fetched map outputs before an in-memory merge is started, expressed as a percentage of memory allocated to storing map outputs in memory. Since map outputs that can't fit in memory can be stalled, setting this high may decrease parallelism between the fetch and merge. Conversely, values as high as 1.0 have been effective for reduces whose input can fit entirely in memory. This parameter influences only the frequency of in-memory merges during the shuffle.
mapreduce.reduce.shuffle.input.buffer.percent	float	The percentage of memory- relative to the maximum heapsize as typically specified in <code>mapreduce.reduce.java.opts</code> - that can be allocated to storing map outputs during the shuffle. Though some memory should be set aside for the framework, in general it is advantageous to set this high enough to store large and numerous map outputs.
mapreduce.reduce.input.buffer.percent	float	The percentage of memory relative to the maximum heapsize in which map outputs may be retained during the reduce. When the reduce begins, map outputs will be merged to disk until those that remain are under the resource limit this defines. By default, all map outputs are merged to disk before the reduce begins to maximize the memory available to the reduce. For less memory-intensive reduces, this should be increased to avoid trips to disk.

Other notes

- If a map output is larger than 25 percent of the memory allocated to copying map outputs, it will be written directly to disk without first staging through memory.
- When running with a combiner, the reasoning about high merge thresholds and large buffers may not hold. For merges started before all map outputs have been fetched, the combiner is run while spilling to disk. In some cases, one can obtain better reduce times by spending resources combining map outputs- making disk spills small and parallelizing spilling and fetching- rather than aggressively increasing buffer sizes.
- When merging in-memory map outputs to disk to begin the reduce, if an intermediate merge is necessary because there are segments to spill and at least `mapreduce.task.io.sort.factor` segments already on disk, the in-memory map outputs will be part of the intermediate merge.

Configured Parameters

The following properties are localized in the job configuration for each task's execution:



Name	Type	Description
mapreduce.job.id	String	The job id
mapreduce.job.jar	String	job.jar location in job directory
mapreduce.job.local.dir	String	The job specific shared scratch space
mapreduce.task.id	String	The task id
mapreduce.task.attempt.id	String	The task attempt id
mapreduce.task.is.map	boolean	Is this a map task
mapreduce.task.partition	int	The id of the task within the job
mapreduce.map.input.file	String	The filename that the map is reading from
mapreduce.map.input.start	long	The offset of the start of the map input split
mapreduce.map.input.length	long	The number of bytes in the map input split
mapreduce.task.output.dir	String	The task's temporary output directory

Note: During the execution of a streaming job, the names of the "mapreduce" parameters are transformed. The dots (.) become underscores (_). For example, mapreduce.job.id becomes mapreduce_job_id and mapreduce.job.jar becomes mapreduce_job_jar. To get the values in a streaming job's mapper/reducer use the parameter names with the underscores.

Task Logs

The standard output (stdout) and error (stderr) streams and the syslog of the task are read by the NodeManager and logged to \${HADOOP_LOG_DIR}/userlogs.

Distributing Libraries

The [DistributedCache](#) can also be used to distribute both jars and native libraries for use in the map and/or reduce tasks. The child-jvm always has its *current working directory* added to the `java.library.path` and `LD_LIBRARY_PATH`. And hence the cached libraries can be loaded via [System.loadLibrary](#)  or [System.load](#) . More details on how to load shared libraries through distributed cache are documented at [Native Libraries](#).

Job Submission and Monitoring

[Job](#) is the primary interface by which user-job interacts with the ResourceManager.

Job provides facilities to submit jobs, track their progress, access component-tasks' reports and logs, get the MapReduce cluster's status information and so on.

The job submission process involves:

1. Checking the input and output specifications of the job.
2. Computing the InputSplit values for the job.
3. Setting up the requisite accounting information for the DistributedCache of the job, if necessary.
4. Copying the job's jar and configuration to the MapReduce system directory on the FileSystem.
5. Submitting the job to the ResourceManager and optionally monitoring it's status.

