

Introducing MapReduce



Riccardo Torlone

Alessio Conte - Luca Menichetti - Fabrizio Rebecca

Summary

- Install Pseudo-Distributed Mode
- Architecture
- Example: Temperature
- Example: Word Count
- Combiner
- Example: Bigram Count (with custom writable)
- How to run a job
- How to run a job on EMR

Summary

- **Install Pseudo-Distributed Mode**
- Architecture
- Example: Temperature
- Example: Word Count
- Combiner
- Example: Bigram Count (with custom writable)
- How to run a job
- How to run a job on EMR

Install Pseudo-Distributed Mode

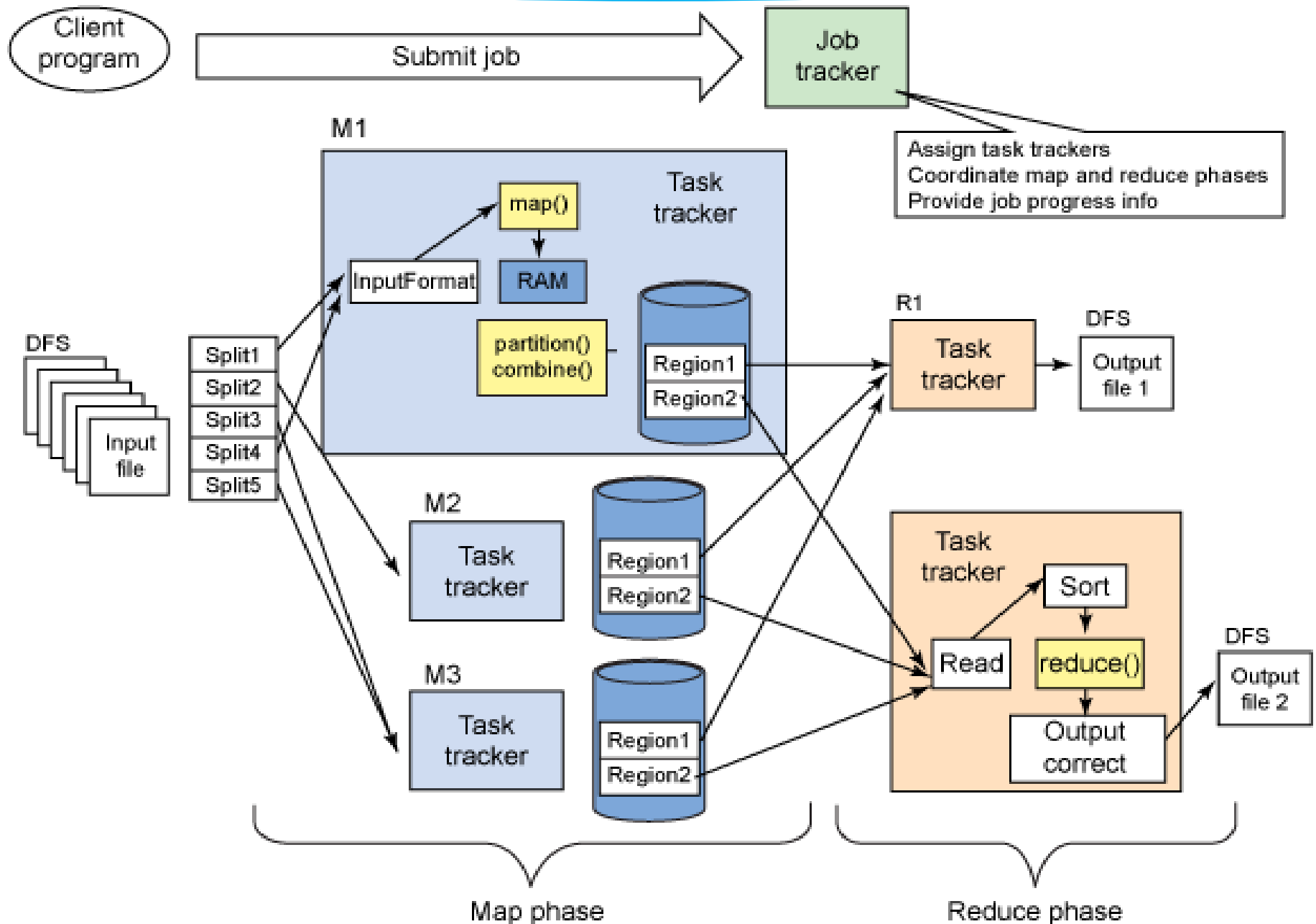
Good luck!

