# Big Data 4

## Week 3: MapReduce

Dr. Nikos Ntarmos

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

School of Computing Science
University of Glasgow

15 October, 2013

# Outline

# MapReduce basics :: Introduction

# Motivation

- Need to process **petabytes** of data ...

- ... on your **multi-thousand commodity PC** cluster ...

- ... and make it **easy**!

$\Rightarrow$ Enter **Map/Reduce**

---

[1] J. Dean and S. Ghemawat. ``MapReduce: Simplified Data Processing on Large Clusters''. In *Proc. OSDI '04*.

# What is MapReduce?

1. MapReduce -- The programming model
   - Functional-like programming $\rightarrow$ Map/reduce functions

2. MapReduce -- The execution framework
   - Input splitting, intermediate result sort/shuffling, task scheduling, ...

3. MapReduce -- The software implementation
   - Google MapReduce, Hadoop, Storm/Kafka, AcB/Cell, ...

# MapReduce premises

- Scale *out*, not *up*.
- Design for fault tolerance (assume failures are the norm).
- Move processing to the data.
- Optimize for **high throughput batch jobs** -- i.e., sequential reads, no random accesses.
- Hide system-related stuff from the app developers $\Rightarrow$ provide automated distribution, parallelization, scheduling, . . .
- Provide mechanisms for progress report and monitoring.
- ``Always'' operates on key-value pairs

# Hadoop Timeline

**2002:** Apache Nutch (open-source web search engine).

**2003:** Google File System (GFS) paper appears in SOSP.

**2004:** Nutch Distributed File System (NDFS) airs.

Google MapReduce paper appears in OSDI.

**2005:** MapReduce is implemented for NDFS.

**2006: Hadoop is born!**

**2008:** Hadoop gets top-level Apache project status

... and breaks the record for fastest TB sort.

**2009:** Hadoop breaks two more records ($<1'$ to sort 500GB, $<3'$ to sort 100 TB).

**2012:** YARN is born!

**2013:** . . .

# Example: Word Count

- ``*Million lyrics dataset*´´

```
Hear the rime of the Ancient Mariner, see his eye as he stops one of three.
Mesmerises one of the wedding guests. Stay here and listen to the nightmares
    of the Sea.
And the music plays on, as the bride passes by. Caught by his spell and the
    Mariner tells his tale.

Driven south to the land of the snow and ice, to a place where nobody's been.
Through the snow fog flies on the albatross. Hailed in God's name, hoping
    good luck it brings.

And the ship sails on, back to the North. Through the fog and ice and the
    albatross follows on.

The mariner kills the bird of good omen. His shipmates cry against what he's
    done,
but when the fog clears, they justify him and make themselves a part of the
    crime.
[...]
```

Q: Compute all unique words in the stream and their number of appearances (*Word Count*).

# Example: Word Count

*Pfft... We can do that with simple shell scripts...*

- Using shell programming, take #1:

```
cat input.txt | tr '[:space:]' '\n' | sort | uniq -c | \
    awk '{print $2" "$1}'
```

Now what? → **Must sort entire input and go over it twice!**

- Using shell programming, take #2:

```
cat input.txt | awk '\
    {for (i = 1; i <= NF; i++) freqs[$i]++;} END \
    {for (word in freqs) print word" "freqs[word];}'
```

What's wrong with this? → **Out of RAM as freqs grows huge! Plus serial execution...**

**Let's combine these two and go all-out distributed!**

MapReduce basics :: The MapReduce paradigm

## The MapReduce paradigm

$$\textbf{map} : (k_1, v_1) \rightarrow [(k_2, v_2)]$$
$$\textbf{reduce} : (k_2, [v_2]) \rightarrow [(k_3, v_3)]$$

```
void Map(int lineNo, String line) {
   foreach (Word w in line)
      emit(tolower(w), 1);
}
```

```
void Reduce(Word w, int[] counts) {
   int sum = 0;
   foreach (int i in counts)
      sum += i;
   emit(w, sum);
}
```

``*And the ship sails on, back to the North, through the fog and ice and the albatross follows on.*''