Job history files are also logged to user specified directory `mapreduce.jobhistory.intermediate-done-dir` and `mapreduce.jobhistory.done-dir`, which defaults to job output directory.

User can view the history logs summary in specified directory using the following command `$ mapred job -history output.jhist` This command will print job details, failed and killed tip details. More details about the job such as successful tasks and task attempts made for each task can be viewed using the following command `$ mapred job -history all output.jhist`

Normally the user uses Job to create the application, describe various facets of the job, submit the job, and monitor its progress.

Job Control

Users may need to chain MapReduce jobs to accomplish complex tasks which cannot be done via a single MapReduce job. This is fairly easy since the output of the job typically goes to distributed file-system, and the output, in turn, can be used as the input for the next job.

However, this also means that the onus on ensuring jobs are complete (success/failure) lies squarely on the clients. In such cases, the various job-control options are:

- `Job.submit()` : Submit the job to the cluster and return immediately.
- `Job.waitForCompletion(boolean)` : Submit the job to the cluster and wait for it to finish.

Job Input

`InputFormat` describes the input-specification for a MapReduce job.

The MapReduce framework relies on the `InputFormat` of the job to:

1. Validate the input-specification of the job.
2. Split-up the input file(s) into logical `InputSplit` instances, each of which is then assigned to an individual Mapper.
3. Provide the `RecordReader` implementation used to glean input records from the logical `InputSplit` for processing by the Mapper.

The default behavior of file-based `InputFormat` implementations, typically sub-classes of `FileInputFormat`, is to split the input into *logical* `InputSplit` instances based on the total size, in bytes, of the input files. However, the `FileSystem` blocksize of the input files is treated as an upper bound for input splits. A lower bound on the split size can be set via `mapreduce.input.fileinputformat.split.minsize`.

Clearly, logical splits based on input-size is insufficient for many applications since record boundaries must be respected. In such cases, the application should implement a `RecordReader`, who is responsible for respecting record-boundaries and presents a record-oriented view of the logical `InputSplit` to the individual task.

`TextInputFormat` is the default `InputFormat`.

If `TextInputFormat` is the `InputFormat` for a given job, the framework detects input-files with the `.gz` extensions and automatically decompresses them using the appropriate `CompressionCodec`. However, it must be noted that compressed files with the above extensions cannot be *split* and each compressed file is processed in its entirety by a single mapper.

InputSplit

`InputSplit` represents the data to be processed by an individual Mapper.

Typically `InputSplit` presents a byte-oriented view of the input, and it is the responsibility of `RecordReader` to process and present a record-oriented view.

`FileSplit` is the default `InputSplit`. It sets `mapreduce.map.input.file` to the path of the input file for the logical split.

RecordReader

`RecordReader` reads <key, value> pairs from an `InputSplit`.

Typically the `RecordReader` converts the byte-oriented view of the input, provided by the `InputSplit`, and presents a record-oriented to the Mapper implementations for processing. `RecordReader` thus assumes the responsibility of processing record boundaries and presents the tasks with keys and values.

Job Output

`OutputFormat` describes the output-specification for a MapReduce job.

The MapReduce framework relies on the `OutputFormat` of the job to:

1. Validate the output-specification of the job; for example, check that the output directory doesn't already exist.
2. Provide the `RecordWriter` implementation used to write the output files of the job. Output files are stored in a `FileSystem`.

`TextOutputFormat` is the default `OutputFormat`.

OutputCommitter

`OutputCommitter` describes the commit of task output for a MapReduce job.

