

Cloud Computing - Project Report

Totaro Group

Zahra Omrani
z.omrani@studenti.unipi.it

Paolo Palumbo
p.palumbo3@studenti.unipi.it

Ettore Ricci
e.ricci32@studenti.unipi.it

June 14, 2024

<https://github.com/Etto48/CloudComputingProject>

Abstract

This report presents the development and performance evaluation of a Java-based application designed to count letter frequencies in large datasets using the Hadoop framework. To provide a comprehensive analysis, we also implemented and tested non-distributed applications in Python and Rust for comparison. The results show that for the sizes of the datasets considered, Java and Python execution times are comparable, while Rust is significantly faster. Both Python and Rust applications used significantly less memory than the Java application.

1 Introduction

Our project for the Cloud Computing course consists of developing a Java application to count the frequency of letters in a large dataset using the Hadoop framework. The main objective of the project is to analyze and compare the performance of the application with different configurations and input sizes and also to compare it with non-distributed implementations in Python and Rust.

- **Mapper with in-mapper combiner:** This mapper stores the letter counts in a vector and emits them in the cleanup method. At the end of the cleanup method, the mapper increments a job counter with the number of emitted letters. (Listing 1)
- **Mapper without in-mapper combiner:** This mapper emits the letter counts (with value 1) for each letter in the map method. At the end of the map method, the mapper increments a job counter with the number of emitted letters. (Listing 2)

2 Mapreduce

MapReduce is a programming model for processing large data sets. Users specify a map function that processes a key/value pair to generate a set of intermediate key/value pairs, and a reduce function that merges all intermediate values associated with the same intermediate key. For our application we needed two subsequent MapReduce jobs: the first one to count the occurrences of each letter in the input text and the second one to calculate the frequency of each letter.

2.1 Job 1: Counting letters

2.1.1 Mapper

For testing purposes, we created two different mappers: one with an in-mapper combiner and one without. For the combiner (when enabled), we used the same code as the reducer.

2.1.2 Reducer

The reducer receives the letter counts from the mappers and sums them up. (Listing 3)

2.2 Job 2: Calculating frequencies

The results of the first job are saved in a temporary file and used as input for the second job. The counter obtained from the first job is stored in a configuration variable that is passed to the reducer (only one is needed because the number of different letters is small). For this job we do not need a combiner so it's always disabled.

2.2.1 Mapper

We used an identity mapper that emits the letter counts as they are.

2.2.2 Reducer

The reducer receives the letter counts from the mappers and calculates the frequency of each letter using the total count that was previously stored in a configuration variable. (Listing 4)

3 Dataset

We used three different datasets:

- **english.txt**: A 1.2GB text file containing text from the books on [Gutenberg project](#).
- **italian.txt**: A 1.3GB text file from the [PAISÀ corpus](#)[1].
- **spanish.txt**: A 130MB text file from the [Leipzig Corpora Collection](#)[2].

We also created 11 files of increasing size (from 100MB to 1.1GB with 100MB steps) to test the performance of the application with different input sizes. These files are called **part_XMB.txt** and contain the first XMBs of the **english.txt** file.

4 Experiments

We wrote our program in a way that allows us to configure it for the different tests with command line arguments:

- `-i` | `--input`: Input file path inside the HDFS.
- `-r` | `--reducers`: Number of reducers to use in the first job.
- `-n` | `--no-combiner`: Disable the combiner in the first job.
- `-m` | `--no-in-mapper-combiner`: Disable the in-mapper combiner in the first job.
- `-o` | `--output`: Output directory path inside the HDFS, by default it's set to **output**.
- `-t` | `--tmp`: Temporary directory path inside the HDFS, by default it's set to **tmp**.

We carried out 21 tests with the configurations shown in Table 1 and Table 2. Originally we ran tests 7, 8 and 9 with 4 reducers, but test 9 could never complete because of a Shuffle Error caused by a Java Heap Space error (Out Of Memory Error). Because of this, we decided to run these tests again with 3 reducers. Only the 3 reducers configuration will be shown in the results. We also ran the Python and Rust applications with the same input file (english.txt).

