

UNIVERSITY OF PISA



Artificial Intelligence and Data Engineering

Cloud Computing

Hadoop Letter Frequency

Authors: Gemelli M. Namaki D. Nocella F.

Academic Year 2023/2024

Contents

1	Introduction	2
2	Algorithm Design	3
2.1	MapReduce with Combiner	3
2.1.1	Pseudocode	3
2.2	MapReduce with In-Mapper Combiner	5
2.2.1	Pseudocode	5
3	Results	7
3.1	Experimental Setup	7
3.2	Performance Analysis	7
3.3	Letter Frequency Analysis	11
4	Conclusions	14
4.1	Performance	14
4.2	Qualitative Analysis	14

Introduction

The objective of this project is to implement a data processing pipeline that can handle substantial data sets, ensuring efficient computation and meaningful data insights. By exploiting the MapReduce paradigm, data processing tasks are split into two main functions: the **Mapper** and the **Reducer**.

The Mapper function processes the input data, emitting **key-value pairs**, while the Reducer aggregates these pairs to produce the final output.

The project is developed using Hadoop framework, which provides an open-source implementation of the MapReduce paradigm. Hadoop allows for the distributed processing of large data sets across clusters of computers using simple programming models. The Hadoop ecosystem also includes other tools, such as HDFS (Hadoop Distributed File System) for distributed storage, and YARN (Yet Another Resource Negotiator) for cluster resource management.

Algorithm Design

The objective of this project is to analyze letter frequency in text documents utilizing Hadoop's MapReduce framework. Specifically, two distinct approaches are adopted to optimize the MapReduce task: the use of a Combiner and the implementation of an In-Mapper Combiner. These methods aim to enhance the efficiency of the MapReduce process by reducing the amount of data transferred between the Mapper and Reducer stages.

MapReduce with Combiner

The Combiner is a mini-reducer that processes the output of the Mapper tasks before passing it to the Reducer. By aggregating the intermediate data locally on the mapper nodes, the Combiner reduces the volume of data shuffled across the network, thus improving the performance of the MapReduce job. It performs its operations on the same node where the mapper is running.

Pseudocode

Algorithm 1 Letter Count with Combiner

Require: Txt file

Ensure: Total count of each letter in the input file

Mapper

```
1: procedure SETUP(context)
2:   normalize  $\leftarrow$  context.getConfiguration().get("normalize")
3:   one  $\leftarrow$  new LongWritable(1)
4:   letterCountKey  $\leftarrow$  new Text("total_letter_count")
5: end procedure

6: procedure MAP(Object key, Text value)
7:   line  $\leftarrow$  Normalize(value.toString(), normalize)  $\triangleright$  Remove accents and set lowercase
8:   for each character  $c$  in line do
9:     Emit("total_letter_count", one)
10:  end for
11: end procedure
```

Combiner & Reducer

```
12: procedure REDUCE(Text key, Iterable<LongWritable> values)
13:   sum  $\leftarrow$  0
14:   for each LongWritable val in values do
15:     sum  $\leftarrow$  sum + val.get()
16:   end for
17:   Emit(key, new LongWritable(sum))
18: end procedure
```

Algorithm 2 Letter Frequency with Combiner

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

Ensure: Frequency of each letter in the input file

Mapper

```
1: procedure SETUP(context)
2:   normalize  $\leftarrow$  context.getConfiguration().get("normalize")
3:   one  $\leftarrow$  new LongWritable(1)
4: end procedure

5: procedure MAP(Object key, Text value)
6:   line  $\leftarrow$  Normalize(value.toString(), normalize)     $\triangleright$  Remove accents and set lowercase
7:   for each character  $c$  in line do
8:     Emit(String.valueOf( $c$ ), one)
9:   end for
10: end procedure
```

Combiner

```
11: procedure REDUCE(Text key, Iterable<LongWritable> values)
12:   sum  $\leftarrow$  0
13:   for each LongWritable val in values do
14:     sum  $\leftarrow$  sum + val.get()
15:   end for
16:   Emit(key, new LongWritable(sum))
17: end procedure
```

