

## Why Use ToolRunner?

---

- You can use `ToolRunner` in MapReduce driver classes
  - This is not required, but is a best practice
- `ToolRunner` uses the `GenericOptionsParser` class internally
  - Allows you to specify configuration options on the command line
  - Also allows you to specify items for the Distributed Cache on the command line (see later)

# How to Implement ToolRunner: Complete Driver

```
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.Job; import org.apache.hadoop.conf.Configured;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;

public class WordCount extends Configured implements Tool {

    public static void main(String[] args) throws Exception {
        int exitCode = ToolRunner.run(new Configuration(), new WordCount(), args);
        System.exit(exitCode);
    }

    public int run(String[] args) throws Exception {
        if (args.length != 2) {
            System.out.printf(
                "Usage: %s [generic options] <input dir> <output dir>\n", getClass().getSimpleName());
            return -1;
        }
        Job job = new Job(getConf());
        job.setJarByClass(WordCount.class);
        job.setJobName("Word Count");

        FileInputFormat.setInputPaths(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.setMapperClass(WordMapper.class);
        job.setReducerClass(SumReducer.class);
        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(IntWritable.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        boolean success = job.waitForCompletion(true);
        return success ? 0 : 1;
    }
}
```

# How to Implement ToolRunner: Imports

```
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.conf.Configured;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;
```

```
public class WordCount
```

```
    public static void main(String[] args) throws Exception {
        int exitCode = ToolRunner.run(new WordCount(), args, new Configuration());
        System.exit(exitCode);
    }
```

```
    public int run(String[] args) throws Exception {
        if (args.length != 2) {
            System.out.printf(
                "Usage: %s [generic options] <input dir> <output dir>\n",
                getClass().getSimpleName());
            return -1;
        }
        Job job = new Job(getConf());
```

Import the relevant classes. We omit the `import` statements in future slides for brevity.



## How to Implement ToolRunner: Driver Class Definition

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

```
    public static void main(String[] args) {
        int exitCode = ToolRunner.run(WordCount.class, args);
        System.exit(exitCode);
    }
```

The driver class implements the `Tool` interface and extends the `Configured` class.

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

```
        if (args.length != 2) {
            System.out.printf(
                "Usage: %s [generic options] <input dir> <output dir>\n", getClass().getSimpleName());
            return -1;
        }
```

```
        Job job = new Job(getConf());
        job.setJarByClass(WordCount.class); job.setJobName("Word Count");
```

```
        FileInputFormat.setInputPaths(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
```

```
        job.setMapperClass(WordMapper.class);
        job.setReducerClass(SumReducer.class);
        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(IntWritable.class);
```

```
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
```

```
        boolean success = job.waitForCompletion(true);
        return success ? 0 : 1;
```

```
    }
```

## How to Implement ToolRunner: Main Method

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

```
    public static void main(String[] args) throws Exception {  
        int exitCode = ToolRunner.run(new Configuration(),  
            new WordCount(), args);  
        System.exit(exitCode);  
    }
```

```
    public int run(String[] args) throws Exception {  
        if (args.length != 2)  
            System.out.printf("Usage: %s [options] %s\n",  
                getClass().getSimpleName(), this)  
            return -1;  
    }
```

The driver main method calls `ToolRunner.run`.

```
    Job job = new Job(getConf());  
    job.setJarByClass(WordCount.class);  
    job.setJobName("Word Count");  
  
    FileInputFormat.setInputPaths(job, new Path(args[0]));  
    FileOutputFormat.setOutputPath(job, new Path(args[1]));  
  
    job.setMapperClass(WordMapper.class);  
    job.setReducerClass(SumReducer.class);  
    job.setMapOutputKeyClass(Text.class);  
    job.setMapOutputValueClass(IntWritable.class);  
  
    job.setOutputKeyClass(Text.class);  
    job.setOutputValueClass(IntWritable.class);  
  
    ...
```

## How to Implement ToolRunner: Run Method

```
public class WordCount {  
  
    public static void main(  
        int exitCode = ToolRunner.  
            new WordCount().run(  
                System.exit(exitCode)  
            )  
    }
```

The driver `run` method creates, configures, and submits the job.

```
    public int run(String[] args) throws Exception {  
        if (args.length != 2) {  
            System.out.printf(  
                "Usage: %s [generic options] <input dir> <output dir>\n", getClass().getSimpleName());  
            return -1;  
        }  
        Job job = new Job(getConf());  
        job.setJarByClass(WordCount.class);  
        job.setJobName("Word Count");  
  
        FileInputFormat.setInputPaths(job, new Path(args[0]));  
        FileOutputFormat.setOutputPath(job, new Path(args[1]));  
  
        ...  
    }
```



