



UNIVERSITY OF PISA

MSc in Artificial Intelligence and Data Engineering

Cloud Computing

Letter Frequency Analysis through Hadoop MapReduce

Project Report

Emanuele Respino

Christian Petruzzella

Alessandro Ciarniello

ACADEMIC YEAR 2023/2024

Contents

1	Algorithm Design	3
	MapReduce with Combiner	4
	In-Mapping Combining	6
2	Results	9
	Experimental Settings	9
	Performance Evaluation	9
	Letter Frequency	11

Chapter 1

Algorithm Design

The goal of this study is to provide the computing of letter frequency, given an input text: a **workflow** is required, allowing firstly the *letter counting* in the text, and then the *frequency computing* of each letter, based on results of the first job. This project provides the development and comparison of **two different workflow implementations**, which pseudo-codes are provided as following.

In *Letter counting*, a **Partitioner** has been implemented to equally distribute keys among each reducer: otherwise, due to the fact that *reducer key* has always the same value, this would result in having all the pairs redirected to the same reducer, nullifying the benefits of exploiting more reducers.

During the project, some *scripts* (bash files and Jupyter notebooks) have been implemented, leading to an easier plotting of results and a faster execution of the workflow with several configurations and different inputs.

MapReduce with Combiner

This is the classic implementation, in which a Combiner function is built to be optionally executed for local aggregation.

Algorithm 1 LetterCount with Combiner

Require: Txt file

Ensure: Total number of characters belonging to English alphabet

Class Mapper

```
1: procedure MAP(key, value, context)
2:   reducerKey  $\leftarrow$  NullWritable.get()
3:   reducerValue  $\leftarrow$  1
4:   CHARACTER_PATTERN  $\leftarrow$  regex("[a-z]", case_insensitive)
5:   // Convert the line to lower case and remove accents
6:   line  $\leftarrow$  removeAccents(value.toString()).toLowerCase()
7:   for each character ch in line do
8:     if CHARACTER_PATTERN matches ch then
9:       context.write(reducerKey, reducerValue)
10:    end if
11:  end for
12: end procedure
```

Class Partitioner

```
1: procedure PARTITION(key, value, numReduceTasks)
2:   return random() * numReduceTasks
3: end procedure
```

Class Reducer

```
1: procedure REDUCE(key, values, context)
2:   reducerKey  $\leftarrow$  NullWritable.get()
3:   reducerValue  $\leftarrow$  0
4:   for each value in values do
5:     reducerValue  $\leftarrow$  reducerValue + value
6:   end for
7:   // Write the output
8:   context.write(reducerKey, reducerValue)
9: end procedure
```

Algorithm 2 LetterFrequency with Combiner

Require: Txt file, Total number of characters in the txt file

Ensure: Frequency of each letter in the input file

Class Mapper

```
1: procedure MAP(key, value, context)
2:   reducerKey  $\leftarrow$  Text()
3:   reducerValue  $\leftarrow$  1
4:   CHARACTER_PATTERN  $\leftarrow$  regex("[a-z]", case_insensitive)
5:   // Convert the line to lower case and remove accents
6:   line  $\leftarrow$  removeAccents(value.toString()).toLowerCase()
7:   for each character ch in line do
8:     if CHARACTER_PATTERN matches ch then
9:       reducerKey.set(ch)
10:      context.write(reducerKey, reducerValue)
11:    end if
12:  end for
13: end procedure
```

Class Combiner

```
1: procedure COMBINE(key, values, context)
2:   sum  $\leftarrow$  0
3:   for each value in values do
4:     sum  $\leftarrow$  sum + value
5:   end for
6:   // Write the output
7:   context.write(key, sum)
8: end procedure
```

Class Reducer

```
1: procedure SETUP(context)
2:   map  $\leftarrow$  HashMap<Text, IntWritable>()
3: end procedure
4: procedure REDUCE(key, values, context)
5:   sum  $\leftarrow$  0
6:   for each value in values do
7:     sum  $\leftarrow$  sum + value
8:   end for
9:   // Write the output
10:  context.write(key, sum / TEXT_LENGTH)
11: end procedure
```

In-Mapping Combining

In-Mapper Combining design pattern provides local aggregation directly in the Mapper procedure, resulting in the complete control of the local aggregation process.