<https://github.com/H4ml3t/Introducing-Hadoop>

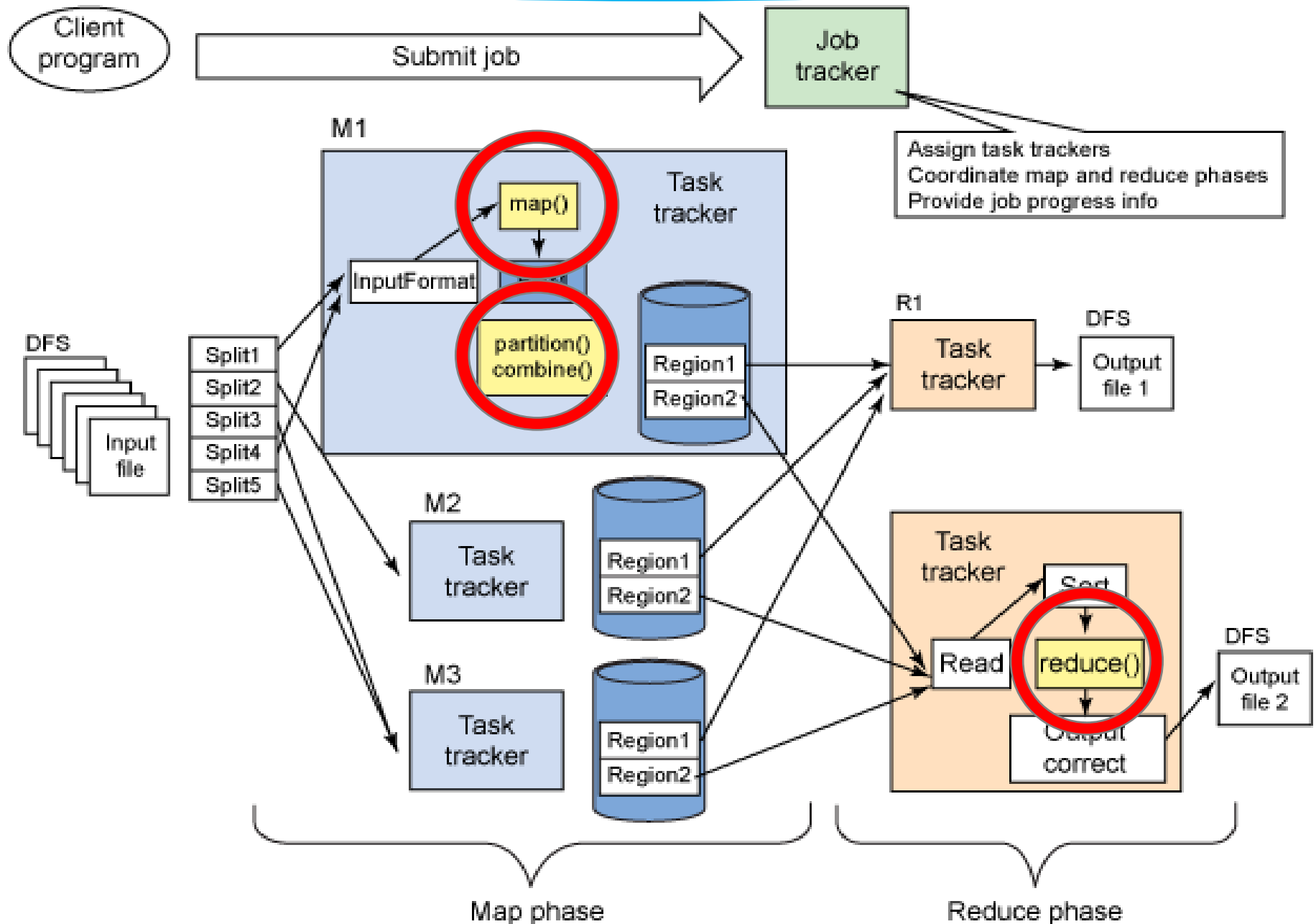
Summary

- Install Pseudo-Distributed Mode
- **Architecture**
- Example: Temperature
- Example: Word Count
- Combiner
- Example: Bigram Count (with custom writable)
- How to run a job
- How to run a job on EMR

Architecture



Architecture



Summary

- Install Pseudo-Distributed Mode
- Architecture
- **Example: Temperature**
- Example: Word Count
- Combiner
- Example: Bigram Count (with custom writable)
- How to run a job
- How to run a job on EMR