## ToolRunner Command Line Options

---

- ToolRunner allows the user to specify configuration options on the command line
- Commonly used to specify Hadoop properties using the `-D` flag
  - Will override any default or site properties in the configuration
  - But will *not* override those set in the driver code

```
$ hadoop jar myjar.jar MyDriver \  
-D mapred.reduce.tasks=10 myinputdir myoutputdir
```

- Note that `-D` options must appear before any additional program arguments
- Can specify an XML configuration file with `-conf`
- Can specify the default filesystem with `-fs uri`
  - Shortcut for `-D fs.default.name=uri`

## **Setting Up and Tearing Down Mappers and Reducers**



## The setup Method

---

- It is common to want your Mapper or Reducer to execute some code before the map or reduce method is called for the first time
  - Initialize data structures
  - Read data from an external file
  - Set parameters
- The setup method is run before the map or reduce method is called for the first time

```
public void setup(Context context)
```

## The `cleanup` Method

---

- Similarly, you may wish to perform some action(s) after all the records have been processed by your Mapper or Reducer
- The `cleanup` method is called before the Mapper or Reducer terminates

```
public void cleanup(Context context) throws  
    IOException, InterruptedException
```

## Passing Parameters

```
public class MyDriverClass {
    public int main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        conf.setInt ("paramname",value);
        Job job = new Job(conf);
        ...
        boolean success = job.waitForCompletion(true);
        return success ? 0 : 1;
    }
}
```

```
public class MyMapper extends Mapper {

    public void setup(Context context) {
        Configuration conf = context.getConfiguration();
        int myParam = conf.getInt("paramname", 0);
        ...
    }
    public void map...
}
```



## Accessing HDFS Programmatically

---

- **In addition to using the command-line shell, you can access HDFS programmatically**
  - Useful if your code needs to read or write ‘side data’ in addition to the standard MapReduce inputs and outputs
  - Or for programs outside of Hadoop which need to read the results of MapReduce jobs
- **Beware: HDFS is not a general-purpose filesystem!**
  - Files cannot be modified once they have been written, for example
- **Hadoop provides the `FileSystem` abstract base class**
  - Provides an API to generic file systems
    - Could be HDFS
    - Could be your local file system
    - Could even be, for example, Amazon S3

## The FileSystem API (1)

---

- In order to use the `FileSystem` API, retrieve an instance of it

```
Configuration conf = new Configuration();  
FileSystem fs = FileSystem.get(conf);
```

- The `conf` object has read in the Hadoop configuration files, and therefore knows the address of the `NameNode`
- A file in HDFS is represented by a `Path` object

```
Path p = new Path("/path/to/my/file");
```

## The FileSystem API (2)

---

- **Some useful API methods:**

- `FSDataOutputStream create(...)`
  - **Extends** `java.io.DataOutputStream`
  - **Provides methods for writing primitives, raw bytes etc**
- `FSDataInputStream open(...)`
  - **Extends** `java.io.DataInputStream`
  - **Provides methods for reading primitives, raw bytes etc**
- `boolean delete(...)`
- `boolean mkdirs(...)`
- `void copyFromLocalFile(...)`
- `void copyToLocalFile(...)`
- `FileStatus[] listStatus(...)`



## The FileSystem API: Directory Listing

---

- Get a directory listing:

```
Path p = new Path("/my/path");

Configuration conf = new Configuration();
FileSystem fs = FileSystem.get(conf);
FileStatus[] fileStats = fs.listStatus(p);

for (int i = 0; i < fileStats.length; i++) {
    Path f = fileStats[i].getPath();

    // do something interesting
}
```

## The FileSystem API: Writing Data

---

- Write data to a file

```
Configuration conf = new Configuration();
FileSystem fs = FileSystem.get(conf);

Path p = new Path("/my/path/foo");

FSDataOutputStream out = fs.create(p, false);

// write some raw bytes
out.write(getBytes());

// write an int
out.writeInt(getInt());

...

out.close();
```

## The Distributed Cache: Motivation

---

- **A common requirement is for a Mapper or Reducer to need access to some 'side data'**
  - Lookup tables
  - Dictionaries
  - Standard configuration values
- **One option: read directly from HDFS in the setup method**
  - Using the API seen in the previous section
  - Works, but is not scalable
- **The Distributed Cache provides an API to push data to all slave nodes**
  - Transfer happens behind the scenes before any task is executed
  - Data is only transferred once to each node, rather
  - Note: Distributed Cache is read-only
  - Files in the Distributed Cache are automatically deleted from slave nodes when the job finishes