The MapReduce framework relies on the `OutputCommitter` of the job to:

1. Setup the job during initialization. For example, create the temporary output directory for the job during the initialization of the job. Job setup is done by a separate task when the job is in PREP state and after initializing tasks. Once the setup task completes, the job will be moved to RUNNING state.
2. Cleanup the job after the job completion. For example, remove the temporary output directory after the job completion. Job cleanup is done by a separate task at the end of the job. Job is declared SUCCEEDED/FAILED/KILLED after the cleanup task completes.
3. Setup the task temporary output. Task setup is done as part of the same task, during task initialization.
4. Check whether a task needs a commit. This is to avoid the commit procedure if a task does not need commit.
5. Commit of the task output. Once task is done, the task will commit its output if required.
6. Discard the task commit. If the task has been failed/killed, the output will be cleaned-up. If task could not cleanup (in exception block), a separate task will be launched with same attempt-id to do the cleanup.

`FileOutputCommitter` is the default `OutputCommitter`. Job setup/cleanup tasks occupy map or reduce containers, whichever is available on the `NodeManager`. And `JobCleanup` task, `TaskCleanup` tasks and `JobSetup` task have the highest priority, and in that order.

Task Side-Effect Files

In some applications, component tasks need to create and/or write to side-files, which differ from the actual job-output files.

In such cases there could be issues with two instances of the same Mapper or Reducer running simultaneously (for example, speculative tasks) trying to open and/or write to the same file (path) on the `FileSystem`. Hence the application-writer will have to pick unique names per task-attempt (using the attemptid, say `attempt_200709221812_0001_m_000000_0`), not just per task.

To avoid these issues the MapReduce framework, when the `OutputCommitter` is `FileOutputCommitter`, maintains a special `${mapreduce.output.fileoutputformat.outputdir}/_temporary/_${taskid}` sub-directory accessible via `${mapreduce.task.output.dir}` for each task-attempt on the `FileSystem` where the output of the task-attempt is stored. On successful completion of the task-attempt, the files in the `${mapreduce.output.fileoutputformat.outputdir}/_temporary/_${taskid}` (only) are *promoted* to `${mapreduce.output.fileoutputformat.outputdir}`. Of course, the framework discards the sub-directory of unsuccessful task-attempts. This process is completely transparent to the application.

The application-writer can take advantage of this feature by creating any side-files required in `${mapreduce.task.output.dir}` during execution of a task via `FileOutputFormat.getWorkOutputPath(Conext)`, and the framework will promote them similarly for successful task-attempts, thus eliminating the need to pick unique paths per task-attempt.

Note: The value of `${mapreduce.task.output.dir}` during execution of a particular task-attempt is actually `${mapreduce.output.fileoutputformat.outputdir}/_temporary/_${taskid}`, and this value is set by the MapReduce framework. So, just create any side-files in the path returned by `FileOutputFormat.getWorkOutputPath(Conext)` from MapReduce task to take advantage of this feature.

The entire discussion holds true for maps of jobs with `reducer=NONE` (i.e. 0 reducers) since output of the map, in that case, goes directly to HDFS.

RecordWriter

`RecordWriter` writes the output `<key, value>` pairs to an output file.

RecordWriter implementations write the job outputs to the FileSystem.

➔ Other Useful Features

Submitting Jobs to Queues

Users submit jobs to Queues. Queues, as collection of jobs, allow the system to provide specific functionality. For example, queues use ACLs to control which users who can submit jobs to them. Queues are expected to be primarily used by Hadoop Schedulers.

Hadoop comes configured with a single mandatory queue, called 'default'. Queue names are defined in the `mapreduce.job.queueName` property of the Hadoop site configuration. Some job schedulers, such as the [Capacity Scheduler](#), support multiple queues.

A job defines the queue it needs to be submitted to through the `mapreduce.job.queueName` property, or through the `Configuration.set(MRJobConfig.QUEUE_NAME, String)` API. Setting the queue name is optional. If a job is submitted without an associated queue name, it is submitted to the 'default' queue.

Counters



Counters represent global counters, defined either by the MapReduce framework or applications. Each Counter can be of any Enum type. Counters of a particular Enum are bunched into groups of type `Counters.Group`.

Applications can define arbitrary Counters (of type Enum) and update them via [Counters.incrCounter\(Enum, long\)](#) or `Counters.incrCounter(String, String, long)` in the map and/or reduce methods. These counters are then globally aggregated by the framework.