Example: Temperature

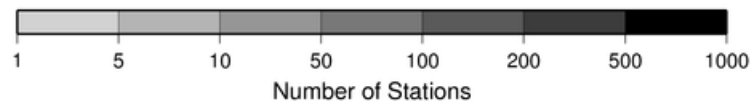
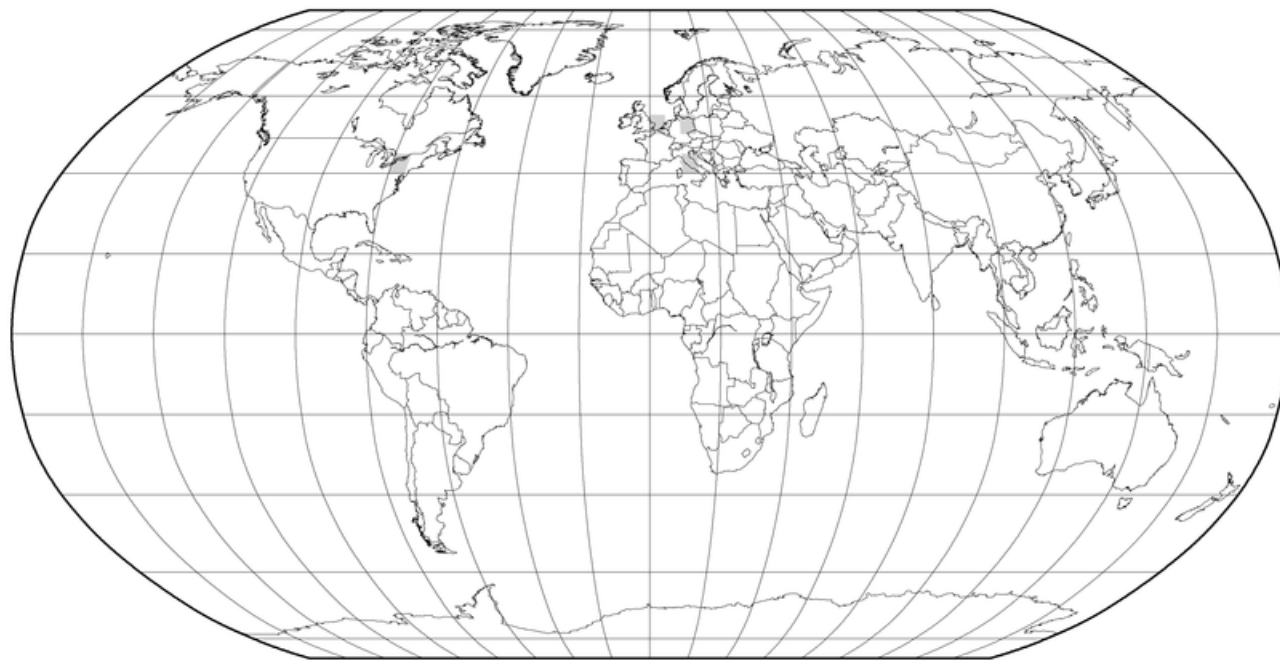
We collected information from <https://www.ncdc.noaa.gov/> (National Climatic Data Center) about the atmosphere.

For each year (starting from 1763), for each month, for each day, for each hours, for each existing meteorological stations in the world we have an entry in a file that specifies: data, time, id, air temperature, pressure, elevation, latitude, longitude...

Example: Temperature

What if we want to know the max temperature in each year?

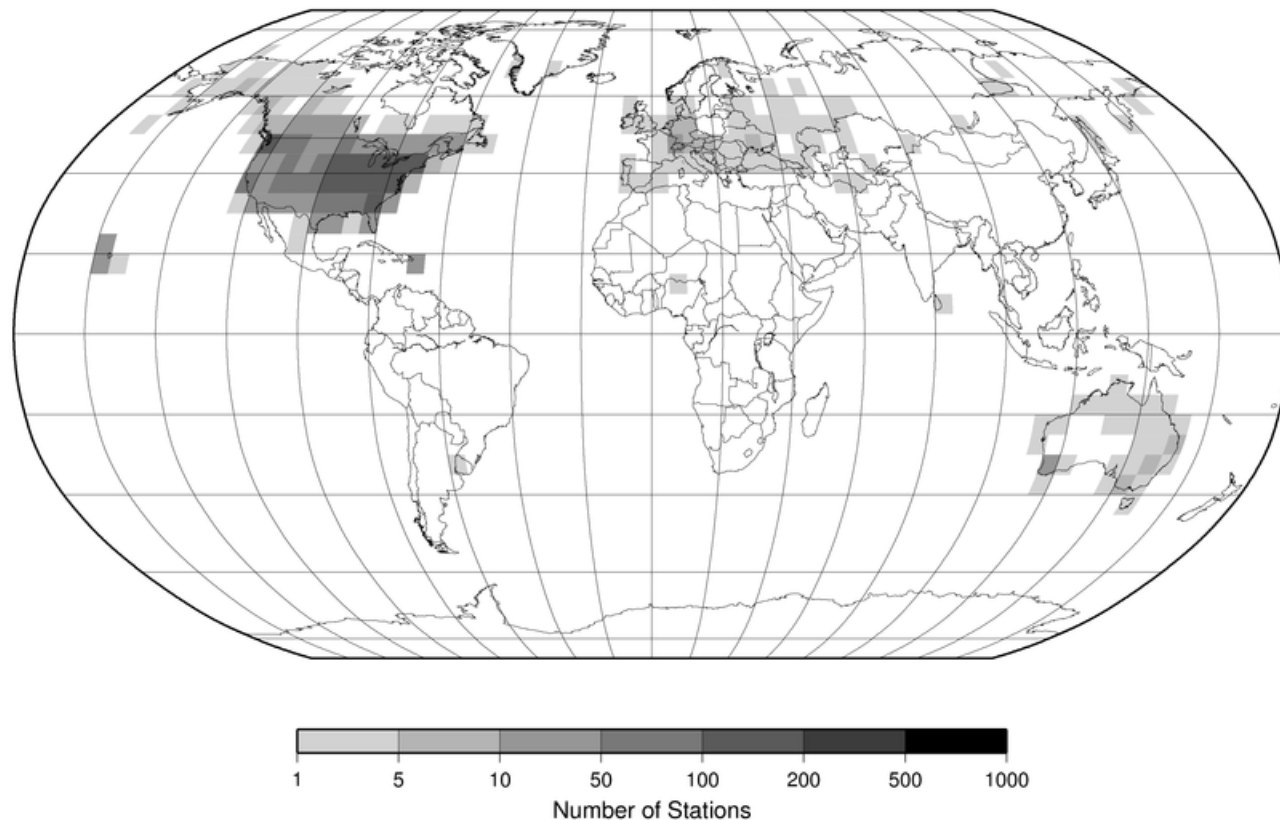
1831–1860 (v3.00–upd–2013020606)