## Using the Distributed Cache: The Difficult Way

---

- Place the files into HDFS
- Configure the Distributed Cache in your driver code

```
Configuration conf = new Configuration();  
DistributedCache.addCacheFile(new URI("/myapp/lookup.dat"), conf);  
DistributedCache.addFileToClassPath(new Path("/myapp/mylib.jar"), conf);  
DistributedCache.addCacheArchive(new URI("/myapp/map.zip", conf));  
DistributedCache.addCacheArchive(new URI("/myapp/mytar.tar", conf));  
DistributedCache.addCacheArchive(new URI("/myapp/mytgz.tgz", conf));  
DistributedCache.addCacheArchive(new URI("/myapp/mytargz.tar.gz", conf));
```

- .jar files added with `addFileToClassPath` will be added to your Mapper or Reducer's classpath
- Files added with `addCacheArchive` will automatically be dearchived/decompressed

## Using the DistributedCache: The Easy Way

---

- If you are using ToolRunner, you can add files to the Distributed Cache directly from the command line when you run the job
  - No need to copy the files to HDFS first
- Use the `-files` option to add files

```
hadoop jar myjar.jar MyDriver -files file1, file2, file3, ...
```

- The `-archives` flag adds archived files, and automatically unarchives them on the destination machines
- The `-libjars` flag adds jar files to the classpath

## Accessing Files in the Distributed Cache

---

- Files added to the Distributed Cache are made available in your task's local working directory
  - Access them from your Mapper or Reducer the way you would read any ordinary local file

```
File f = new File("file_name_here");
```



## Reusable Classes for the New API

---

- The `org.apache.hadoop.mapreduce.lib.*/*` packages contain a library of Mappers, Reducers, and Partitioners supporting the new API
- **Example classes:**
  - `InverseMapper` – Swaps keys and values
  - `RegexMapper` – Extracts text based on a regular expression
  - `IntSumReducer`, `LongSumReducer` – Add up all values for a key

# Strategies for Debugging MapReduce Code

## Introduction to Debugging

---

- **Debugging MapReduce code is difficult!**
  - Each instance of a Mapper runs as a separate task
    - Often on a different machine
  - Difficult to attach a debugger to the process
  - Difficult to catch 'edge cases'
- **Very large volumes of data mean that unexpected input is likely to appear**
  - Code which expects all data to be well-formed is likely to fail



## Common-Sense Debugging Tips

---

- **Code defensively**
  - Ensure that input data is in the expected format
  - Expect things to go wrong
  - Catch exceptions
- **Start small, build incrementally**
- **Make as much of your code as possible Hadoop-agnostic**
  - Makes it easier to test
- **Write unit tests**
- **Test locally whenever possible**
  - With small amounts of data
- **Then test in pseudo-distributed mode**
- **Finally, test on the cluster**

## Testing Strategies

---

- **When testing in pseudo-distributed mode, ensure that you are testing with a similar environment to that on the real cluster**
  - Same amount of RAM allocated to the task JVMs
  - Same version of Hadoop
  - Same version of Java
  - Same versions of third-party libraries

## **Testing MapReduce Code Locally Using LocalJobRunner**



## Testing Locally (1)

---

- **Hadoop can run MapReduce in a single, local process**
  - Does not require any Hadoop daemons to be running
  - Uses the local filesystem instead of HDFS
  - Known as LocalJobRunner mode
- **This is a very useful way of quickly testing incremental changes to code**

## Testing Locally (2)

- To run in LocalJobRunner mode, add the following lines to the driver code:

```
Configuration conf = new Configuration();  
conf.set("mapred.job.tracker", "local");  
conf.set("fs.default.name", "file:///");
```

- Or set these options on the command line if your driver uses ToolRunner

- `-fs` is equivalent to `-D fs.default.name`
- `-jt` is equivalent to `-D mapred.job.tracker`
- e.g.

```
$ hadoop jar myjar.jar MyDriver -fs=file:/// -jt=local \  
indir outdir
```

## Testing Locally (3)

---

- **Some limitations of LocalJobRunner mode:**
  - Distributed Cache does not work
  - The job can only specify a single Reducer
  - Some 'beginner' mistakes may not be caught
    - For example, attempting to share data between Mappers will work, because the code is running in a single JVM



## LocalJobRunner Mode in Eclipse (1)

---

- **Eclipse on the course VM runs Hadoop code in LocalJobRunner mode from within the IDE**
  - This is Hadoop's default behavior when no configuration is provided
- **This allows rapid development iterations**
  - 'Agile programming'

