

# Data Intensive Computing – Handout 8

## Hadoop

Hadoop 1.2.1 is installed in `/HADOOP` directory. The JobTracker web interface is available at `http://dlrc:50030`, the NameNode web interface is available at `http://dlrc:50070`.

To conveniently access the web interfaces remotely, you can use `ssh -D localhost:2020 -p 11422 ufallab.ms.mff.cuni.cz` to open SOCKS5 proxy server forwarding requests to the remote site, and use it as a proxy server in a browser, for example as `chromium --proxy-server=socks://localhost:2020`.

## Java Interface

### Running Java Jobs

Java jobs are executing using `hadoop jar filename.jar` command. The additional arguments are passed to the `filename.jar` program. Usually the arguments are `[options] input_path output_path` with the following options:

- `-D mapred.map.tasks=N` set specified number of map tasks
- `-D mapred.reduce.tasks=N` set specified number of reduce tasks, can be 0
- `-D mapreduce.inputformat.class=cls` set specified input format, can be one of
  - `org.apache.hadoop.mapreduce.lib.input.TextInputFormat`
  - `org.apache.hadoop.mapreduce.lib.input.SequenceFileInputFormat`
  - `org.apache.hadoop.mapreduce.lib.input.KeyValueTextInputFormat`
- `-D mapreduce.outputformat.class=cls` set specified input format, can be one of
  - `org.apache.hadoop.mapreduce.lib.output.TextOutputFormat`
  - `org.apache.hadoop.mapreduce.lib.output.SequenceFileOutputFormat`

### Running Java Jobs Locally

Java jobs can be executed locally by using `hadoop-local` command. In that case, mappers and reducers are executed sequentially, directly in the running process.

## Manupulating HDFS

The HDFS content can be manipulated using `hadoop fs` command, which prints nice usage when executed without arguments. Some simple commands:

- `hadoop fs -ls /`
- `hadoop fs -rmr /users/straka/test /`
- `hadoop fs -get /users/straka/test . /`

## Java Examples

All mentioned examples are available in `/home/straka/java/examples`.

# Grep AB

## GrepAB.java

```
import java.io.IOException;

import org.apache.hadoop.conf.*;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.*;
import org.apache.hadoop.mapreduce.lib.output.*;
import org.apache.hadoop.util.*;

public class GrepAB extends Configured implements Tool {
    // Mapper
    public static class TheMapper extends Mapper<Text, Text, Text, Text>{
        public void map(Text key, Text value, Context context) throws IOException,
            InterruptedException {
            if (key.toString().startsWith("AB"))
                context.write(key, value);
        }
    }

    // Job configuration
    public int run(String[] args) throws Exception {
        if (args.length < 2) {
            System.err.printf("Usage: %s.jar -in-path -out-path\n",
                this.getClass().getName());
            return 1;
        }

        Job job = new Job(getConf(), this.getClass().getName());
        job.setJarByClass(this.getClass());
        job.setMapperClass(TheMapper.class);
        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(Text.class);

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

        return job.waitForCompletion(true) ? 0 : 1;
    }

    // Main method
    public static void main(String[] args) throws Exception {
        int res = ToolRunner.run(new GrepAB(), args);
        System.exit(res);
    }
}
```

## GrepAB.sh

```
hadoop fs -rmr /users/$USER/GrepAB
hadoop jar GrepAB.jar \
    -D mapreduce.inputformat.class=org.apache.hadoop.mapreduce.lib.input.SequenceFileInputFormat \
    -D mapred.reduce.tasks=0 \
    /data/wiki/cs-medium /users/$USER/GrepAB
```

## Compiling Examples

Examples can be compiled using either:

- supplied Makefile: just run `make` in the source directory and a `GrepAB.jar` is created. Use `make clean` to remove the compiled classes and java archive.
- ant: create a directory and copy `/home/straka/java/ant/build.xml` to it. Also create `src` subdirectory and copy `GrepAB.java` to it. Then run `ant` in the original directory to create `GrepAB.jar` in the `bin` directory.