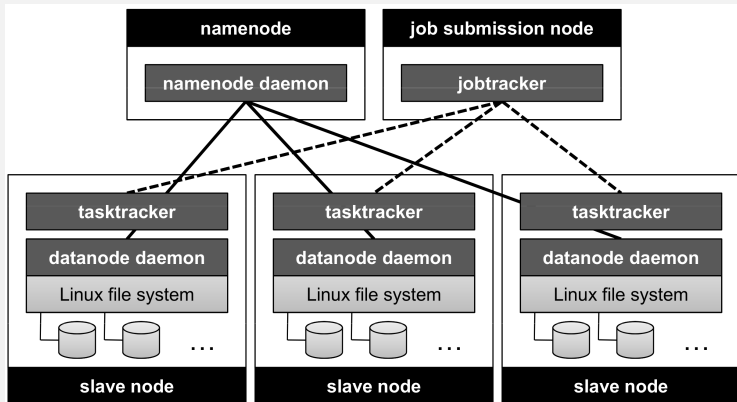| | | | | | | | |
|---|---|---|---|---|---|---|---|
| and | 1 | | albatross | 1 | | albatross | 1 |
| the | 1 | | and | 1,1,1 | | and | 3 |
| ship | 1 | | back | 1 | | back | 1 |
| sails | 1 | | fog | 1 | | fog | 1 |
| on | 1 | | follows | 1 | | follows | 1 |
| back | 1 | $\Rightarrow$ | ice | 1 | $\Rightarrow$ | ice | 1 |
| to | 1 | | north | 1 | | north | 1 |
| the | 1 | | on | 1,1,1 | | on | 2 |
| north | 1 | | sails | 1 | | sails | 1 |
| through | 1 | | ship | 1 | | ship | 1 |
| the | 1 | | the | 1,1,1,1 | | the | 4 |
| fog | 1 | | through | 1 | | through | 1 |
| ... | | | to | 1 | | to | 1 |

# MapReduce processing steps



- Input parsing & splitting → *InputFormat*
- Map function execution → *Mapper, Combiner*
- Intermediate result partitioning, sorting, and storing → *Partitioner, Merge*
- Intermediate data transfer to reducers → *Shuffle*
- Reduce function execution → *Reducer*
- Final output storage → *OutputFormat*

# MapReduce infrastructure components



- Jobtracker-tasktrackers communication over **heartbeats**
- **Predefined slots** for mappers & reducers on tasktrackers

Anatomy of a MapReduce job :: Execution stages

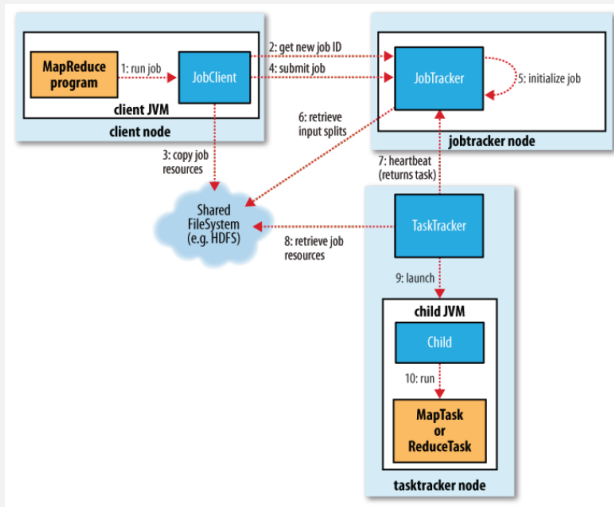# Execution stages

- Execution stages:
    1. Job submission
    2. Job initialization
    3. Map tasks scheduling
    4. Execution of mappers
    5. Execution of combiners
    6. Intermediate output sorting
    7. Reduce task scheduling
    8. Execution of reducers
    9. Job completion/cleanup

# Execution stages



---

[1]

[1]Source: *Hadoop: The definitive guide, 3rd Edition.*

Anatomy of a MapReduce job :: Job submission

# Job submission

1. Contact the jobtracker and get a new job ID
2. Configure all client-writable job settings
3. Define the class names implementing the mappers, reducers, combiners, ...
4. Define the types of keys, values, intermediate results, ...
5. Compute input splits
6. Upload class files (jar), conf files, split keys, etc. to HDFS
7. Let the jobtracker know everything is ready (i.e., ``submit'' the job)

# Input handling

**Goal: High throughput, parallel, fault tolerant data crunching**

$\Rightarrow$ Allow for as much parallelism as possible

$\Rightarrow$ Minimize the need for synchronization across mappers

$\Rightarrow$ Lose only a small part of the computation when a node goes down

**Solution: Split up the input; let each mapper get its own part**

# Input splitting

- Default strategy: split files into **as many equally-sized parts** as there are mappers
- Splits should (by default) be at most one HDFS block large **(lower limit on # splits)**
- Splits are precomputed & stored along with the data files on HDFS
- Splits are **data agnostic**