Algorithm 3 LetterCount with In-Mapper Combining

Require: Txt file

Ensure: Total number of characters belonging to English alphabet

Class Mapper

```
1: procedure SETUP(context)
2:   map  $\leftarrow$  HashMap<NullWritable, IntWritable>()
3: end procedure

4: procedure MAP(key, value, context)
5:   // Convert the line to lower case and remove accents
6:   line  $\leftarrow$  removeAccents(value.toString()).toLowerCase()
7:   for each character ch in line do
8:     if CHARACTER_PATTERN matches ch then
9:       if map contains NullWritable then
10:        map[NullWritable]  $\leftarrow$  map[NullWritable] + 1
11:      else
12:        map[NullWritable]  $\leftarrow$  1
13:      end if
14:    end if
15:  end for
16: end procedure

17: procedure CLEANUP(context)
18:   for each entry in map do
19:     context.write(entry.getKey(), entry.getValue())
20:   end for
21: end procedure
```

Class Partitioner

```
1: function PARTITION(key, value, numReduceTasks)
2:   return random() * numReduceTasks
3: end function
```

Class Reducer

```
1: procedure REDUCE(key, values, context)
2:   reducerKey  $\leftarrow$  NullWritable.get()
3:   reducerValue  $\leftarrow$  0
4:   for each value in values do
5:     reducerValue  $\leftarrow$  reducerValue + value
6:   end for
7:   context.write(reducerKey, reducerValue)
8: end procedure
```

Algorithm 4 LetterFrequency with In-Mapping Combining

Require: Txt file, Total number of characters in the txt file

Ensure: Frequency of each letter in the input file

Class Mapper

```
1: procedure SETUP(context)
2:   map  $\leftarrow$  HashMap<Text, IntWritable>()
3: end procedure

4: procedure MAP(key, value, context)
5:   // Convert the line to lower case and remove accents
6:   line  $\leftarrow$  removeAccents(value.toString()).toLowerCase()
7:   for each character ch in line do
8:     if CHARACTER_PATTERN matches ch then
9:       if map contains Text(ch) then
10:        map[Text(ch)]  $\leftarrow$  map[Text(ch)] + 1
11:      else
12:        map[Text(ch)]  $\leftarrow$  1
13:      end if
14:    end if
15:  end for
16: end procedure

17: procedure CLEANUP(context)
18:   for each entry in map do
19:     context.write(entry.getKey(), entry.getValue())
20:   end for
21: end procedure
```

Class Reducer

```
1: procedure SETUP(context)
2:   Configuration conf  $\leftarrow$  context.getConfiguration()
3:   TEXT_LENGTH  $\leftarrow$  parseInteger(conf.get("textLength"))
4: end procedure

5: procedure REDUCE(key, values, context)
6:   sum  $\leftarrow$  0
7:   for each value in values do
8:     sum  $\leftarrow$  sum + value
9:   end for
10:  context.write(key, sum / TEXT_LENGTH)
11: end procedure
```

Chapter 2

Results

Experimental Settings

Two different analysis have been provided:

- **Performances:** Executing Map-Reduce workflow with different configurations leads to a better comprehension of its behavior:
 - *Type of workflow:* with Combiner, In-Mapper Combining;
 - *Size of input data:* 15.2KB, 2.93MB, 2.82GB (randomly generated);
 - *Number of Reducers:* 1, 2, 3.
- **Letter frequency:** using a fixed configuration, the letter frequency of several languages have been compared, using the same text as reference, in order to understand similarities and differences between each other:
 - *Type of workflow:* In-Mapper Combining;
 - *Number of Reducers:* 3;
 - *Input text:* The Bible (~4.2MB), from Bible SuperSearch API;
 - *Languages:* Afrikaans, Albanian, Czech, Dutch, English, Finnish, French, German, Hungarian, Indonesian, Italian, Maori, Polish, Portuguese, Romanian, Spanish, Turkish.

Performance Evaluation

For the performance evaluation, 5 parameters have been considered among the ones Hadoop provides at the end of each Map-Reduce job:

- **Total time spent by all map tasks** (ms);
- **Total time spent by all reduce tasks** (ms);

- CPU time spent (ms);
- Garbage Collection (GC) time elapsed (ms);
- Reduce shuffle bytes.