DistributedCache

DistributedCache distributes application-specific, large, read-only files efficiently.

DistributedCache is a facility provided by the MapReduce framework to cache files (text, archives, jars and so on) needed by applications.

Applications specify the files to be cached via urls ([hdfs://](#)  in the Job. The DistributedCache assumes that the files specified via [hdfs://](#)  urls are already present on the FileSystem.

The framework will copy the necessary files to the worker node before any tasks for the job are executed on that node. Its efficiency stems from the fact that the files are only copied once per job and the ability to cache archives which are un-archived on the workers.

DistributedCache tracks the modification timestamps of the cached files. Clearly the cache files should not be modified by the application or externally while the job is executing.

DistributedCache can be used to distribute simple, read-only data/text files and more complex types such as archives and jars. Archives (zip, tar, tgz and tar.gz files) are *un-archived* at the worker nodes. Files have *execution permissions* set.

The files/archives can be distributed by setting the property `mapreduce.job.cache.{files | archives}`. If more than one file/archive has to be distributed, they can be added as comma separated paths. The properties can also be set by APIs [Job.addCacheFile\(URI\)](#)/ [Job.addCacheArchive\(URI\)](#) and [Job.setCacheFiles\(URI\[\]\)](#)/ [Job.setCacheArchives\(URI\[\]\)](#) where URI is of the form `hdfs://host:port/absolute-path#link-name`. In Streaming, the files can be distributed through command line option `-cacheFile/-cacheArchive`.

The DistributedCache can also be used as a rudimentary software distribution mechanism for use in the map and/or reduce tasks. It can be used to distribute both jars and native libraries. The [Job.addArchiveToClassPath\(Path\)](#) or [Job.addFileToClassPath\(Path\)](#) api can be used to cache files/jars and also add them to the *classpath* of child-jvm. The same can be done by setting the configuration properties `mapreduce.job.classpath.{files | archives}`. Similarly the cached files that are symlinked into the working directory of the task can be used to distribute native libraries and load them.

Private and Public DistributedCache Files

DistributedCache files can be private or public, that determines how they can be shared on the worker nodes.

- "Private" DistributedCache files are cached in a `localdirectory` private to the user whose jobs need these files. These files are shared by all tasks and jobs of the specific user only and cannot be accessed by jobs of other

users on the workers. A DistributedCache file becomes private by virtue of its permissions on the file system where the files are uploaded, typically HDFS. If the file has no world readable access, or if the directory path leading to the file has no world executable access for lookup, then the file becomes private.

- “Public” DistributedCache files are cached in a global directory and the file access is setup such that they are publicly visible to all users. These files can be shared by tasks and jobs of all users on the workers. A DistributedCache file becomes public by virtue of its permissions on the file system where the files are uploaded, typically HDFS. If the file has world readable access, AND if the directory path leading to the file has world executable access for lookup, then the file becomes public. In other words, if the user intends to make a file publicly available to all users, the file permissions must be set to be world readable, and the directory permissions on the path leading to the file must be world executable.

Profiling

Profiling is a utility to get a representative (2 or 3) sample of built-in java profiler for a sample of maps and reduces.

User can specify whether the system should collect profiler information for some of the tasks in the job by setting the configuration property `mapreduce.task.profile`. The value can be set using the api `Configuration.set(MRJobConfig.TASK_PROFILE, boolean)`. If the value is set `true`, the task profiling is enabled. The profiler information is stored in the user log directory. By default, profiling is not enabled for the job.

Once user configures that profiling is needed, she/he can use the configuration property `mapreduce.task.profile.{maps|reduces}` to set the ranges of MapReduce tasks to profile. The value can be set using the api `Configuration.set(MRJobConfig.NUM_{MAP|REDUCE}_PROFILES, String)`. By default, the specified range is 0–2.