Reducer

```
18: procedure SETUP(context)
19:   letterCount  $\leftarrow$  context.getConfiguration().getLong("letterCountValue", 1)
20: end procedure

21: procedure REDUCE(Text key, Iterable<LongWritable> values)
22:   sum  $\leftarrow$  0
23:   for each LongWritable val in values do
24:     sum  $\leftarrow$  sum + val.get()
25:   end for
26:   freq  $\leftarrow$  (double) sum / (double) letterCount
27:   Emit(key, new DoubleWritable(freq))
28: end procedure
```

MapReduce with In-Mapper Combiner

The In-Mapper Combiner combines the mapping and combining steps within the Mapper itself. This method involves accumulating the results in a data structure within the Mapper, which is then emitted at the end of the mapping phase. This approach minimizes the overhead of multiple data passes by efficiently combining intermediate results within the Mapper, reducing the need for external Combiner steps and further optimizing network usage and processing time.

Pseudocode

Algorithm 3 Letter Count with In-Mapper Combiner

Require: Txt file

Ensure: Total count of each letter in the input file

CountMapper

```
1: procedure SETUP(Context context)
2:   normalize  $\leftarrow$  context.getConfiguration().get("normalize")
3:   map  $\leftarrow$  {}
4:   letterCountKey  $\leftarrow$  new Text("total_letter_count")
5: end procedure

6: procedure MAP(Object key, Text value)
7:   line  $\leftarrow$  Normalize(value.toString(), normalize)
8:   for each character  $c$  in line do
9:     map{letterCountKey}  $\leftarrow$  map{letterCountKey} + 1
10:  end for
11: end procedure

12: procedure CLEANUP(Context context)
13:   for each entry  $\langle k, v \rangle$  in map do
14:     Emit( $k, v$ )  $\triangleright$  Emit key and count map{key}
15:   end for
16: end procedure
```

CountReducer

```
17: procedure REDUCE(Text key, Iterable<LongWritable> values)
18:   sum  $\leftarrow$  0
19:   for each value in values do
20:     sum  $\leftarrow$  sum + value
21:   end for
22:   Emit(key, sum)
23: end procedure
```

Algorithm 4 Letter Frequency with In-Mapper Combiner

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

Ensure: Frequency of each letter in the input file

CountMapper

```
1: procedure SETUP(Context context)
2:   normalize  $\leftarrow$  context.getConfiguration().get("normalize")
3:   map  $\leftarrow$  {}
4: end procedure

5: procedure MAP(Object key, Text value)
6:   line  $\leftarrow$  Normalize(value.toString(), normalize)
7:   for each character  $c$  in line do
8:     map{ $c$ }  $\leftarrow$  map{ $c$ } + 1
9:   end for
10: end procedure

11: procedure CLEANUP(Context context)
12:   for each entry  $< k, v >$  in map do                                 $\triangleright$  Emit key and count map{key}
13:     Emit( $k, v$ )
14:   end for
15: end procedure
```

CountReducer

```
16: procedure SETUP(Context context)
17:   letterCount  $\leftarrow$  context.getConfiguration().getLong("letterCountValue", 1)
18: end procedure

19: procedure REDUCE(Text key, Iterable<LongWritable> values)
20:   sum  $\leftarrow$  0
21:   for each value in values do
22:     sum  $\leftarrow$  sum + value
23:   end for
24:   freq  $\leftarrow$  (double) sum / (double) letterCount
25:   Emit(key, freq)
26: end procedure
```

Results

Experimental Setup

Executing the MapReduce workflow with different inputs, and configurations:

- **Input size:** the size of the input file is varied to evaluate the performance of the MapReduce workflow.
 - **Paradise Lost** ~ 310 kB
 - **Moby Dick** ~ 421kB
 - **Frankenstein** ~ 440 kB
 - **Divina Commedia** ~ 600 kB
 - **Gerusalemme Liberata** ~ 691 kB
 - **Promessi Sposi** ~ 1.440 kB
 - **Test file** (random generated sequence of char) of ~ 800 MB.
- **Number of mappers:** Hadoop handles this step
- **Number of reducers:** From one up to three reducers are used for letter frequency

Performance Analysis

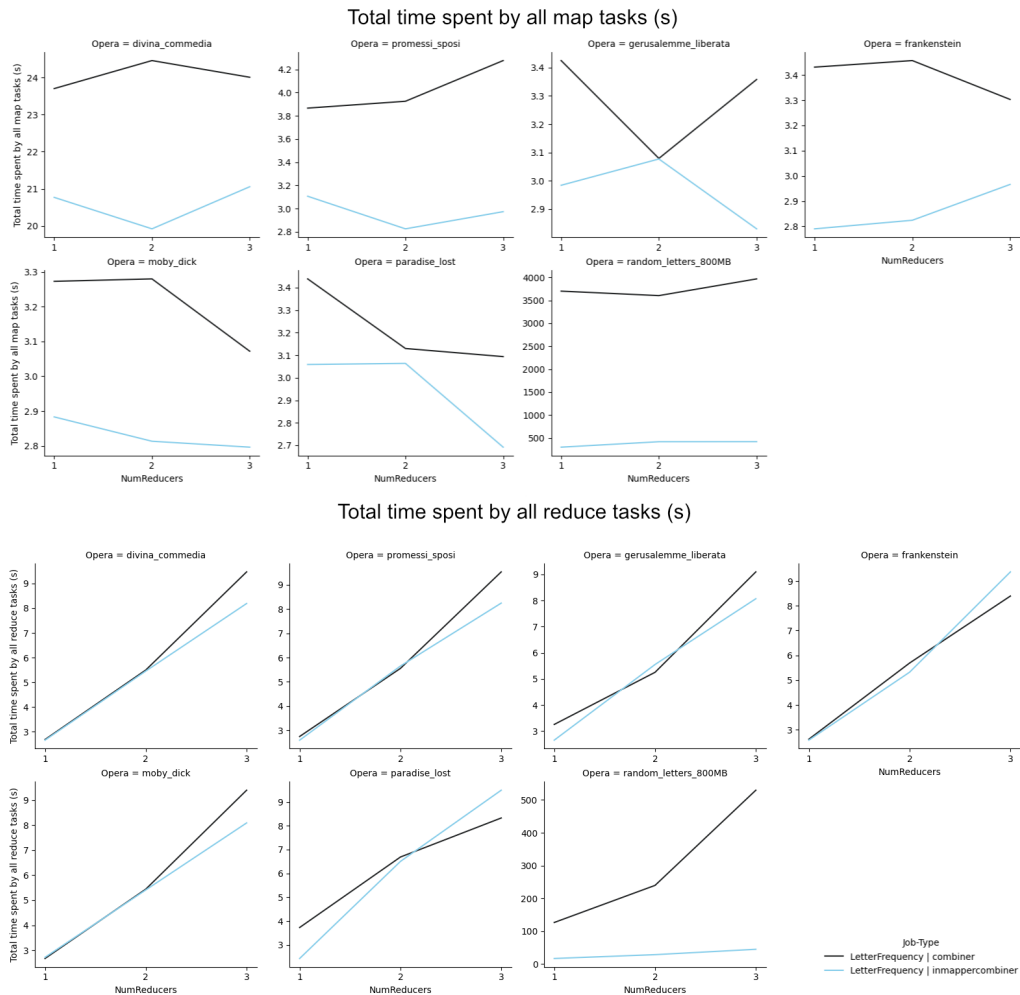


Figure 1: TotalTimeMapReduce

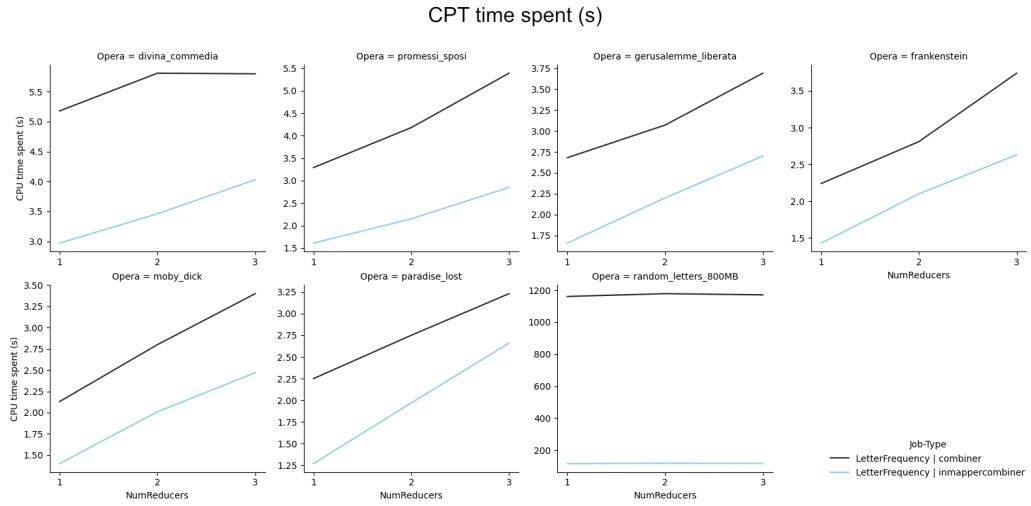


Figure 2: Cpu Time