Example: Temperature

What if we want to know the max temperature in each year?

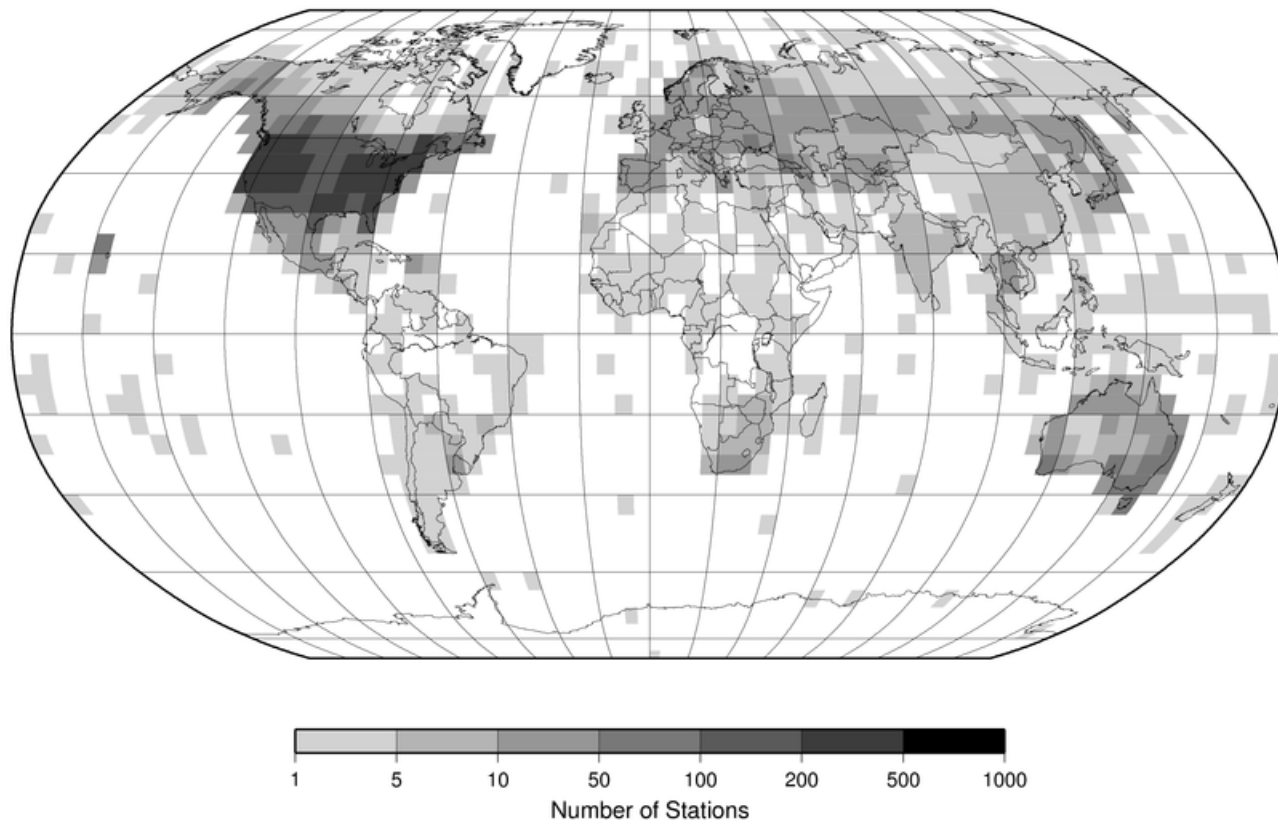
1891–1920 (v3.00–upd–2013020606)



Example: Temperature

What if we want to know the max temperature in each year?