User can also specify the profiler configuration arguments by setting the configuration property `mapreduce.task.profile.params`. The value can be specified using the api `Configuration.set(MRJobConfig.TASK_PROFILE_PARAMS, String)`. If the string contains a `%s`, it will be replaced with the name of the profiling output file when the task runs. These parameters are passed to the task child JVM on the command line. The default value for the profiling parameters is –
`agentlib:hprof=cpu=samples,heap=sites,force=n,thread=y,verbose=n,file=%s`.

Debugging

The MapReduce framework provides a facility to run user-provided scripts for debugging. When a MapReduce task fails, a user can run a debug script, to process task logs for example. The script is given access to the task’s stdout and stderr outputs, syslog and jobconf. The output from the debug script’s stdout and stderr is displayed on the console diagnostics and also as part of the job UI.

In the following sections we discuss how to submit a debug script with a job. The script file needs to be distributed and submitted to the framework.

How to distribute the script file:

The user needs to use [DistributedCache](#) to *distribute* and *symlink* to the script file.

How to submit the script:

A quick way to submit the debug script is to set values for the properties `mapreduce.map.debug.script` and `mapreduce.reduce.debug.script`, for debugging map and reduce tasks respectively. These properties can also be set by using APIs [Configuration.set\(MRJobConfig.MAP_DEBUG_SCRIPT, String\)](#) and [Configuration.set\(MRJobConfig.REDUCE_DEBUG_SCRIPT, String\)](#). In streaming mode, a debug script can be submitted with the command-line options `–mapdebug` and `–reduceddebug`, for debugging map and reduce tasks respectively.






The arguments to the script are the task’s stdout, stderr, syslog and jobconf files. The debug command, run on the node where the MapReduce task failed, is:
`$script $stdout $stderr $syslog $jobconf`

Pipes programs have the `c++` program name as a fifth argument for the command. Thus for the pipes programs the command is
`$script $stdout $stderr $syslog $jobconf $program`

Default Behavior:

For pipes, a default script is run to process core dumps under `gdb`, prints stack trace and gives info about running threads.

Data Compression

Hadoop MapReduce provides facilities for the application-writer to specify compression for both intermediate map-outputs and the job-outputs i.e. output of the reduces. It also comes bundled with [CompressionCodec](#) implementation for the [zlib](#)  compression algorithm. The [gzip](#) , [bzip2](#) , [snappy](#) , and [lz4](#)  file format are also supported.

Hadoop also provides native implementations of the above compression codecs for reasons of both performance (zlib) and non-availability of Java libraries. More details on their usage and availability are available [here](#).

Intermediate Outputs

Applications can control compression of intermediate map-outputs via the `Configuration.set(MRJobConfig.MAP_OUTPUT_COMPRESS, boolean)` api and the `CompressionCodec` to be used via the `Configuration.set(MRJobConfig.MAP_OUTPUT_COMPRESS_CODEC, Class)` api.

Job Outputs

Applications can control compression of job-outputs via the `FileOutputFormat.setCompressOutput(Job, boolean)` api and the `CompressionCodec` to be used can be specified via the `FileOutputFormat.setOutputCompressorClass(Job, Class)` api.

If the job outputs are to be stored in the [SequenceFileOutputFormat](#), the required `SequenceFile.CompressionType` (i.e. `RECORD` / `BLOCK` - defaults to `RECORD`) can be specified via the `SequenceFileOutputFormat.setOutputCompressionType(Job, SequenceFile.CompressionType)` api.

Skipping Bad Records

Hadoop provides an option where a certain set of bad input records can be skipped when processing map inputs. Applications can control this feature through the [SkipBadRecords](#) class.

This feature can be used when map tasks crash deterministically on certain input. This usually happens due to bugs in the map function. Usually, the user would have to fix these bugs. This is, however, not possible sometimes. The bug may be in third party libraries, for example, for which the source code is not available. In such cases, the task never completes successfully even after multiple attempts, and the job fails. With this feature, only a small portion of data surrounding the bad records is lost, which may be acceptable for some applications (those performing statistical analysis on very large data, for example).