The results of those parameters for each configuration are shown in the following charts.

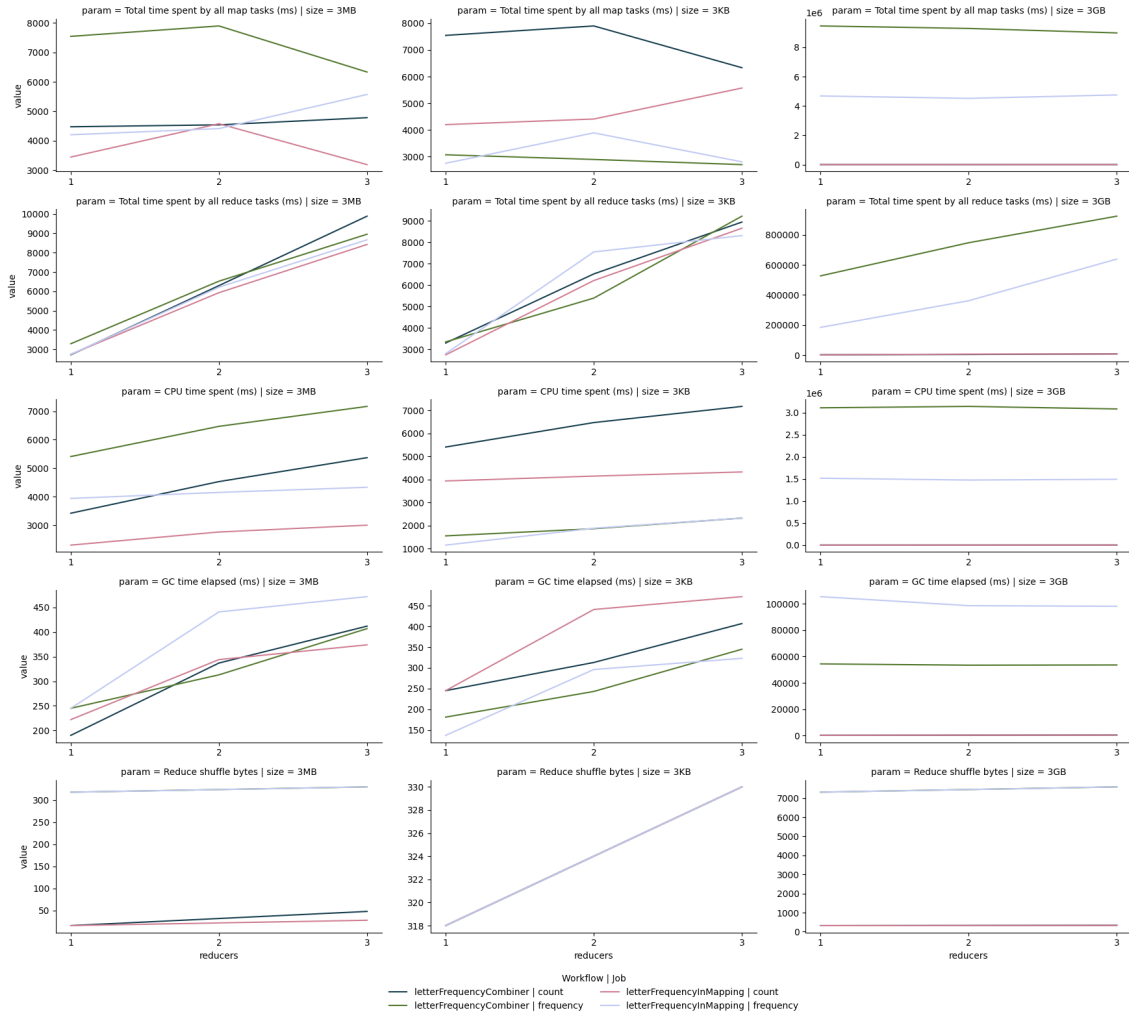


Figure 2.1: Performance comparison on some parameters for each configuration

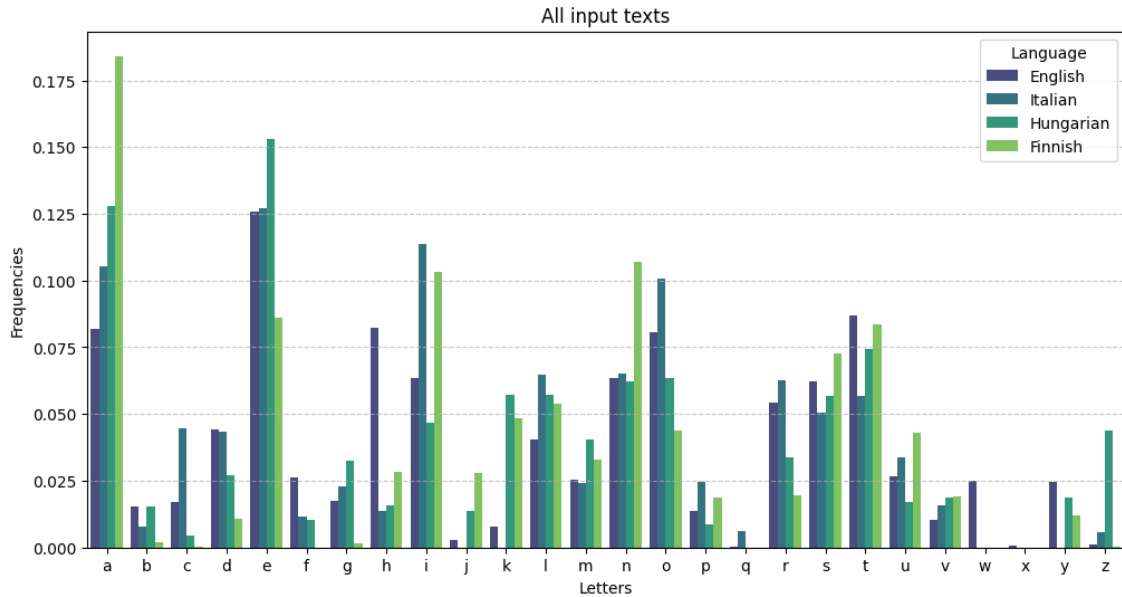
Workflow. The results obtained show that the usage of In-Mapper Combining reduces execution time, especially with larger input files. On the other hand this implementation leads to a much higher *GC time elapsed* value, showing that it is more memory consuming rather than the implementation with the Combiner option.

Reducers. Considering the specific purpose of this study, it is shown that increasing the number of reducers do not lead to an overall significant improvement. Actually, for input texts of lower dimension, using a higher number reducers is less time efficient.

Letter Frequency

The analysis has been conducted for all the languages previously mentioned.

A more detailed comparison is reported, highlighting the differences and similarities among four European languages: **English**, **Italian**, **Hungarian** and **Finnish**.



This report shows some interesting results:

- **Most common letters:** *E* is the most common letter in English, Italian and Hungarian, while for Finnish it is the letter *A*: this can indicate differences in the vowel structure between this language and the others.
- **Rare letters:** *Q*, *W*, *X*, *Z* tends to be rare or absent, except for *Z* in Hungarian and *W* in English.
- **Peaks in frequencies:** *C* is relatively more used in Italian rather than in the other languages, as for *H* in English and *N* in Finnish.

The following graphs focus on the most common letters for each language.