1981–2010 (v3.00–upd–2013020606)



Temperature - data format

File entries example:

...

```
0057332130999991958010103004+51317+028783FM-  
12+017199999V0203201N00721004501CN0100001N9-00211-01391102681  
0057332130999991959010103004+51317+028783FM-  
12+017199999V0203201N00721004501CN0100001N9-+00651-01391102681  
0057332130999991960010103004+51317+028783FM-  
12+017199999V0203201N00721004501CN0100001N9-+00541-01391102681
```

...



Year



Degrees
Celsius x10

Temperature – logical data flow

INPUT

0067011990...
0043011990...
0043011990...
0043012650...
0043012650...



(0, 0067011990...)
(106, 0043011990...)
(212, 0043011990...)
(318, 0043012650...)
(424, 0043012650...)

MAP



(1950, 0)
(1950, 22)
(1950, -11)
(1949, 111)
(1949, 78)

SHUFFLE & SORT



(1950, 0)
(1950, 22)
(1950, -11)
(1949, 111)
(1949, 78)

(1949, [111, 78])
(1950, [0, 22, -11])

REDUCE



(1949, 111)
(1950, 22)

Temperature – Mapper

map (k1, v1) → [(k2, v2)]

```
import org.apache.hadoop.io.*;  
import org.apache.hadoop.mapred.*;
```

```
public class MaxTemperatureMapper extends MapReduceBase  
    implements Mapper<LongWritable, Text, Text, IntWritable> {
```

```
    private static final int MISSING = 9999;
```

```
    public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable>  
        output, Reporter reporter) throws IOException {
```

```
        String line = value.toString();  
        String year = line.substring(15, 19);  
        int airTemperature;
```

map (k1, v1) → [(k2, v2)]

```
        if (line.charAt(87) == '+') {
```

Temperature – Reducer (1)

```
import org.apache.hadoop.io.*;  
import org.apache.hadoop.mapred.*;
```

reduce (k2, [v2]) → [(k3, v3)]

```
public class MaxTemperatureReducer extends MapReduceBase  
    implements Reducer<Text, IntWritable, Text, DoubleWritable> {  
  
    public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text,  
        DoubleWritable> output, Reporter reporter) throws IOException {  
  
        int maxVal = Integer.MIN_VALUE;  
  
        while (values.hasNext()) {  
            maxVal = Math.max(maxVal, values.next().get());  
        }  
  
        output.collect(key, new DoubleWritable(((double)maxVal)/10));  
    }  
}
```


Temperature – Reducer (2)

```
import org.apache.hadoop.io.*;  
import org.apache.hadoop.mapred.*;
```

reduce (k2, [v2]) → [(k3, v3)]

```
public class MaxTemperatureReducer extends MapReduceBase  
    implements Reducer<Text, IntWritable, Text, DoubleWritable> {  
  
    public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text,  
        DoubleWritable> output, Reporter reporter) throws IOException {  
  
        int maxValue = Integer.MIN_VALUE;  
  
        while (values.hasNext()) {  
            maxValue = Math.max(maxValue, values.next().get());  
        }  
  
        output.collect(key, new DoubleWritable(((double)maxValue)/10));  
    }  
}
```

Temperature - Job

```
public class MaxTemperature {
```

```
    public static void main(String[] args) throws IOException {
```

```
        JobConf conf = new JobConf(MaxTemperature.class);  
        conf.setJobName("Max temperature");
```

```
        FileInputFormat.addInputPath(conf, new Path("temperature.txt"));  
        FileOutputFormat.setOutputPath(conf, new Path("output"));
```

```
        conf.setMapperClass(MaxTemperatureMapper.class);  
        conf.setReducerClass(MaxTemperatureReducer.class);
```

```
        conf.setMapOutputKeyClass(Text.class);  
        conf.setMapOutputValueClass(IntWritable.class);
```

```
        conf.setOutputKeyClass(Text.class);  
        conf.setOutputValueClass(DoubleWritable.class);
```

```
        JobClient.runJob(conf);
```