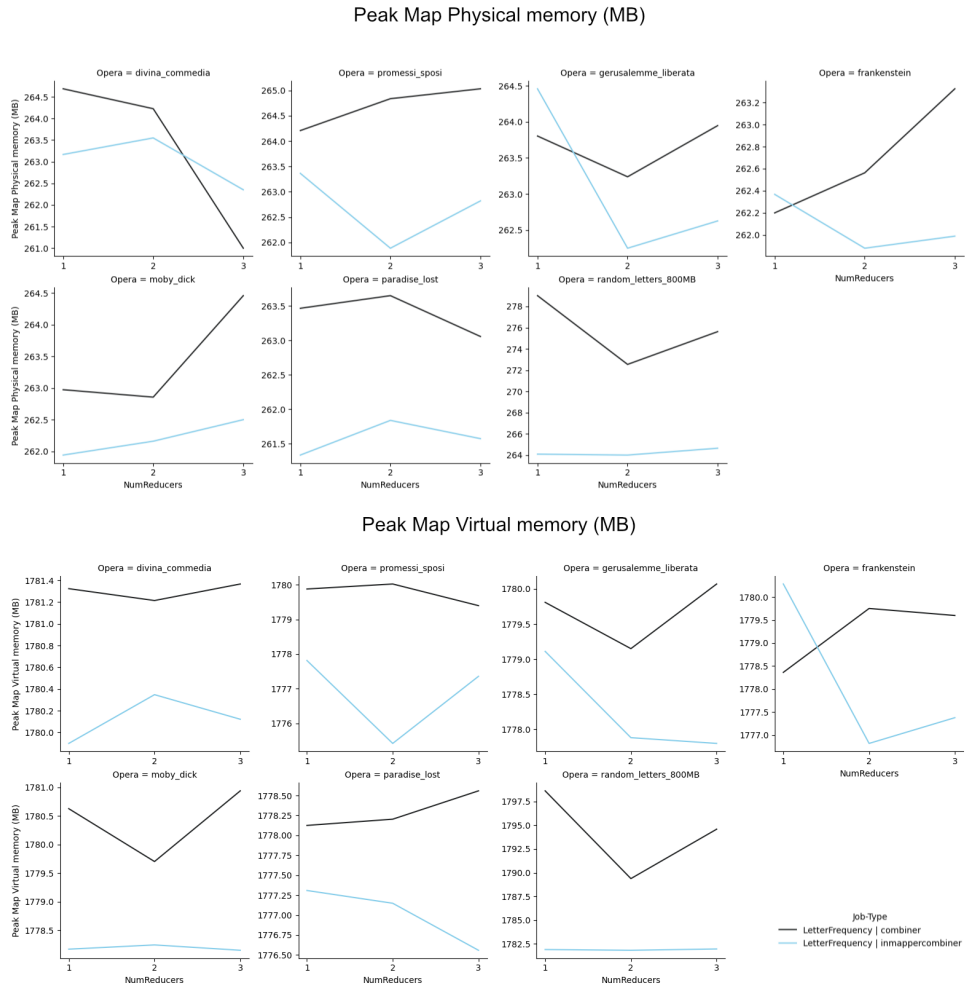


Figure 3: Peak Map Memory

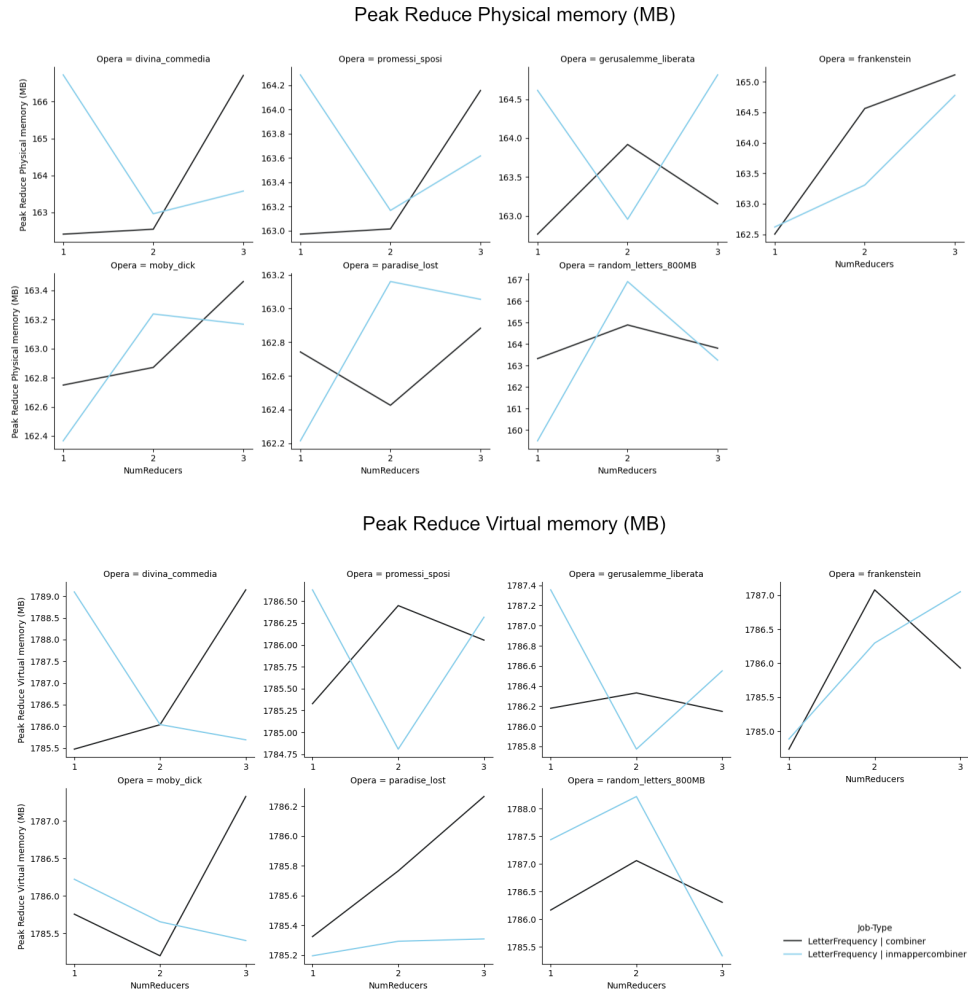


Figure 4: Peak Reduce Memory

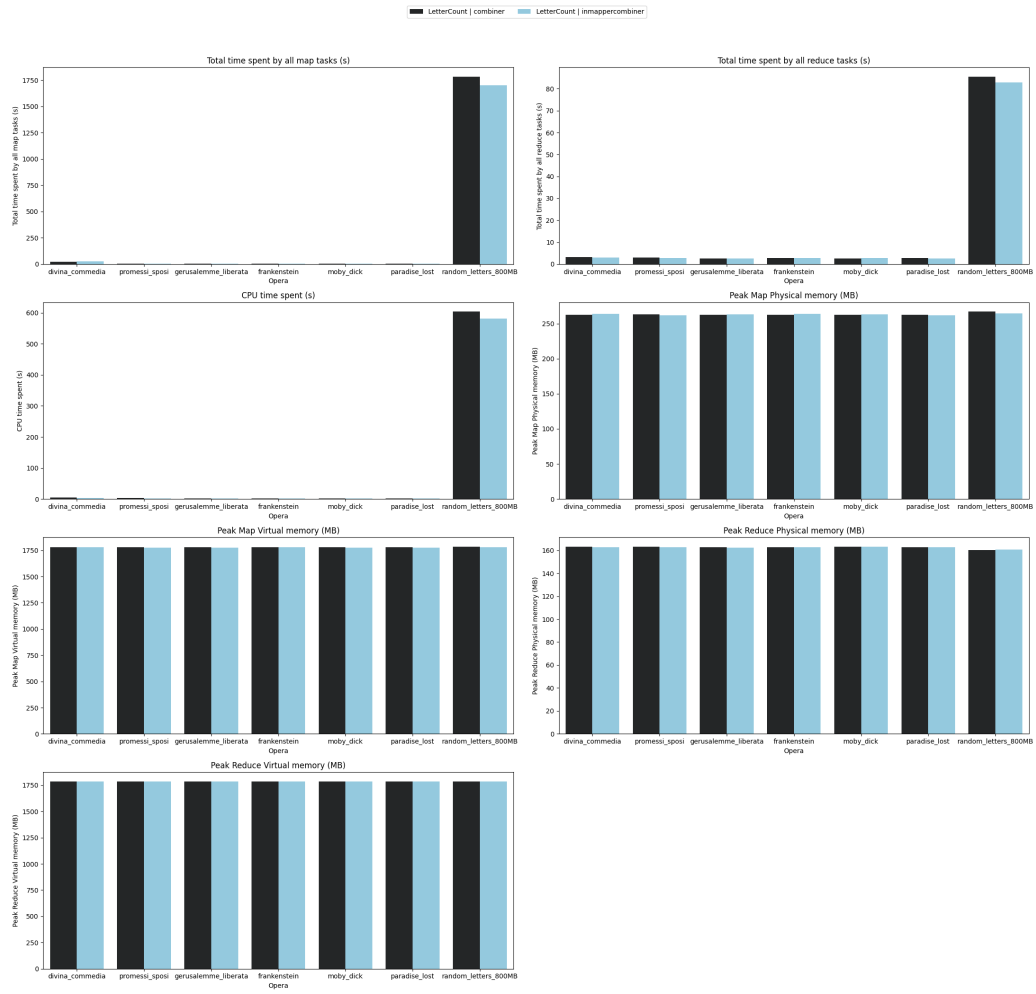


Figure 5: Performance Letter Count