## LocalJobRunner Mode in Eclipse (2)

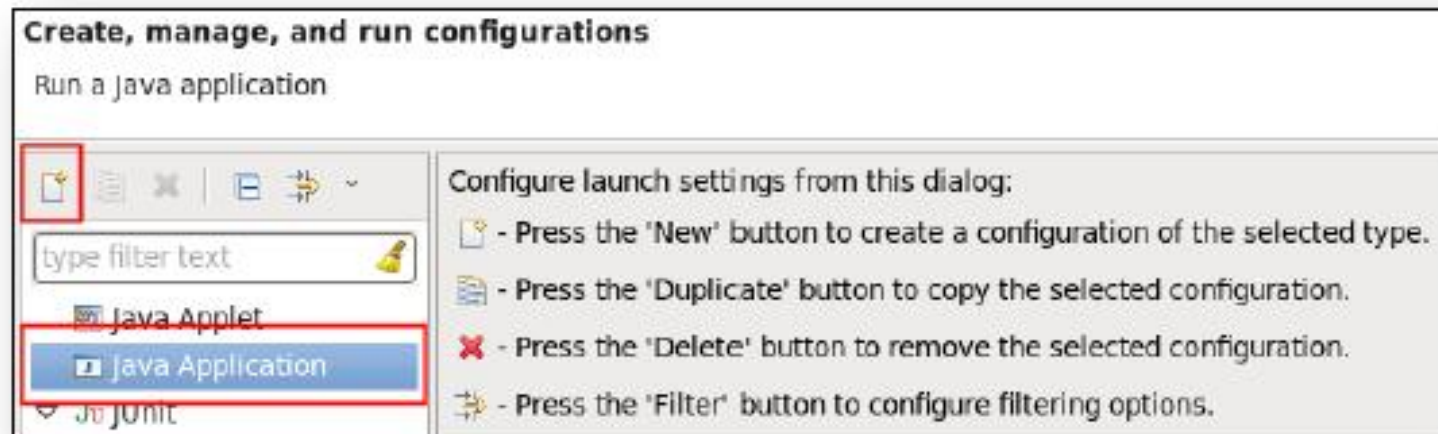
---

- Specify a Run Configuration

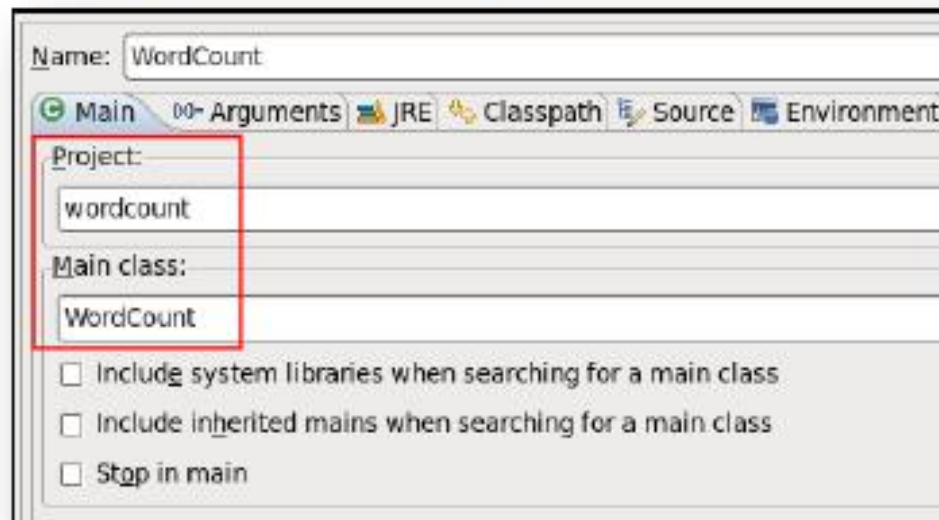


## LocalJobRunner Mode in Eclipse (3)

- **Select Java Application, then select the New button**



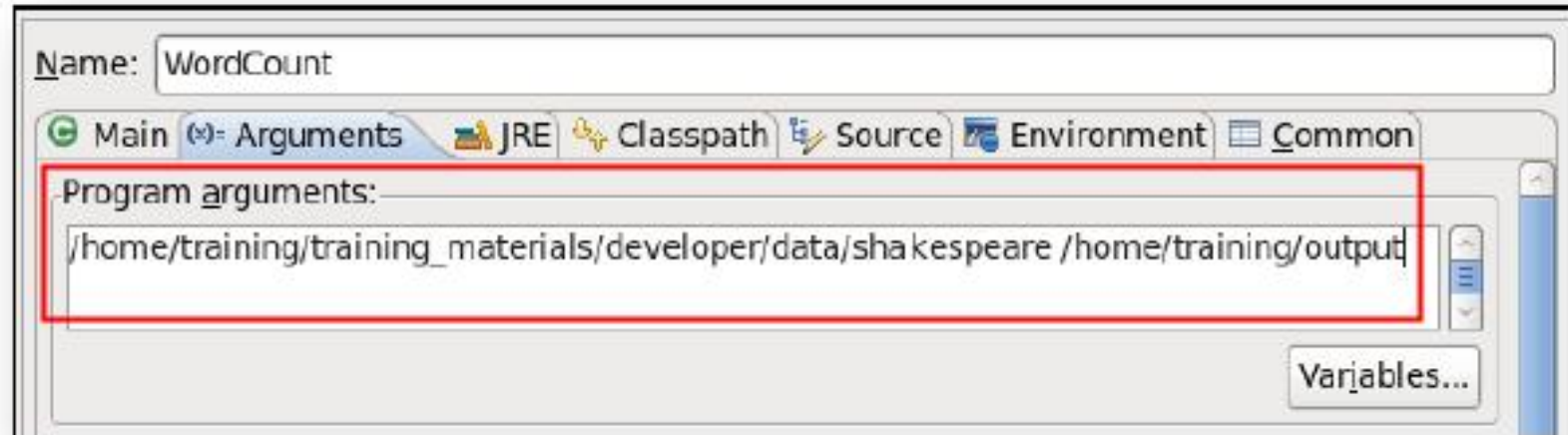
- **Verify that the Project and Main Class fields are pre-filled correctly**





