

# TP3 (Bonus): Developing MapReduce applications

Alessio Pagliari and Shadi Ibrahim

March 22, 2022

*After studying in the first TPs how to operate the Hadoop platform using examples, we will now implement MapReduce applications to solve several types of problems: including numerical calculation and analysis examples.*

## Exercise 1: My Wordcount application

As a simple illustration of the Map and Reduce functions, Figure 1 shows the pseudo-code and the algorithm and illustrates the process steps using the widely used “Wordcount” example. The Wordcount application counts the number of occurrences of each word in a large collection of documents. The steps of this process are briefly described as follows: The input is read (typically from a distributed file system) and broken up into key/value pairs (e.g., the Map function emits a word and its associated count of occurrence, which is just “1”). The pairs are partitioned into groups for processing, and they are sorted according to their key as they arrive for reduction. Finally, the key/value pairs are reduced, once for each unique key in the sorted list, to produce a combined result (e.g., the Reduce function sums all the counts emitted for a particular word).

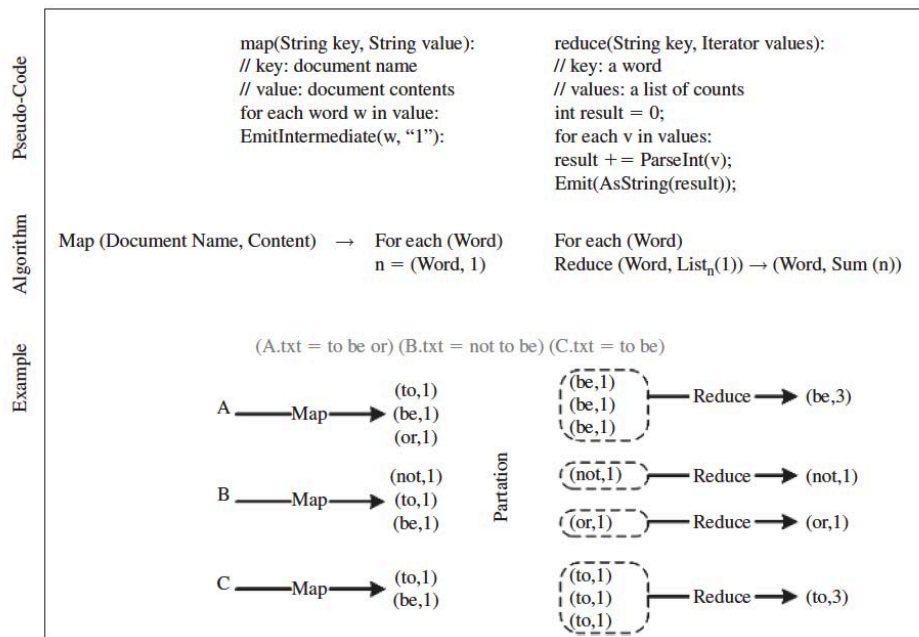


Figure 1

### Question 1.1

To use Eclipse to edit your MapReduce programs, you should add Hadoop libraries to your project. To do so, after creating a new project in your workspace (MyWordCount), add the two jar files given in the TP resources (hadoop-common-3.3.2.jar and hadoop-mapreduce-client-core-3.3.2.jar) to the Java Build Path:

- Select Build Path from the context menu (right click on the file) project then Configure Build Path.
- In the Libraries tab, click Add External JARs;
- Add the two Hadoop library files;

Now you are ready to create your first MapReduce program.

### Question 1.2

Copy the file **WordCount1.java** (you can find it in the resources) to your Eclipse project. By integrating this file in Eclipse, make sure the compilation goes well (no red cross).

### Question 1.3

Create you Jar file **MyWordCount.jar**.

### Question 1.4

Copy **MyWordCount.jar** to Hadoop folder and then run the wordcount with the data set 1000MB.

```
hadoop jar MyWordCount.jar MyWordCount input output
```

## Exercise 2: Modify your WordCount application

### Question 2.1

Change the number of Reduce task by adding:

```
job.setNumReduceTasks(X);
```

Create your Jar file and Run the job with 2 and 4 reducers.

### Question 2.2

Change the maximum size of map split size (the default size is 64 MB 67108864 bytes)

```
job.getConfiguration().setLong("mapred.max.split.size", X);
```

Create your Jar file and Run the job with two values  $n \times 512$ .

### Question 2.3

Create a new project (MyWourdCount1), copy WordCount1.java (named WordCount1.java in the resources) , now add and delete this line:

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

Run the program with and without the combiner function.

#### Question 2.4

Create a new project that gives you the Line Count (you can find an example in the resources).

#### Question 2.5

Create a new project that contains four wordcount applications (2 Reducers, 4 Reducers, 4 Reducers + combiner, 4 Reducers + combiner + split size = 5MB) and one linecount

//Run all with 2000MB file!

### Exercise 3 (Advanced): Developing your $\Pi$ application

This application estimates the value of pi based on sampling. The estimator first generates random points in a  $1 \times 1$  area. The map phase checks for each pair if it falls inside a 1-diameter circle; the reduce phase computes the ratio between the number of points inside the circle and the ones outside the circle. This ratio gives an estimate for the value of pi.

#### Question 3.1

To do this exercise, you will use the skeleton MyPiEstimator.java and the the HaltonSequence.java class (in the resources folder): A Java MapReduce program that will take the number of parameter maps, Reduces the number of and the number of samples per map.

- Define the input files for each map.
- Retrieve and parse the output files.
- Calculate and show the value of Pi.

#### Question 3.2

Write a Java code for the Map method.

#### Question 3.3

Write a Java code for the Reduce method.