**What happens if a word is split in half?**

**Let the InputFormat decide!**

## Input splitting

**abstract class org.apache.hadoop.mapreduce.InputFormat**
(subclasses under `org.apache.hadoop.mapreduce.lib.input`)

$\rightarrow$ `FileInputFormat`: Read data from files on HDFS; **byte-oriented operations**

$\rightarrow$ `TextInputFormat`: Treat input as text files; line-oriented operation

$\rightarrow$ `KeyValueTextInputFormat`: Treat input as files storing key-value pairs, one per line; key-values separated by a special character (default: tab)

$\rightarrow$ `SequenceFileInputFormat`: Treat input as binary files storing key-value pairs (i.e., `org.apache.hadoop.io.SequenceFile`'s)

$\rightarrow$ `SequenceFileAsTextInputFormat`: Converts keys/values to Strings by means of their `toString()` method

$\rightarrow$ `SequenceFileAsBinaryInputFormat`: Reads keys/values in their raw format

$\rightarrow$ `CombineFileInputFormat`: Process multiple small files as a single large file

$\rightarrow$ `CombineTextInputFormat`: Combine multiple text files

$\rightarrow$ `CombineSequenceFileInputFormat`: Combine multiple `SequenceFile`'s

$\rightarrow$ `NLineInputFormat`: Treat input as text files, create a split every N lines

$\rightarrow$ `DBInputFormat`: Read data from a SQL table (over JDBC)

$\rightarrow$ `CompositeInputFormat`: ``Join'' multiple data sources

# Anatomy of a MapReduce job :: Job initialization

# Job initialization

- The client has now submitted the new job to the Jobtracker.
- $\rightarrow$ The job is added to the scheduler queue
- $\rightarrow$ The scheduler picks it up and *initializes it*
  1. Retrieves the precomputed input splits
  2. Creates one *map* task per split
  3. Creates as many *reduce* tasks as requested by the client
  4. Creates a *job setup* task for each tasktracker
     ($\rightarrow$ creates intermediate and final output directories, if needed)
  5. Creates a *job cleanup* task for each tasktracker
     ($\rightarrow$ removes intermediate output files and directories)

Anatomy of a MapReduce job :: Task assignment

# Task assignment

- Next step: assign tasks to tasktrackers
- The jobtracker selects which tasktracker will take on the next task
- **? How does it know which tasktracker is idle?**
- → Through jobtracker-tasktracker **heartbeats**
- **? How does it know which task to run where?**
- → That's a work for the **Scheduler**
  - FIFO / Fair / Capacity Scheduler
  - Data-local, rack-local, ``remote'' tasks
- ... plus **Speculative Execution** ...

Anatomy of a MapReduce job :: Task execution

# Task execution

- Now tasktrackers take the lead:
  1. Create local working dir for the task
  2. Transfer/unpack the job files to the local dir
  3. Launch a new JVM to execute the task code
  4. Monitor the progress of children JVM ...
     ... where progress is:
       - Reading/writing an input/output record
       - Setting the status description
       - Incrementing a counter
       - Explicitly calling the Reporter's `progress()` method
  5. Aggregate progress reports and send them to the jobtracker
  6. JVMs optionally reused -- but **not** across jobs

Anatomy of a MapReduce job :: Map tasks

# Map tasks

- Map tasks execute the map function in a loop
- Separate invocation for each input record ... but this can be overriden
  (through `TaskInputOutputContext.nextKeyValue()`)
- Output is *emitted*, through `Context.write()`
  - Initially stored in memory
  - Periodically written out to disk
    ... or wherever requested ...

# Mapper output handling

- Mappers strive for data locality. What about reducers?
- → Each reducer may need data from **all** mappers!
- → Mappers pre-partition their output to separate files
- Partitions defined by the Partitioner
    - **HashPartitioner**, `BinaryPartitioner`, `KeyFieldBasedPartitioner`, `TotalOrderPartitioner`,...
- Each partition sorted in-memory before being written to disk
- Finally, all partitions merged/sorted into a single file on disk

Anatomy of a MapReduce job :: Combiners

# Combiners

- **Combiner $\approx$ (mini) Reducer**
    - Executed prior to storing map outputs on disk (maybe also later)
    - Works with data already in main memory $\rightarrow$ **fast** data access
    - Pre-processes output data just like a reducer would
    - Speeds up later phases (sorting/shuffling/transferring/reducing)
        - $\Rightarrow$ Lower intermediate results' disk space usage
        - $\Rightarrow$ Lower shuffling phase network usage
        - $\Rightarrow$ ... and more ...