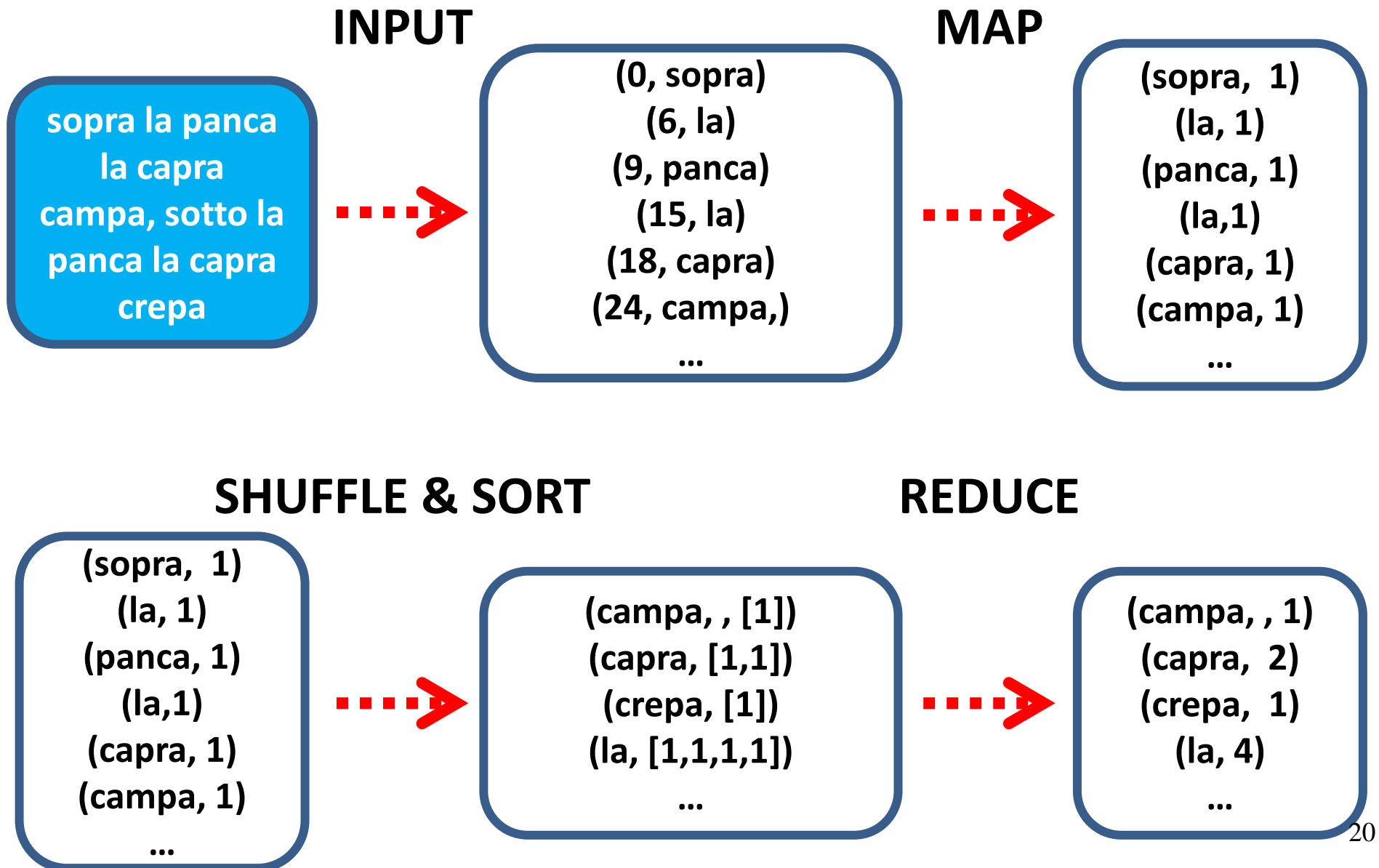
```
    }
```

```
}
```

Summary

- Install Pseudo-Distributed Mode
- Architecture
- Example: Temperature
- **Example: Word Count**
- Combiner
- Example: Bigram Count (with custom writable)
- How to run a job
- How to run a job on EMR

Word Count – logical data flow



Word Count – Mapper (1)

```
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.Mapper;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {

    private static final 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);
        }
    }
}
```

Word Count – Mapper (2)

```
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.Mapper;

public class WordCountMapper extends Mapper<LongWritable, Text, Text, IntWritable> {

    private static final 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);
        }
    }
}
```

Word Count – Reducer (1)

```
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.Reducer;

public class WordCountReducer 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 – Reducer (2)

```
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.Reducer;

public class WordCountReducer 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 – Job

```
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;

public class WordCount {

    public static void main(String[] args) throws Exception {

        Job job = new Job(new Configuration(), "WordCount");
        job.setJarByClass(WordCount.class);

        job.setMapperClass(WordCountMapper.class);
        job.setReducerClass(WordCountReducer.class);

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

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

        job.waitForCompletion(true);
    }
}
```