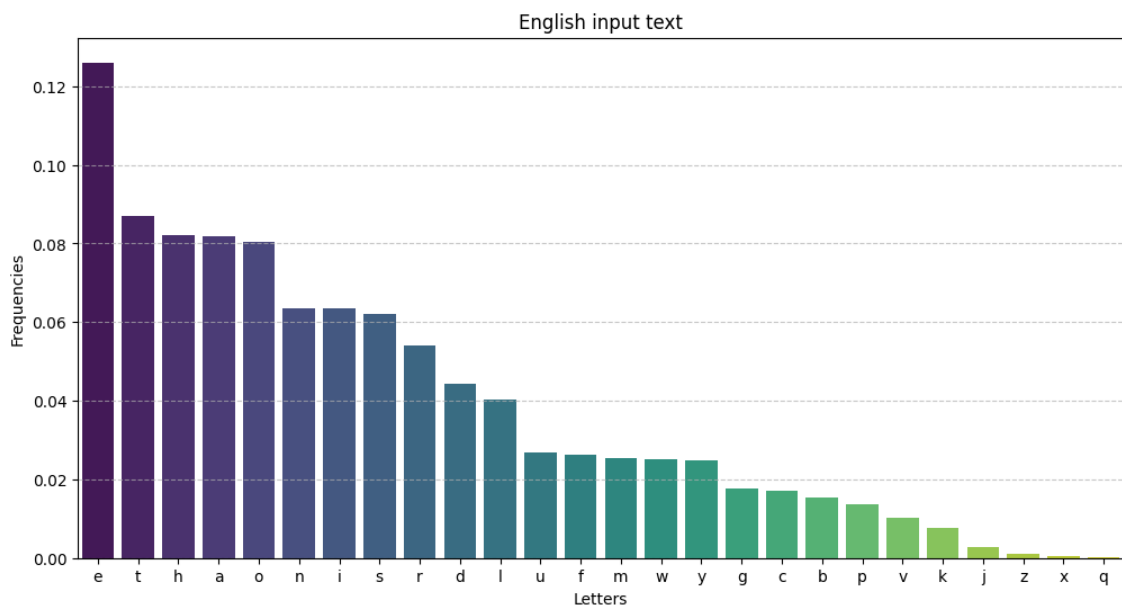


Figure 2.2: English most-common letters

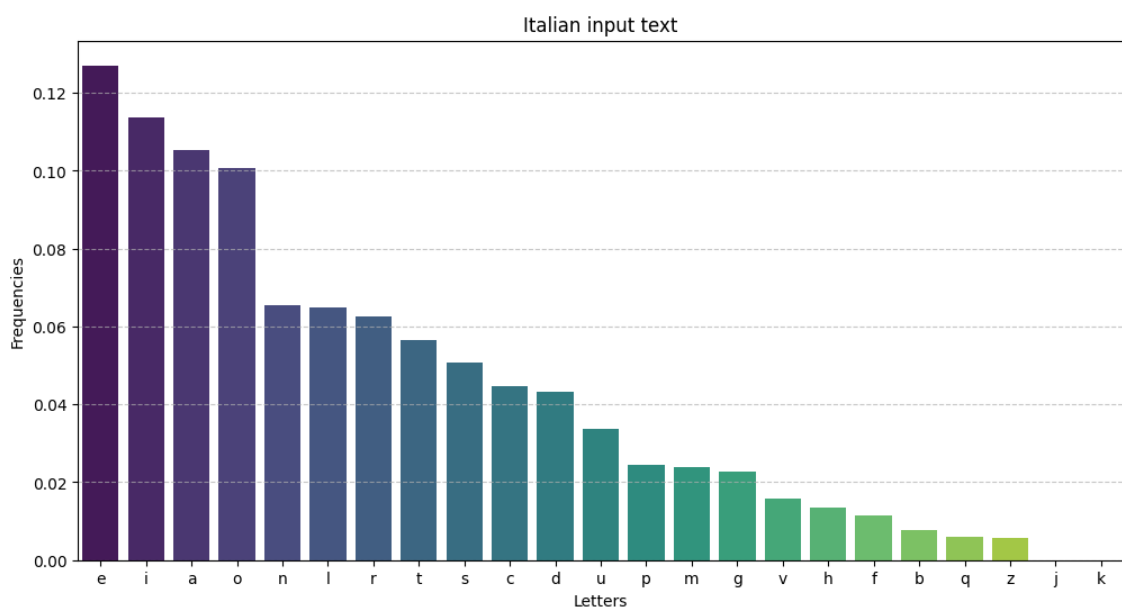


Figure 2.3: Italian most-common letters