Anatomy of a MapReduce job :: Reduce tasks

# Reduce tasks

1. As map tasks finish, the jobtracker informs tasktrackers to launch reducers

2. Each reducer first fetches relevant map output partitions

   ... in parallel, to increase throughput

3. Partitions are merge/sorted, in rounds

4. The final round output is directly fed to the reducer function...

   ... and its output directly written to HDFS (or wherever instructed by the job's `OutputFormat`)

# Anatomy of a MapReduce job :: Job output

# Job output

**abstract class org.apache.hadoop.mapreduce.OutputFormat**
(subclasses under org.apache.hadoop.mapreduce.lib.output)

→ FileOutputFormat: Write data to files on HDFS

   → TextOutputFormat: Write plain text files

   → MapFileOutputFormat: Write indexed key-value files

   → SequenceFileOutputFormat: Write output in
org.apache.hadoop.io.SequenceFile's

     → SequenceFileAsBinaryOutputFormat: Write output in raw format

→ NullOutputFormat: Send all output to /dev/null

→ LazyOutputFormat: Writes data out in a *lazy* manner

→ MultipleOutputs: Writes data to multiple files, possibly using a different
OutputFormat for each

Extras :: Compression

# Compression

- Input/(intermediate) output data can/should be compressed
  **Why?**
  - More bytes per HDFS block/partition $\Rightarrow$ Lower storage/network overhead
  - Less I/Os across the dataset $\Rightarrow$ Higher throughput
  - Higher CPU overhead...
    ... but it pays out for I/O-intensive workloads
- Multiple codecs already available: ZLib, GZip, BZip2, LZO, Snappy, LZ4, ...

  **... but what about splitting???**
  - Only Bzip2 splittable by default
  - Can be worked around, but not for the faint at heart
    (e.g., see `https://github.com/kevinweil/hadoop-lzo`)

Extras :: Distributed cache

# Distributed cache

- Allows to distribute (read-only) files to tasktrackers
  - Can be used for jars, text, archives, conf files, etc.
  - Compressed files unarchived automatically
  - Jars/libs (can be) added to the task's classpath
- Files are copied only once per job to each tasktracker
- Can be *private* or *public*
- Interesting ``hack'' to share state across tasks!

Extras :: Counters

# Counters

- Hadoop/MapReduce further supplies ``Counters''
- Built-in counters for:
    - HDFS and FileInputFormat/FileOutputFormat profiling
    - Map/reduce tasks' progress
    - Job progress
- Also support for user-defined counters
- All counters automatically aggregated at the jobtracker

Extras :: Streaming/pipes

# Streaming/pipes

- Hadoop is written in Java...
- ...but can ``interface'' with other languages as well
- Pipes
  - $\rightarrow$ Use C++ for mapper/reducer implementations
- Streaming
  - Use any language able of I/O on stdin/stdout!

MapReduce failure handling :: Task failure

# Task failure

- Task failure detected by tasktracker
  - Runtime exception/other user code bug $\to$ error reported to tasktracker by child JVM
  - Child JVM dies $\to$ tasktracker informed by the OS
  - Task is hung $\to$ detected through timeouts in reporting
- The tasktracker informs the jobtracker of the failure
- The jobtracker reschedules the failed task (to a different tasktracker, if possible)
$\to$ Multiple attempts per task
- Jobs with many failed tasks are eventually terminated
- ... or bad input records are marked and skipped

MapReduce failure handling :: Tasktracker failure

# Tasktracker failure

- Tasktracker failure detected by jobtracker
  - Crashed or slow tasktracker $\rightarrow$ missing heartbeats
  - Failing or misconfigured tasktracker $\rightarrow$ multiple task failures
- Tasktracker removed from pool of ``live'' tasktrackers
- Tasks initially assigned to the failed tasktracker are rescheduled
- Failed/blacklisted tasktrackers reconsidered after some time

# MapReduce failure handling :: Jobtracker failure

# Jobtracker failure

### *The death of all!!!*

(... well, not really, but close enough ...)

## ... but wait, YARN is on its way!!!

Putting it all together :: Word Count galore

Word Count v0 :: Built-in mappers/reducers

# Word Count v0

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

Word Count v1 :: User-defined mappers/reducers

# Word Count v1

```java
public class WordCount extends Configured implements Tool {
    static class Map extends org.apache.hadoop.mapreduce.Mapper<LongWritable,
        Text, Text, IntWritable> {
        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();

        public void map(LongWritable key, Text value, Context context) throws
            IOException, InterruptedException {
            String line = value.toString();
            StringTokenizer tokenizer = new StringTokenizer(line);
            while (tokenizer.hasMoreTokens()) {
                word.set(tokenizer.nextToken());
                context.write(word, one);
            }
        }
    }

    public static class Reduce extends Reducer<Text, IntWritable, Text,
        IntWritable> {
        public void reduce(Text key, Iterable<IntWritable> values, Context
            context) throws IOException, InterruptedException {
            int sum = 0;
            for (IntWritable value: values)
                sum += value.get();
            context.write(key, new IntWritable(sum));
        }
    }
```