Test ID	Arguments
0	<code>-i english.txt</code> <code>-r 1</code>
1	<code>-i english.txt</code> <code>-r 2</code>
2	<code>-i english.txt</code> <code>-r 4</code>
3	<code>-i english.txt</code> <code>-r 8</code>
4	<code>-i english.txt</code> <code>-r 1</code> <code>--no-combiner</code>
5	<code>-i english.txt</code> <code>-r 1</code> <code>--no-in-mapper-combiner</code>
6	<code>-i english.txt</code> <code>-r 1</code> <code>--no-in-mapper-combiner</code> <code>--no-combiner</code>
7	<code>-i english.txt</code> <code>-r 3</code> <code>--no-combiner</code>
8	<code>-i english.txt</code> <code>-r 3</code> <code>--no-in-mapper-combiner</code>
9	<code>-i english.txt</code> <code>-r 3</code> <code>--no-in-mapper-combiner</code> <code>--no-combiner</code>

Table 1: General tests

Test ID	Arguments
10	<code>-i part_100MB.txt</code> <code>-r 1</code>
11	<code>-i part_200MB.txt</code> <code>-r 1</code>
12	<code>-i part_300MB.txt</code> <code>-r 1</code>
13	<code>-i part_400MB.txt</code> <code>-r 1</code>
14	<code>-i part_500MB.txt</code> <code>-r 1</code>
15	<code>-i part_600MB.txt</code> <code>-r 1</code>
16	<code>-i part_700MB.txt</code> <code>-r 1</code>
17	<code>-i part_800MB.txt</code> <code>-r 1</code>
18	<code>-i part_900MB.txt</code> <code>-r 1</code>
19	<code>-i part_1000MB.txt</code> <code>-r 1</code>
20	<code>-i part_1100MB.txt</code> <code>-r 1</code>