## LocalJobRunner Mode in Eclipse (4)

- **Specify values in the Arguments tab**
  - Local input and output files
  - Any configuration options needed when your job runs

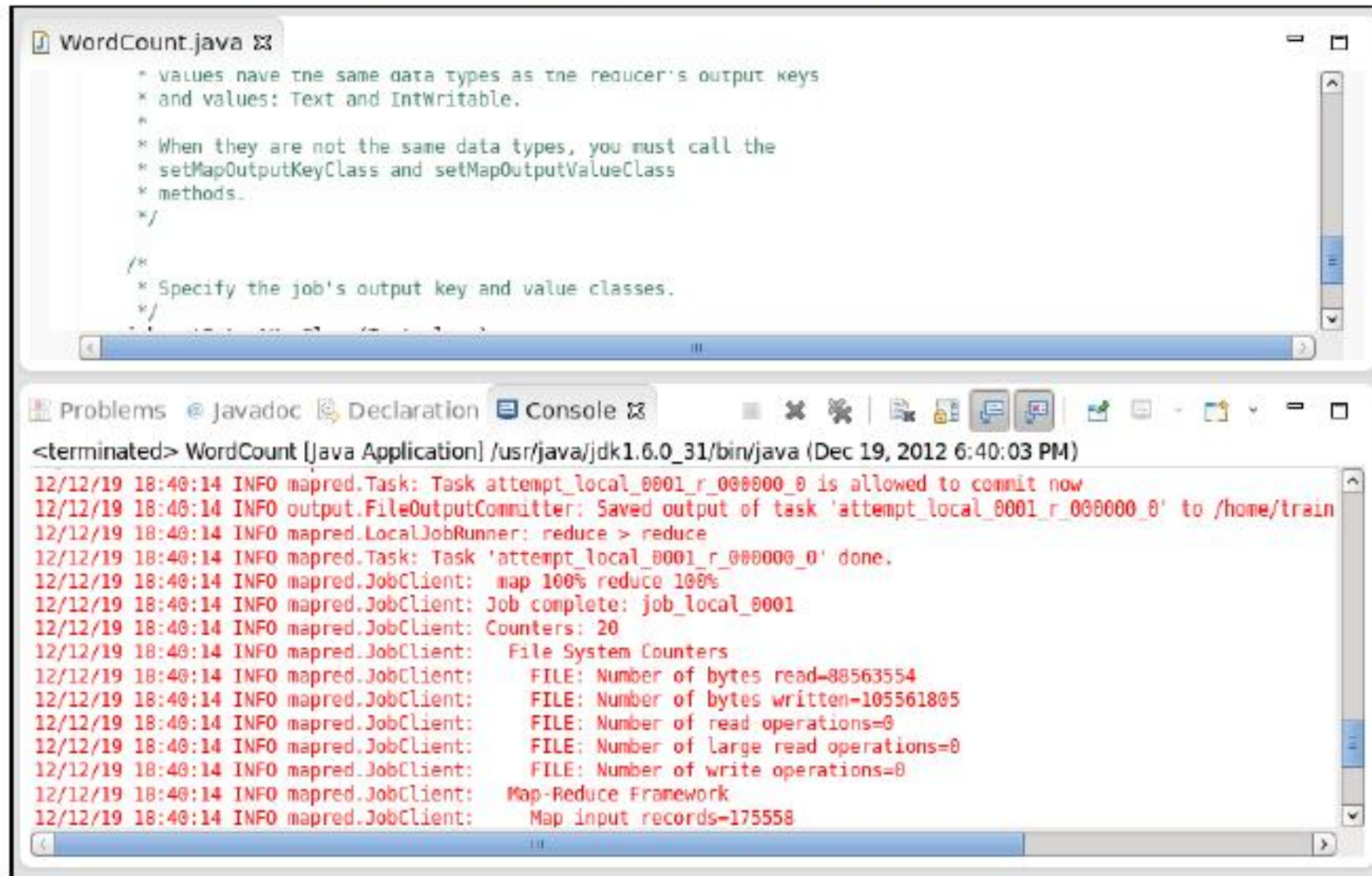


- **Define breakpoints if desired**
- **Execute the application in run mode or debug mode**