By default this feature is disabled. For enabling it, refer to [SkipBadRecords.setMapperMaxSkipRecords\(Configuration, long\)](#) and [SkipBadRecords.setReducerMaxSkipGroups\(Configuration, long\)](#).

With this feature enabled, the framework gets into 'skipping mode' after a certain number of map failures. For more details, see [SkipBadRecords.setAttemptsToStartSkipping\(Configuration, int\)](#). In 'skipping mode', map tasks maintain the range of records being processed. To do this, the framework relies on the processed record counter. See [SkipBadRecords.COUNTER_MAP_PROCESSED_RECORDS](#) and [SkipBadRecords.COUNTER_REDUCE_PROCESSED_GROUPS](#). This counter enables the framework to know how many records have been processed successfully, and hence, what record range caused a task to crash. On further attempts, this range of records is skipped.

The number of records skipped depends on how frequently the processed record counter is incremented by the application. It is recommended that this counter be incremented after every record is processed. This may not be possible in some applications that typically batch their processing. In such cases, the framework may skip additional records surrounding the bad record. Users can control the number of skipped records through [SkipBadRecords.setMapperMaxSkipRecords\(Configuration, long\)](#) and [SkipBadRecords.setReducerMaxSkipGroups\(Configuration, long\)](#). The framework tries to narrow the range of skipped records using a binary search-like approach. The skipped range is divided into two halves and only one half gets executed. On subsequent failures, the framework figures out which half contains bad records. A task will be re-executed till the acceptable skipped value is met or all task attempts are exhausted. To increase the number of task attempts, use [Job.setMaxMapAttempts\(int\)](#) and [Job.setMaxReduceAttempts\(int\)](#)

Skipped records are written to HDFS in the sequence file format, for later analysis. The location can be changed through [SkipBadRecords.setSkipOutputPath\(JobConf, Path\)](#).

Example: WordCount v2.0

Here is a more complete WordCount which uses many of the features provided by the MapReduce framework we discussed so far.

This needs the HDFS to be up and running, especially for the DistributedCache-related features. Hence it only works with a [pseudo-distributed](#) or [fully-distributed](#) Hadoop installation.