Summary

- Install Pseudo-Distributed Mode
- Architecture
- Example: Temperature
- Example: Word Count
- **Combiner**
- Example: Bigram Count (with custom writable)
- How to run a job
- How to run a job on EMR

Word Count: Without Combiner

MAP OUTPUT

(wordA,1)
(wordA,1)
(wordB,1)
(wordB,1)
(wordB,1)

(wordA,1)
(wordA,1)
(wordA,1)
(wordB,1)
(wordC,1)

REDUCER INPUT

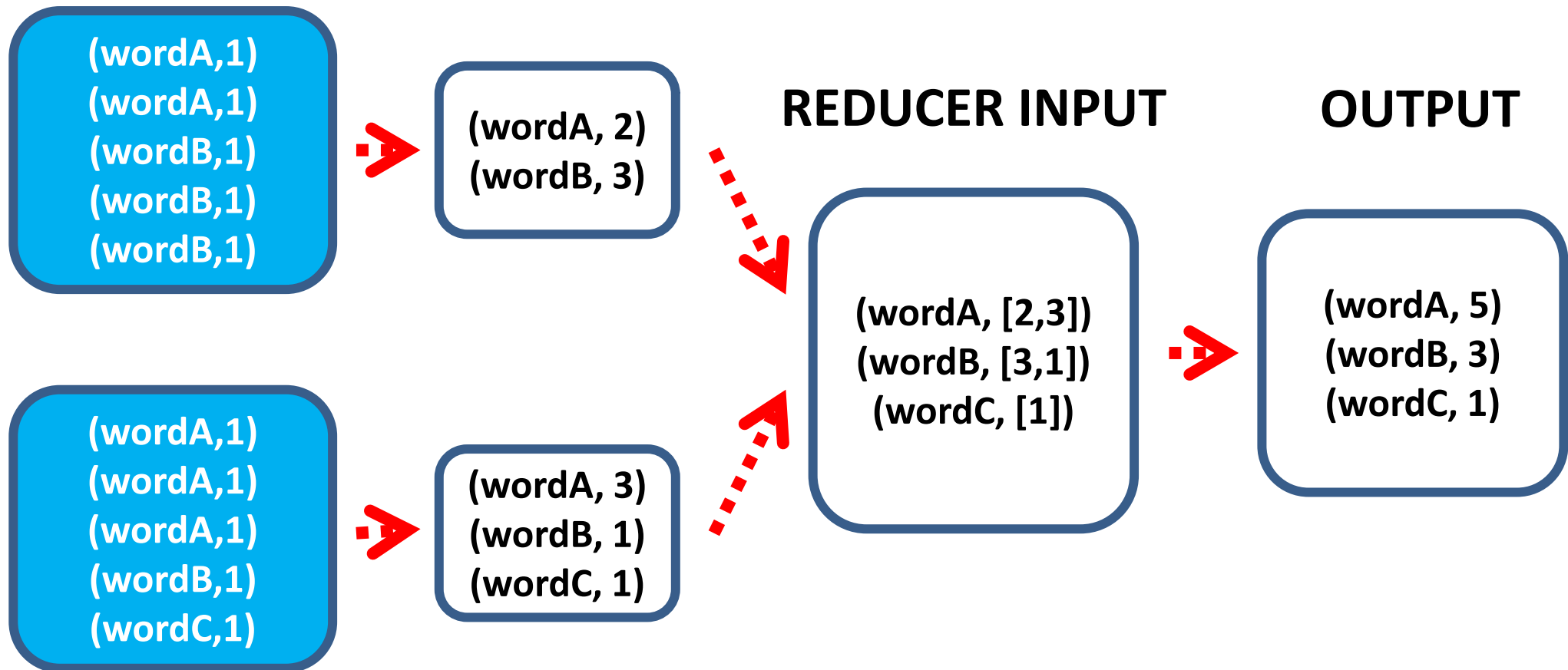
(wordA, [1,1,1,1,1])
(wordB, [1,1,1])
(wordC, [1])

OUTPUT

(wordA, 5)
(wordB, 3)
(wordC, 1)

Word Count: With Combiner

MAP OUTPUT



COMBINER OUTPUT

Combiner

In the job configuration:

...

```
job.setMapperClass(WordCountMapper.class);
```

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

```
job.setReducerClass(WordCountReducer.class);
```

...

Watch for types!

Summary

- Install Pseudo-Distributed Mode
- Architecture
- Example: Temperature
- Example: Word Count
- Combiner
- **Example: Bigram Count (with custom writable)**
- How to run a job
- How to run a job on EMR

Bigram Count

“A bigram is a couple of two consecutive words”

In this sentence there are 8 bigrams:

- A bigram
- bigram is
- is a
- a couple
- ...