## LocalJobRunner Mode in Eclipse (5)

- Review output in the Eclipse console window



The screenshot shows the Eclipse IDE interface. The top editor window displays the `WordCount.java` file with Javadoc comments. The bottom console window shows the output of the program execution.

```
WordCount.java
/*
 * values have the same data types as the reducer's output keys
 * and values: Text and IntWritable.
 *
 * When they are not the same data types, you must call the
 * setMapOutputKeyClass and setMapOutputValueClass
 * methods.
 */
/*
 * Specify the job's output key and value classes.
 */
```

Problems Javadoc Declaration Console

```
<terminated> WordCount [Java Application] /usr/java/jdk1.6.0_31/bin/java (Dec 19, 2012 6:40:03 PM)
12/12/19 18:40:14 INFO mapred.Task: Task attempt_local_0001_r_000000_0 is allowed to commit now
12/12/19 18:40:14 INFO output.FileOutputCommitter: Saved output of task 'attempt_local_0001_r_000000_0' to /home/train
12/12/19 18:40:14 INFO mapred.LocalJobRunner: reduce > reduce
12/12/19 18:40:14 INFO mapred.Task: Task 'attempt_local_0001_r_000000_0' done.
12/12/19 18:40:14 INFO mapred.JobClient: map 100% reduce 100%
12/12/19 18:40:14 INFO mapred.JobClient: Job complete: job_local_0001
12/12/19 18:40:14 INFO mapred.JobClient: Counters: 20
12/12/19 18:40:14 INFO mapred.JobClient: File System Counters
12/12/19 18:40:14 INFO mapred.JobClient: FILE: Number of bytes read=88563554
12/12/19 18:40:14 INFO mapred.JobClient: FILE: Number of bytes written=105561805
12/12/19 18:40:14 INFO mapred.JobClient: FILE: Number of read operations=0
12/12/19 18:40:14 INFO mapred.JobClient: FILE: Number of large read operations=0
12/12/19 18:40:14 INFO mapred.JobClient: FILE: Number of write operations=0
12/12/19 18:40:14 INFO mapred.JobClient: Map-Reduce Framework
12/12/19 18:40:14 INFO mapred.JobClient: Map input records=175558
```

# Writing and Viewing Log Files



## Before Logging: `stdout` and `stderr`

---

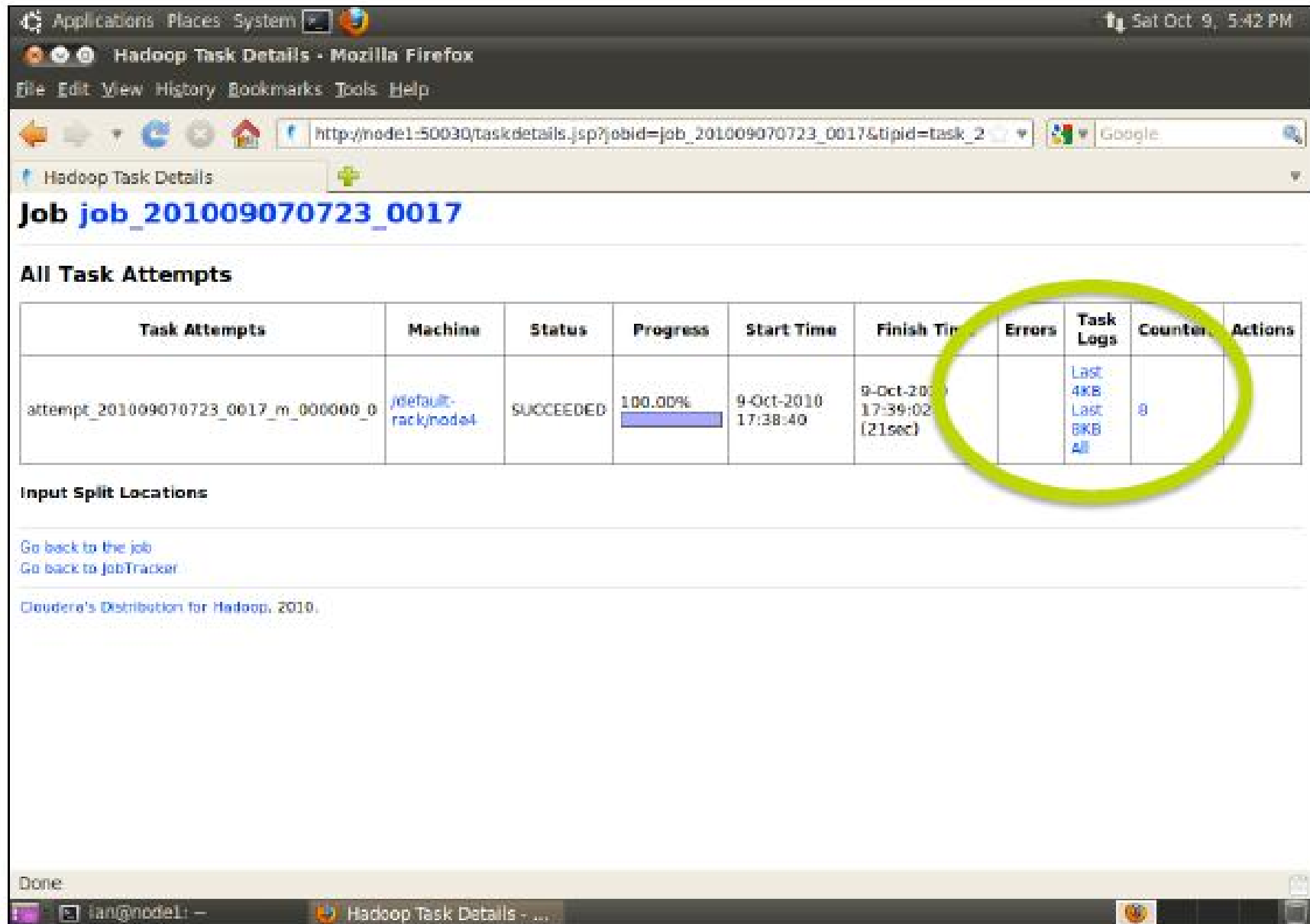
- Tried-and-true debugging technique: write to `stdout` or `stderr`
- If running in `LocalJobRunner` mode, you will see the results of `System.err.println()`
- If running on a cluster, that output will not appear on your console
  - Output is visible via Hadoop's Web UI