Source Code

```
import java.io.BufferedReader;
import java.io.FileReader;
import java.io.IOException;
import java.net.URI;
import java.util.ArrayList;
import java.util.HashSet;
import java.util.List;
import java.util.Set;
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;
import org.apache.hadoop.mapreduce.Counter;
import org.apache.hadoop.util.GenericOptionsParser;
import org.apache.hadoop.util.StringUtils;

public class WordCount2 {

    public static class TokenizerMapper
        extends Mapper<Object, Text, Text, IntWritable>{

        static enum CountersEnum { INPUT_WORDS }

        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();

        private boolean caseSensitive;
        private Set<String> patternsToSkip = new HashSet<String>();

        private Configuration conf;
        private BufferedReader fis;

        @Override
        public void setup(Context context) throws IOException,
            InterruptedException {
            conf = context.getConfiguration();
            caseSensitive = conf.getBoolean("wordcount.case.sensitive", true);
            if (conf.getBoolean("wordcount.skip.patterns", false)) {
                URI[] patternsURIs = Job.getInstance(conf).getCacheFiles();
                for (URI patternsURI : patternsURIs) {
                    Path patternsPath = new Path(patternsURI.getPath());
                    String patternsFileName = patternsPath.getName().toString();
                    parseSkipFile(patternsFileName);
                }
            }
        }

        private void parseSkipFile(String fileName) {
            try {
                fis = new BufferedReader(new FileReader(fileName));
                String pattern = null;
                while ((pattern = fis.readLine()) != null) {
                    patternsToSkip.add(pattern);
                }
            }
        }
    }
}
```

```

    } catch (IOException ioe) {
        System.err.println("Caught exception while parsing the cached file '"
            + StringUtils.stringifyException(ioe));
    }
}

@Override
public void map(Object key, Text value, Context context
    ) throws IOException, InterruptedException {
    String line = (caseSensitive) ?
        value.toString() : value.toString().toLowerCase();
    for (String pattern : patternsToSkip) {
        line = line.replaceAll(pattern, "");
    }
    StringTokenizer itr = new StringTokenizer(line);
    while (itr.hasMoreTokens()) {
        word.set(itr.nextToken());
        context.write(word, one);
        Counter counter = context.getCounter(CountersEnum.class.getName(),
            CountersEnum.INPUT_WORDS.toString());
        counter.increment(1);
    }
}

public static class IntSumReducer
    extends Reducer<Text,IntWritable,Text,IntWritable> {
    private IntWritable result = new IntWritable();

    public void reduce(Text key, Iterable<IntWritable> values,
        Context context
        ) throws IOException, InterruptedException {

        int sum = 0;
        for (IntWritable val : values) {
            sum += val.get();
        }
        result.set(sum);
        context.write(key, result);
    }
}

public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    GenericOptionsParser optionParser = new GenericOptionsParser(conf, args);
    String[] remainingArgs = optionParser.getRemainingArgs();
    if ((remainingArgs.length != 2) && (remainingArgs.length != 4)) {
        System.err.println("Usage: wordcount <in> <out> [-skip skipPatternFile]");
        System.exit(2);
    }
    Job job = Job.getInstance(conf, "word count");
    job.setJarByClass(WordCount2.class);
    job.setMapperClass(TokenizerMapper.class);
    job.setCombinerClass(IntSumReducer.class);
    job.setReducerClass(IntSumReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);

    List<String> otherArgs = new ArrayList<String>();
    for (int i=0; i < remainingArgs.length; ++i) {
        if ("--skip".equals(remainingArgs[i])) {
            job.addCacheFile(new Path(remainingArgs[++i]).toUri());
            job.getConfiguration().setBoolean("wordcount.skip.patterns", true);
        } else {
            otherArgs.add(remainingArgs[i]);
        }
    }
    FileInputFormat.addInputPath(job, new Path(otherArgs.get(0)));
    FileOutputFormat.setOutputPath(job, new Path(otherArgs.get(1)));
}

```

```
        System.exit(job.waitForCompletion(true) ? 0 : 1);  
    }  
}
```

Sample Runs

Sample text-files as input:

```
$ bin/hadoop fs -ls /user/joe/wordcount/input/  
/user/joe/wordcount/input/file01  
/user/joe/wordcount/input/file02  
  
$ bin/hadoop fs -cat /user/joe/wordcount/input/file01  
Hello World, Bye World!  
  
$ bin/hadoop fs -cat /user/joe/wordcount/input/file02  
Hello Hadoop, Goodbye to hadoop.
```

Run the application:

```
$ bin/hadoop jar wc.jar WordCount2 /user/joe/wordcount/input /user/joe/wordcount/output
```

Output:

```
$ bin/hadoop fs -cat /user/joe/wordcount/output/part-r-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.

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

Run it again, this time with more options:

```
$ bin/hadoop jar wc.jar WordCount2 -Dwordcount.case.sensitive=true /user/joe/wordco
```

As expected, the output:

```
$ bin/hadoop fs -cat /user/joe/wordcount/output/part-r-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 wc.jar WordCount2 -Dwordcount.case.sensitive=false /user/joe/wordco
```

Sure enough, the output:

```
$ bin/hadoop fs -cat /user/joe/wordcount/output/part-r-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 setup method of the Mapper (and Reducer) implementations.
- 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.
- Demonstrates the utility of the GenericOptionsParser to handle generic Hadoop command-line options.
- Demonstrates how applications can use Counters and how they can set application-specific status information passed to the map (and reduce) method.

Java and JNI are trademarks or registered trademarks of Oracle America, Inc. in the United States and other countries.