Bigram Count

Example:

sopra la panca la
capra campà,
sotto la panca la
capra crepa



campà, sotto 1
capra campà, 1
capra crepa 1
la capra 2
la panca 2
panca la 2
sopra la 1
sotto la 1

To represent a Bigram we want to use a custom Writable type (BigramWritable)

Bigram Count

```
import org.apache.hadoop.io.*;

public class BigramWritable implements WritableComparable<BigramWritable> {

    private Text leftBigram;
    private Text rightBigram;

    public BigramWritable(Text left, Text right){
        this.leftBigram = left;
        this.rightBigram = right;
    }

    @Override
    public void readFields(DataInput in) throws IOException {
        leftBigram = new Text(in.readUTF());
        rightBigram = new Text(in.readUTF());
    }

    @Override
    public void write(DataOutput out) throws IOException {
        out.writeUTF(leftBigram.toString());
        out.writeUTF(rightBigram.toString());
    }
}
```

Bigram Count

...

```
public void set(Text prev, Text cur) {  
    leftBigram = prev;  
    rightBigram = cur;  
}
```

```
@Override  
public int hashCode() {  
    return leftBigram.hashCode() + rightBigram.hashCode();  
}
```

```
@Override  
public boolean equals(Object o) {  
    if (o instanceof BigramWritable) {  
        BigramWritable bigram = (BigramWritable) o;  
        return leftBigram.equals(bigram.leftBigram) &&  
            rightBigram.equals(bigram.rightBigram);  
    }  
    return false;  
}
```

...

Bigram Count

...

@Override

```
public int compareTo(BigramWritable tp) {  
    int cmp = leftBigram.compareTo(tp.leftBigram);  
    if (cmp != 0) {  
        return cmp;  
    }  
    return rightBigram.compareTo(tp.rightBigram);  
}  
}
```

Bigram Count - Mapper

```
public class BigramCountMapper extends Mapper<LongWritable, Text, BigramWritable,  
IntWritable> {
```

```
    private static final IntWritable ONE = new IntWritable(1);
```

```
    private static final BigramWritable BIGRAM = new BigramWritable();
```

```
    @Override
```

```
    public void map(LongWritable key, Text value, Context context)  
        throws IOException, InterruptedException {
```

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

```
        String prev = null;
```

```
        StringTokenizer itr = new StringTokenizer(line);
```

```
        while (itr.hasMoreTokens()) {
```

```
            String cur = itr.nextToken();
```

```
            if (prev != null) {
```

```
                BIGRAM.set( new Text(prev), new Text(cur) );
```

```
                context.write(BIGRAM, ONE);
```

```
            }
```

```
            prev = cur;
```

```
        }
```

```
    }
```

Bigram Count - Reducer

```
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.Reducer;

public class BigramCountReducer extends Reducer<BigramWritable,IntWritable,Text, IntWritable> {

    private final static IntWritable SUM = new IntWritable();

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

        int sum = 0;
        Iterator<IntWritable> iter = values.iterator();
        while (iter.hasNext()) {
            sum += iter.next().get();
        }
        SUM.set(sum);
        context.write(new Text( key.toString() ), SUM);
    }
}
```

Bigram Count - Job

```
import org.apache.hadoop.conf.Configured;
import org.apache.hadoop.util.ToolRunner;
import org.apache.hadoop.util.Tool;

public class BigramCount extends Configured implements Tool {

    private BigramCount() {}
    ...
    public int run(String[] args) throws Exception {
        ...
        Job job = new Job(getConf());
        job.setJobName(BigramCount.class.getSimpleName());
        job.setJarByClass(BigramCount.class);
        ...
        job.waitForCompletion(true);
    }

    public static void main(String[] args) throws Exception {
        ToolRunner.run(new BigramCount(), args);
    }
}
```

Summary

- Install Pseudo-Distributed Mode
- Architecture
- Example: Temperature
- Example: Word Count
- Combiner
- Example: Bigram Count (with custom writable)
- **How to run a job**
- How to run a job on EMR

How to run a job

How to make a jar file from eclipse:

- 1) Right click on the project folder
- 2) Export -> Java -> Runnable Jar
- 3) Select from the launch configuration the current project
- 4) Select a destination to export the .jar file
- 5) Select «Extract required libraries ..» and click Finish

How to copy a file in HDFS (not needed in standalone mode):

```
hadoop fs -copyFromLocal <src> <localDestination>
```


How to run a job

How to run a MapReduce job

```
hadoop jar <filejar> [arguments..]
```

Summary

- Install Pseudo-Distributed Mode
- Architecture
- Example: Temperature
- Example: Word Count
- Combiner
- Example: Bigram Count (with custom writable)
- How to run a job
- **How to run a job on EMR**

SD DEMO time!