Table 2: Input split tests

5 Results

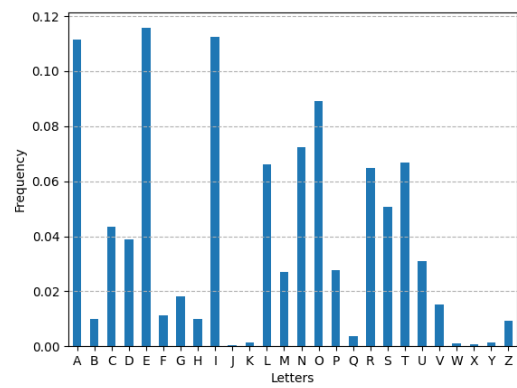


Figure 2: Italian letter frequency

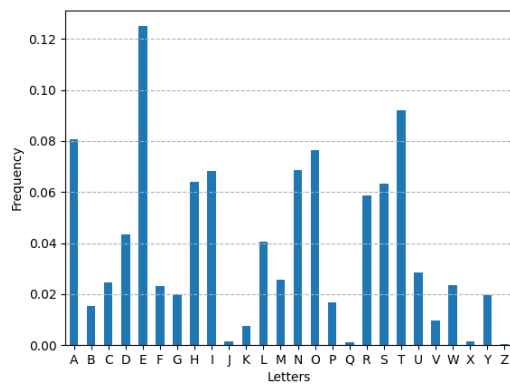


Figure 1: English letter frequency

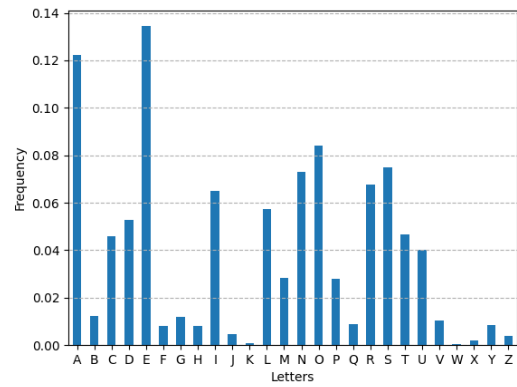


Figure 3: Spanish letter frequency

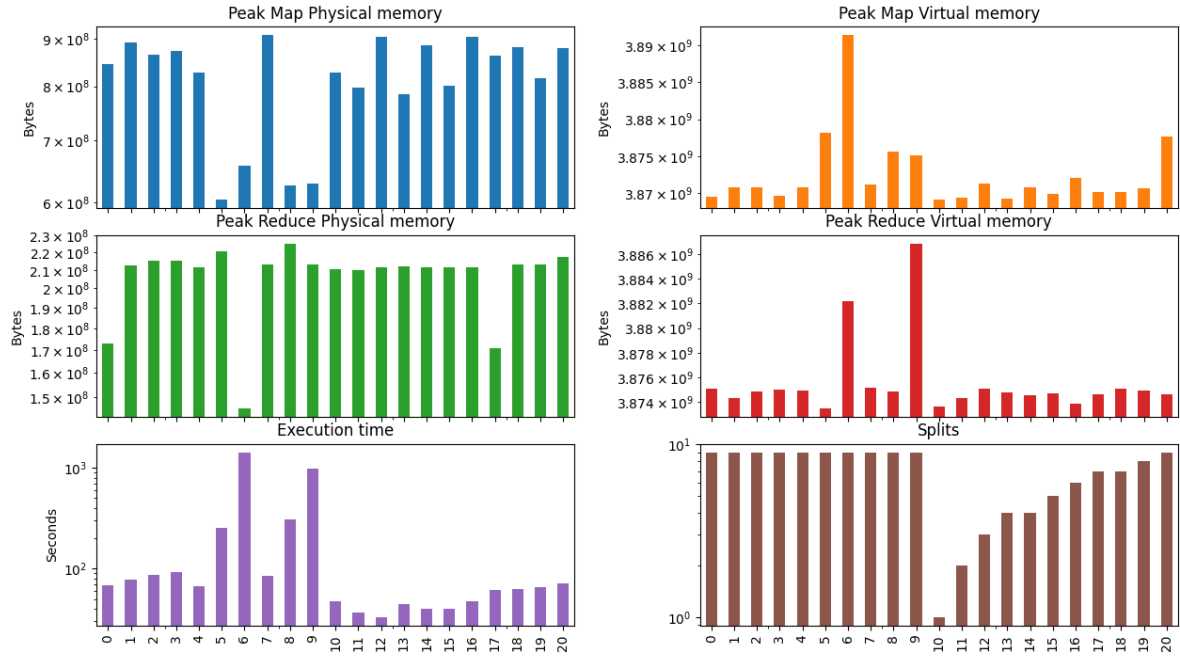


Figure 4: Tests results, the x-axis represents the test ID.

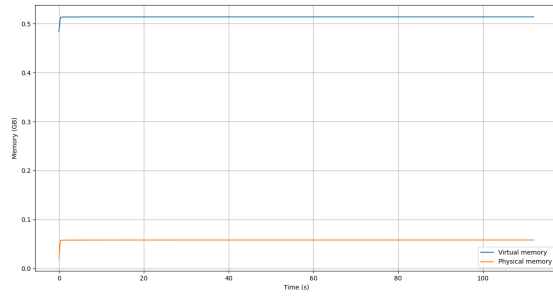


Figure 5: Python baseline memory usage (GB) over time (s)

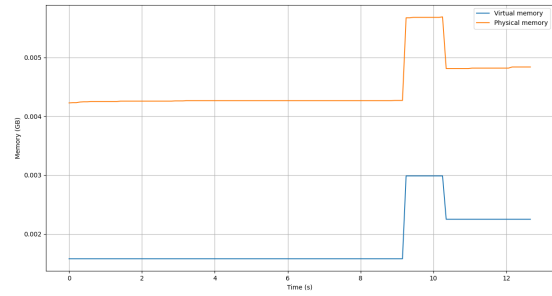


Figure 6: Rust baseline memory usage (GB) over time (s)

6 Conclusions

7 Listings

Listing 1: Mapper with in-mapper combiner

```

1 package it.unipi;
2
3 import java.io.IOException;
4 import java.util.Vector;
5
6 import org.apache.commons.lang3.StringUtils;
7 import org.apache.hadoop.io.LongWritable;
8 import org.apache.hadoop.io.Text;
9
10 public class Mapper extends org.apache.hadoop.mapreduce.Mapper<Object, Text,
11     Char, LongWritable> {
12     private final static int LETTER_COUNT = 26;
13     private final Vector<Integer> combinerCounters = new Vector<Integer>();

```

```

13     private long sum = 0;
14     private final Char key = new Char();
15     private final LongWritable value = new LongWritable();
16
17     @Override
18     public void setup(Context context) {
19         sum = 0;
20         combinerCounters.clear();
21         for (int i = 0; i < LETTER_COUNT; i++) {
22             combinerCounters.add(0);
23         }
24     }
25
26     @Override
27     public void cleanup(Context context) throws IOException,
28         InterruptedException {
29         for (int i = 0; i < LETTER_COUNT; i++) {
30             key.set((char)('a' + i));
31             value.set(combinerCounters.get(i));
32             context.write(key, value);
33         }
34         context.getCounter("LetterFreq", "total").increment(sum);
35     }
36
37     @Override
38     public void map(Object key, Text value, Context context) throws IOException,
39         InterruptedException {
40         StringUtils.stripAccents(value.toString()).chars().forEach(c -> {
41             int lowerC = Character.toLowerCase(c);
42             if (lowerC >= 'a' && lowerC <= 'z') {
43                 int index = lowerC - 'a';
44                 combinerCounters.set(index, combinerCounters.get(index) + 1);
45                 sum++;
46             }
47         });
48     }
49 }