When using the `build.xml` on different class, update `target` property in `build.xml`.

When compiling somewhere else than on `dlrc`, copy `/HADOOP/hadoop-core-1.2.1.jar` and update `hadoop-core.jar` property in `build.xml`.

- Eclipse: copy `build.xml` and `hadoop-core-1.2.1.jar` from `/home/straka/java/eclipse` to a local directory. Copy `GrepAB.java` to the `src` subdirectory. Then run Eclipse, select File – New – Project – Java Project from Existing Ant Buildfile, name the project, choose the `build.xml`, and check “Link to the buildfile in the file system”. The new project will use the original directory as both source and output directory.

## Hadoop Documentation

We use Hadoop 1.2.1. The API documentation can be found at <http://hadoop.apache.org/docs/r1.2.1/api/>.

## Word Count example

Simple example using a reducer and combiner.

WordCount.java

```
import java.io.IOException;
import java.util.StringTokenizer;

import org.apache.hadoop.conf.*;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.*;
import org.apache.hadoop.mapreduce.lib.output.*;
import org.apache.hadoop.util.*;

public class WordCount extends Configured implements Tool {
    // Mapper
    public static class TheMapper extends Mapper<Text, Text, Text, IntWritable>{
        private Text word = new Text();
        private IntWritable one = new IntWritable(1);

        public void map(Text key, Text value, Context context) throws IOException,
            InterruptedException {
            StringTokenizer st = new StringTokenizer(value.toString(), "\t\n\r\f,.;?![]()'\`");
            while (st.hasMoreTokens()) {
                word.set(st.nextToken());
                context.write(word, one);
            }
        }
    }
}
```

```

// Reducer
public static class TheReducer extends Reducer<Text, IntWritable, Text,
    IntWritable> {
    private IntWritable sumWritable = new IntWritable();

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

// Job configuration
public int run(String[] args) throws Exception {
    if (args.length < 2) {
        System.err.printf("Usage: %s.jar -in-path -out-path\n",
            this.getClass().getName());
        return 1;
    }

    Job job = new Job(getConf(), this.getClass().getName());
    job.setJarByClass(this.getClass());
    job.setMapperClass(TheMapper.class);
    job.setMapOutputKeyClass(Text.class);
    job.setMapOutputValueClass(IntWritable.class);
    job.setCombinerClass(TheReducer.class);
    job.setReducerClass(TheReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);

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

    return job.waitForCompletion(true) ? 0 : 1;
}

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

```

#### WordCount.sh

```

hadoop fs -rmr /users/$USER/WordCount
hadoop jar WordCount.jar \
    -D mapreduce.inputformat.class=org.apache.hadoop.mapreduce.lib.input.SequenceFileInputFormat
    /data/wiki/cs-medium /users/$USER/WordCount

```

## Partitioner

Using Java Hadoop, partitioner can be specified easily:

#### Sort.java

```

import java.io.IOException;

```

```

import org.apache.hadoop.conf.*;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.*;
import org.apache.hadoop.mapreduce.lib.output.*;
import org.apache.hadoop.util.*;

public class Sorting extends Configured implements Tool {
    public static class TheMapper extends Mapper<Text, Text, IntWritable, Text> {
        public void map(Text key, Text value, Context context) throws IOException,
            InterruptedException {
            // INSERT CODE HERE
        }
    }

    public static class ThePartitioner extends Partitioner<IntWritable, Text> {
        public int getPartition(IntWritable key, Text value, int numPartitions) {
            // INSERT CODE HERE
        }
    }

    public static class TheReducer extends Reducer<IntWritable, Text,
        IntWritable, Text> {
        public void reduce(IntWritable key, Iterable<Text> values, Context context)
            throws IOException, InterruptedException {
            // INSERT CODE HERE
        }
    }

    public int run(String[] args) throws Exception {
        if (args.length < 2) {
            System.err.printf("Usage: %s.jar -in-path -out-path",
                this.getClass().getName());
            return 1;
        }

        Job job = new Job(getConf(), this.getClass().getName());

        job.setJarByClass(this.getClass());
        job.setMapperClass(TheMapper.class);
        job.setMapOutputKeyClass(IntWritable.class);
        job.setMapOutputValueClass(Text.class);
        job.setPartitionerClass(ThePartitioner.class);
        job.setReducerClass(TheReducer.class);
        job.setOutputKeyClass(IntWritable.class);
        job.setOutputValueClass(Text.class);

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

        return job.waitForCompletion(true) ? 0 : 1;
    }

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