## Aside: The Hadoop Web UI

---

- **All Hadoop daemons contain a Web server**
  - Exposes information on a well-known port
- **Most important for developers is the JobTracker Web UI**
  - `http://<job_tracker_address>:50030/`
  - `http://localhost:50030/` if running in pseudo-distributed mode
- **Also useful: the NameNode Web UI**
  - `http://<name_node_address>:50070/`

# Aside: The Hadoop Web UI (cont'd)



The screenshot shows a web browser window displaying the Hadoop Task Details page for job\_201009070723\_0017. The page title is "Job job\_201009070723\_0017". Below the title, there is a section titled "All Task Attempts" which contains a table with the following data:

Task Attempts	Machine	Status	Progress	Start Time	Finish Time	Errors	Task Logs	Counter	Actions
attempt_201009070723_0017_m_000000_0	/default-rack/node4	SUCCEEDED	100.00% <div><div></div></div>	9-Oct-2010 17:38:40	9-Oct-2010 17:39:02 (21sec)		<a href="#">Last 4KB</a> <a href="#">Last 8KB</a> <a href="#">All</a>	0	

Below the table, there is a section titled "Input Split Locations" with links to "Go back to the job" and "Go back to jobTracker". At the bottom, there is a footer that reads "Cloudera's Distribution for Hadoop, 2010."

A yellow circle highlights the "Errors", "Task Logs", and "Counter" columns of the table.

## Logging: Better Than Printing

---

- **println statements rapidly become awkward**
  - Turning them on and off in your code is tedious, and leads to errors
- **Logging provides much finer-grained control over:**
  - What gets logged
  - When something gets logged
  - How something is logged





## Logging With `log4j`

---

- Hadoop uses `log4j` to generate all its log files
- Your Mappers and Reducers can also use `log4j`
  - All the initialization is handled for you by Hadoop
- Add the `log4j.jar-<version>` file from your CDH distribution to your classpath when you reference the `log4j` classes

```
import org.apache.log4j.Level;
import org.apache.log4j.Logger;

class FooMapper implements Mapper {
    private static final Logger LOGGER =
        Logger.getLogger (FooMapper.class.getName());
    ...
}
```

## Logging With `log4j` (cont'd)

---

- Simply send strings to loggers tagged with severity levels:

```
LOGGER.trace("message");  
LOGGER.debug("message");  
LOGGER.info("message");  
LOGGER.warn("message");  
LOGGER.error("message");
```

- Beware expensive operations like concatenation
  - To avoid performance penalty, make it conditional like this:

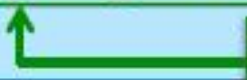
```
if (LOGGER.isDebugEnabled()) {  
    LOGGER.debug("Account info:" + acct.getReport());  
}
```