```

Listing 2: Mapper without in-mapper combiner

```

1 package it.unipi;
2
3 import java.io.IOException;
4
5 import org.apache.commons.lang3.StringUtils;
6 import org.apache.hadoop.io.LongWritable;
7 import org.apache.hadoop.io.Text;
8
9 public class MapperNoCombiner extends org.apache.hadoop.mapreduce.Mapper<Object
10     , Text, Char, LongWritable> {
11     private final Char key = new Char();
12     private final LongWritable one = new LongWritable(1);
13
14     @Override
15     public void map(Object key, Text value, Context context) throws IOException,
16         InterruptedException {
17         long sum = 0;
18         String normalizedInput = StringUtils.stripAccents(value.toString());
19         for (int i = 0; i < normalizedInput.length(); i++) {
20             int lowerC = Character.toLowerCase(normalizedInput.charAt(i));
21             if (lowerC >= 'a' && lowerC <= 'z') {
22                 this.key.set((char)lowerC);
23                 context.write(this.key, this.one);
24                 sum++;
25             }
26         }
27     }
28 }

```

```

23         }
24     };
25     context.getCounter("LetterFreq", "total").increment(sum);
26 }
27 }

```

Listing 3: Reducer

```

1 package it.unipi;
2
3 import java.io.IOException;
4
5 import org.apache.hadoop.io.LongWritable;
6
7 public class Reducer extends org.apache.hadoop.mapreduce.Reducer<
8     Char,
9     LongWritable,
10    Char,
11    LongWritable>
12 {
13     private final LongWritable result = new LongWritable();
14     public void reduce(Char key, Iterable<LongWritable> values, Context context
15         )
16         throws IOException, InterruptedException {
17         long sum = 0;
18         for (LongWritable val : values) {
19             sum += val.get();
20         }
21         result.set(sum);
22         context.write(key, result);
23     }
24 }

```

Listing 4: ReducerFrequency class, used to calculate the frequency of each letter

```

1 package it.unipi;
2
3 import java.io.IOException;
4
5 import org.apache.hadoop.io.DoubleWritable;
6 import org.apache.hadoop.io.LongWritable;
7
8 public class ReducerFrequency extends org.apache.hadoop.mapreduce.Reducer<
9     Char,
10    LongWritable,
11    Char,
12    DoubleWritable>
13 {
14     private long total;
15
16     @Override
17     public void setup(Context context) {
18         total = context.getConfiguration().getLong("total", 0);
19         if (total == 0) {
20             throw new RuntimeException("Total count is zero");
21         }
22     }
23
24     private final DoubleWritable result = new DoubleWritable();
25     public void reduce(Char key, Iterable<LongWritable> values, Context context
26         )
27         throws IOException, InterruptedException
28     {
29         long sum = 0;

```

```

29         for (LongWritable val : values) {
30             sum += val.get();
31         }
32         result.set((double)sum / total);
33         context.write(key, result);
34     }
35 }

```

Listing 5: Char class, used as a custom key

```

1 package it.unipi;
2
3 import java.io.DataInput;
4 import java.io.DataOutput;
5 import java.io.IOException;
6
7 import org.apache.hadoop.io.WritableComparable;
8
9 public class Char implements WritableComparable<Char> {
10     private char c;
11
12     public Char() {
13         c = 0;
14     }
15
16     public Char(char c) {
17         this.c = c;
18     }
19
20     public void set(char c) {
21         this.c = c;
22     }
23
24     public char get() {
25         return c;
26     }
27
28     @Override
29     public void readFields(DataInput arg0) throws IOException {
30         c = arg0.readChar();
31     }
32
33     @Override
34     public void write(DataOutput arg0) throws IOException {
35         arg0.writeChar(c);
36     }
37
38     @Override
39     public int compareTo(Char o) {
40         return Character.compare(c, o.c);
41     }
42
43     @Override
44     public boolean equals(Object o) {
45         if (o instanceof Char) {
46             return c == ((Char)o).c;
47         }
48         return false;
49     }
50
51     @Override
52     public int hashCode() {
53         return Character.hashCode(c);
54     }

```

```
55
56     @Override
57     public String toString() {
58         return Character.toString(c);
59     }
60 }
```

References

- [1] V. Lyding, E. Stemle, C. Borghetti, M. Brunello, S. Castagnoli, F. Dell’Orletta, H. Dittmann, A. Lenci, and V. Pirrelli, “The PAISÀ corpus of Italian web texts,” in *Proceedings of the 9th Web as Corpus Workshop (WaC-9)*, F. Bildhauer and R. Schäfer, Eds. Gothenburg, Sweden: Association for Computational Linguistics, Apr. 2014, pp. 36–43. [Online]. Available: <https://aclanthology.org/W14-0406>
- [2] T. Eckart and U. Quasthoff, *Statistical Corpus and Language Comparison on Comparable Corpora*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2013, pp. 151–165. [Online]. Available: https://doi.org/10.1007/978-3-642-20128-8_8