```

# Tasks

Solve the following tasks. Solution for each task is a Java Hadoop source processing the Wikipedia source data and producing required results.

## Simple Tokenizer

You can use `StringTokenizer(sentence, " \t\n\r\f,.;?![]()'\\"")` as a simple tokenizer:

```
StringTokenizer st = new StringTokenizer(data, " \t\n\r\f,.;?![]()'\\"");
while (st.hasMoreTokens())
    System.out.println(st.nextToken());
```

## Wikipedia Data

The Wikipedia Data available from `/dlrc_share/data/wiki/` are available also in HDFS:

- `/data/wiki/cs`: Czech Wikipedia data (Sep 2009), 195MB, 124k articles
- `/data/wiki/en`: English Wikipedia data (Sep 2009), 4.9GB, 2.9M articles

All data is stored in a record-compressed sequence files, with article names as keys and article texts as values, in UTF-8 encoding.

<i>Task</i>	<i>Points</i>	<i>Description</i>
<code>java_unique_words</code>	2	Create a list of unique words used in the articles using Java Hadoop. Convert them to lowercase to ignore case.
<code>java_anagrams</code>	2	Two words are anagrams if one is a letter permutation of the other (ignoring case). For a given input, find all anagram classes that contain at least <i>A</i> words. Output each anagram class on a separate line.
<code>java_sort</code>	3	You are given data consisting of (31-bit integer, string data) pairs. These are available in plain text format: <ul style="list-style-type: none"><li>• <code>/data/numbers-txt/numbers-small</code>: 3MB</li><li>• <code>/data/numbers-txt/numbers-medium</code>: 184MB</li><li>• <code>/data/numbers-txt/numbers-large</code>: 916MB</li></ul> You can assume that the integers are uniformly distributed. Your task is to sort these data, comparing the key numerically and not lexicographically. The lines in the output must be the same as in the input, only in different order. Your solution should work for TBs of data. For that reason, you must use multiple reducers. If your job is executed using <i>r</i> reducers, the output consists of <i>r</i> files, which when concatenated would produce sorted (key, value) pairs. In other words, each of the output files contains sorted (integer, data) pairs and all keys in one file are either smaller or larger than in other file. Your solution should work for any number of reducer specified. You can use the <code>Sort.java</code> example template.