## log4j Configuration

---

- Node-wide configuration for log4j is stored in `/etc/hadoop/conf/log4j.properties`
- Override settings for your application in your own `log4j.properties`
  - Can change global log settings with `hadoop.root.log` property
  - Can override log level on a per-class basis, e.g.

```
log4j.logger.org.apache.hadoop.mapred.JobTracker=WARN  
log4j.logger.com.mycompany.myproject.FooMapper=DEBUG
```



Full class name

- Or set the level programmatically:

```
LOGGER.setLevel(Level.WARN);
```



## Setting Logging Levels for a Job

---

- You can tell Hadoop to set logging levels for a job using configuration properties

- mapred.map.child.log.level
  - mapred.reduce.child.log.level

- Examples

- Set the logging level to DEBUG for the Mapper

```
$ hadoop jar myjob.jar MyDriver \  
-Dmapred.map.child.log.level=DEBUG indir outdir
```

- Set the logging level to WARN for the Reducer

```
$ hadoop jar myjob.jar MyDriver \  
-Dmapred.reduce.child.log.level=WARN indir outdir
```



## Where Are Log Files Stored?

---

- **Log files are stored on the machine where the task attempt ran**
  - Location is configurable
  - By default:  
`/var/log/hadoop-0.20-mapreduce/  
userlogs/${task.id}/syslog`
- **You will often not have ssh access to a node to view its logs**
  - Much easier to use the JobTracker Web UI
    - Automatically retrieves and displays the log files for you

## Restricting Log Output

---

- **If you suspect the input data of being faulty, you may be tempted to log the (key, value) pairs your Mapper receives**
  - Reasonable for small amounts of input data
  - Caution! If your job runs across 500GB of input data, you could be writing up to 500GB of log files!
  - Remember to think at scale...
- **Instead, wrap vulnerable sections of code in `try { . . . }` blocks**
  - Write logs in the `catch { . . . }` block
    - This way only critical data is logged



## Aside: Throwing Exceptions

---

- You could throw exceptions if a particular condition is met

- For example, if illegal data is found

```
throw new RuntimeException("Your message here");
```



- Usually not a good idea

- Exception causes the task to fail
  - If a task fails four times, the entire job will fail

# Reusing Objects



## Reuse of Objects is Good Practice (1)

---

- **It is generally good practice to reuse objects**
  - Instead of creating many new objects
- **Example: Our original WordCount Mapper code**

```
public class WordMapper extends Mapper<LongWritable, Text, Text, IntWritable>
{
    @Override
    public void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException {

        String word = value.toString();

        for (String w : word.split(" "))
        {
            if (w.length() > 0)
            {
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}
```

Each time the `map()` method is called, we create a new `Text` object and a new `IntWritable` object.

## Reuse of Objects is Good Practice (2)

---

- Instead, this is better practice:

```
public class WordMapper extends Mapper<LongWritable, Text, Text, IntWritable>
{
    private final static IntWritable one = new IntWritable(1);
    private Text wordObject = new Text();
```

```
@Override
```

```
public
```

Create objects for the key and value outside of your `map()` method

```
    String line = value.toString();
```

```
    for (String word : line.split("\\W+")) {
```

```
        if (word.length() > 0) {
```

```
            wordObject.set(word);
```

```
            context.write(wordObject, one);
```

```
        }
```

```
    }
```

```
}
```

```
}
```

## Reuse of Objects is Good Practice (3)

---

- Instead, this is better practice:

```
public class WordMapper extends Mapper<LongWritable, Text, Text, IntWritable>
{
    private final static IntWritable one = new IntWritable(1);
    private Text wordObject = new Text();
```

```
@Override
protected
```

Within the `map()` method, populate the objects and write them out. Hadoop will take care of serializing the data so it is perfectly safe to re-use the objects.

```
    if (word.length() > 0) {
        wordObject.set(word);
        context.write(wordObject, one);
    }
}
}
```

## Object Reuse: Caution!

---

- **Hadoop re-uses objects all the time**
- **For example, each time the Reducer is passed a new value, the same object is reused**
- **This can cause subtle bugs in your code**
  - For example, if you build a list of value objects in the Reducer, each element of the list will point to the same underlying object
    - Unless you do a deep copy



## Map-Only MapReduce Jobs

---

- **There are many types of job where only a Mapper is needed**
- **Examples:**
  - Image processing
  - File format conversion
  - Input data sampling
  - ETL

## Creating Map-Only Jobs

---

- To create a Map-only job, set the number of Reducers to 0 in your Driver code

```
job.setNumReduceTasks(0);
```

- Call the `Job.setOutputKeyClass` and `Job.setOutputValueClass` methods to specify the output types
  - *Not* the `Job.setMapOutputKeyClass` and `Job.setMapOutputValueClass` methods
- Anything written using the `Context.write` method in the Mapper will be written to HDFS
  - Rather than written as intermediate data
  - One file per Mapper will be written