# Word Count v1 (cont.)

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

Word Count v2 ::

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

# Word Count v2

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

# Word Count v2

```java
public class WordCount extends Configured implements Tool {
    public static class Map extends Mapper<LongWritable, Text, Text, IntWritable>{
        static enum Counters { INPUT_WORDS }
        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();
        private boolean caseSensitive = true;
        private Set<String> patternsToSkip = new HashSet<String>();
        private long numRecords = 0;
        private String inputFile;

        public void setup(Context context) {
            Configuration conf = context.getConfiguration();
            caseSensitive = conf.getBoolean("wordcount.case.sensitive", true);
            inputFile = conf.get("map.input.file");

            if (conf.getBoolean("wordcount.skip.patterns", false)) {
                Path[] patternsFiles = new Path[0];
                try {
                    patternsFiles = DistributedCache.getLocalCacheFiles(conf);
                } catch (IOException ioe) {
                    System.err.println("Caught exception while getting cached files:
                        " + StringUtils.stringifyException(ioe));
                }
                for (Path patternsFile : patternsFiles)
                    parseSkipFile(patternsFile);
            }
        }
```

# Word Count v2 (cont.)

```
27        public void map(LongWritable key, Text value, Context context) throws
              IOException, InterruptedException {
28          String line = caseSensitive ? value.toString() :
                value.toString().toLowerCase();
29
30          for (String pattern : patternsToSkip)
31            line = line.replaceAll(pattern, "");
32
33          StringTokenizer tokenizer = new StringTokenizer(line);
34          while (tokenizer.hasMoreTokens()) {
35            word.set(tokenizer.nextToken());
36            context.write(word, one);
37            context.getCounter(Counters.INPUT_WORDS).increment(1);
38          }
39
40          if ((++numRecords % 100) == 0)
41            context.setStatus("Finished processing " + numRecords + " records "
                + "from the input file: " + inputFile);
42        }
```

## Word Count v2 (cont.)

```
43        private void parseSkipFile(Path patternsFile) {
44            try {
45                BufferedReader fis = new BufferedReader(new
                      FileReader(patternsFile.toString()));
46                String pattern = null;
47                while ((pattern = fis.readLine()) != null)
48                    patternsToSkip.add(pattern);
49                fis.close();
50            } catch (IOException ioe) {
51                System.err.println("Caught exception while parsing the cached file
                      '" + patternsFile + "' : " +
                      StringUtils.stringifyException(ioe));
52            }
53        }
54    }
55
56    public static class Reduce extends Reducer<Text, IntWritable, Text,
          IntWritable> {
57        public void reduce(Text key, Iterable<IntWritable> values, Context
              context) throws IOException, InterruptedException {
58            int sum = 0;
59            for (IntWritable value: values)
60                sum += value.get();
61            context.write(key, new IntWritable(sum));
62        }
63    }
```

## Word Count v2 (cont.)

```
64    public int run(String[] args) throws Exception {
65        Job job = new Job();
66        job.setJobName("WordCount-v3");
67        job.setJarByClass(WordCount.class);
68        job.setMapperClass(Map.class);
69        job.setCombinerClass(Reduce.class);
70        job.setReducerClass(Reduce.class);
71        job.setInputFormatClass(TextInputFormat.class);
72        FileInputFormat.setInputPaths(job, new Path(other_args.get(0)));
73        job.setOutputKeyClass(Text.class);
74        job.setOutputValueClass(IntWritable.class);
75        job.setOutputFormatClass(TextOutputFormat.class);
76        FileOutputFormat.setOutputPath(job, new Path(other_args.get(1)));
77
78        List<String> other_args = new ArrayList<String>();
79        for (int i = 0; i < args.length; ++i) {
80            if ("-skip".equals(args[i])) {
81                DistributedCache.addCacheFile(new Path(args[++i]).toUri(),
                     job.getConfiguration());
82                job.getConfiguration().setBoolean("wordcount.skip.patterns", true);
83            } else
84                other_args.add(args[i]);
85        }
86
87        job.submit();
88        return job.waitForCompletion(true) ? 0 : 1;
89    }
```

# Word Count v2 (cont.)

```
90    public static void main(String[] args) throws Exception {
91        System.exit(ToolRunner.run(new Configuration(), new WordCount(), args));
92    }
93 }
```