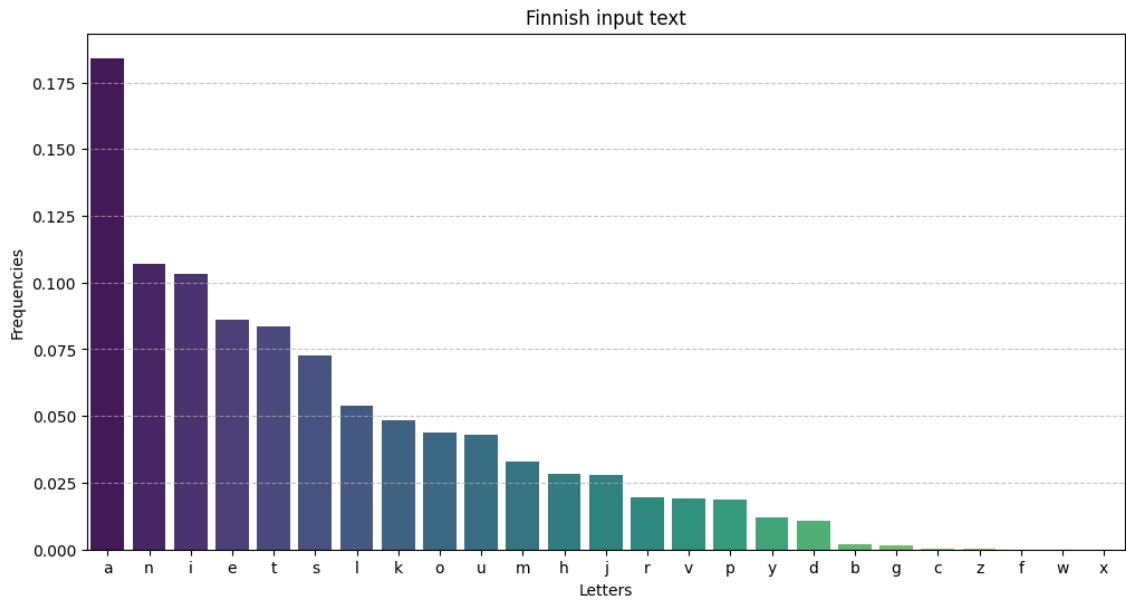


Figure 2.4: Finnish most-common letters

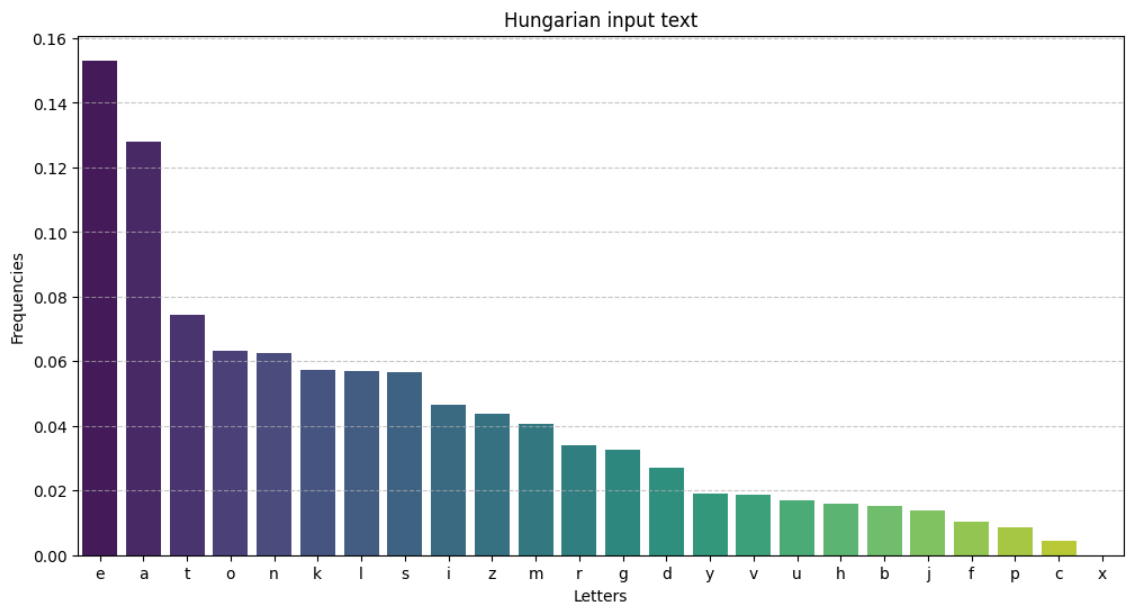


Figure 2.5: Hungarian most-common letters