<i>Task</i>	<i>Points</i>	<i>Description</i>
java_nonuniform_sort	4	<p>Improve the <code>java_sort</code> to handle nonuniform data. You can use the following exponentially distributed data:</p> <ul style="list-style-type: none"> <li>• <code>/data/numbers-txt/nonuniform-small</code>: 3MB</li> <li>• <code>/data/numbers-txt/nonuniform-medium</code>: 160MB</li> <li>• <code>/data/numbers-txt/nonuniform-large</code>: 797MB</li> </ul> <p>Assume we want to produce <math>r</math> output files. One of the solutions is to perform two Hadoop jobs:</p> <ul style="list-style-type: none"> <li>• Go through the data and sample only a small fraction of the keys. As there are not so many of them, they can fit in one reducer.</li> <li>• Find best <math>r - 1</math> separators using the sampled data.</li> <li>• Run the second pass using the computed separators.</li> </ul>
java_inverted_index	2	<p>Compute inverted index in Java Hadoop – for every lowercased word from the articles, compute <i>(article name, ascending positions of occurrences as word indices)</i> pairs.</p> <p>The output should be a file with one word on a line in the following format:</p> <pre>word \t articleName \t spaceSeparatedOccurrences...</pre> <p>You will get 2 additional points if the articles will be numbered using consecutive integers. In that case, the output is ascending <i>(article id, ascending positions of occurrences as word indices)</i> pairs, together with a file containing list of articles representing this mapping (the article on line <math>i</math> is the article with id <math>i</math>).</p>
java_no_references	3	<p>An article <math>A</math> is said to reference article <math>B</math>, if it contains <math>B</math> as a token (ignoring case).</p> <p>Run a Java Hadoop job which counts for each article how many references there exists for the given article (summing all references in a single article).</p> <p>You will get one extra point if the result is sorted by the number of references (you are allowed to use 1 reducer in the sorting phase).</p>
java_wordsim_index	4	<p>In order to implement word similarity search, compute for each form with at least three occurrences all <i>contexts</i> in which it occurs, including their number of occurrences. List the contexts in ascending order.</p> <p>Given <math>N</math> (either 1, 2, 3 or 4), the <i>context</i> of a form occurrence is <math>N</math> forms preceding this occurrence and <math>N</math> forms following this occurrence (ignore sentence boundaries, use empty words when article boundaries are reached).</p> <p>The output should be a file with one form on a line in the following format:</p> <pre>form \t context \t counts...</pre>

<i>Task</i>	<i>Points</i>	<i>Description</i>																
java_wordsim_find	4	<p>Let <math>S</math> be given natural number. Using the index created in <code>java_wordsim_index</code>, find for each form <math>S</math> most similar forms. The similarity of two forms is computed using <i>cosine similarity</i> as <math>\frac{C_A \cdot C_B}{ C_A  \cdot  C_B }</math>, where <math>C_F</math> is a vector of occurrences of form <math>F</math> contexts.</p> <p>The output should be a file with one form on a line in the following format:  <code>form \t most similar form \t cosine similarity...</code></p>																
java_kmeans	6	<p>Implement K-means clustering algorithm as described on <a href="http://en.wikipedia.org/wiki/K-means_clustering#Standard_algorithm">http://en.wikipedia.org/wiki/K-means_clustering#Standard_algorithm</a>.</p> <p>The user specifies number of iterations and the program run specified number of K-means clustering algorithm iterations. You can use the following training data. Each line contains space separated coordinates of one points. The coordinates in one input naturally have the same dimension.</p> <table border="1"> <thead> <tr> <th><i>HDFS path</i></th> <th><i>Points</i></th> <th><i>Dimension</i></th> <th><i>Clusters</i></th> </tr> </thead> <tbody> <tr> <td><code>/data/points-txt/small</code></td> <td>10000</td> <td>50</td> <td>50</td> </tr> <tr> <td><code>/data/points-txt/medium</code></td> <td>100000</td> <td>100</td> <td>100</td> </tr> <tr> <td><code>/data/points-txt/large</code></td> <td>500000</td> <td>200</td> <td>200</td> </tr> </tbody> </table> <p>You will get 2 additional points if the algorithm stops when none of the centroid positions change more than given <math>\epsilon</math>.</p>	<i>HDFS path</i>	<i>Points</i>	<i>Dimension</i>	<i>Clusters</i>	<code>/data/points-txt/small</code>	10000	50	50	<code>/data/points-txt/medium</code>	100000	100	100	<code>/data/points-txt/large</code>	500000	200	200
<i>HDFS path</i>	<i>Points</i>	<i>Dimension</i>	<i>Clusters</i>															
<code>/data/points-txt/small</code>	10000	50	50															
<code>/data/points-txt/medium</code>	100000	100	100															
<code>/data/points-txt/large</code>	500000	200	200															