Letter Frequency Analysis

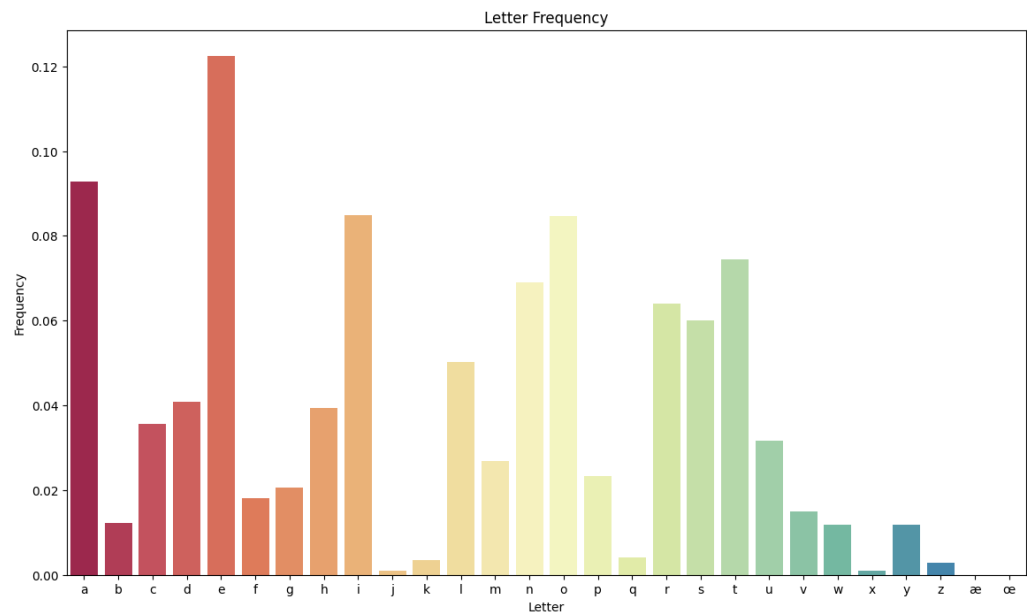


Figure 6: Letter Frequency Distribution

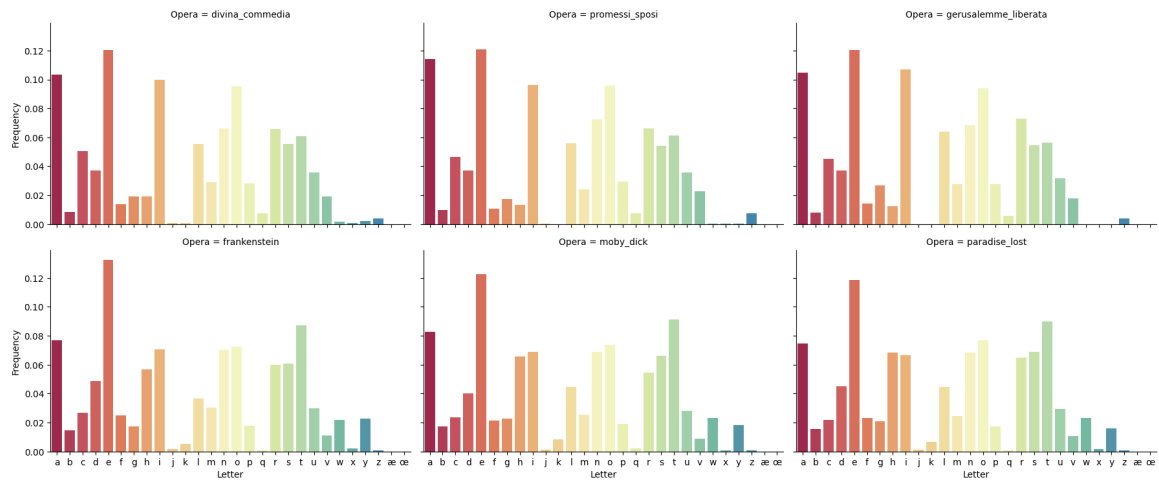


Figure 7: Letter Frequency Distribution for different Operas

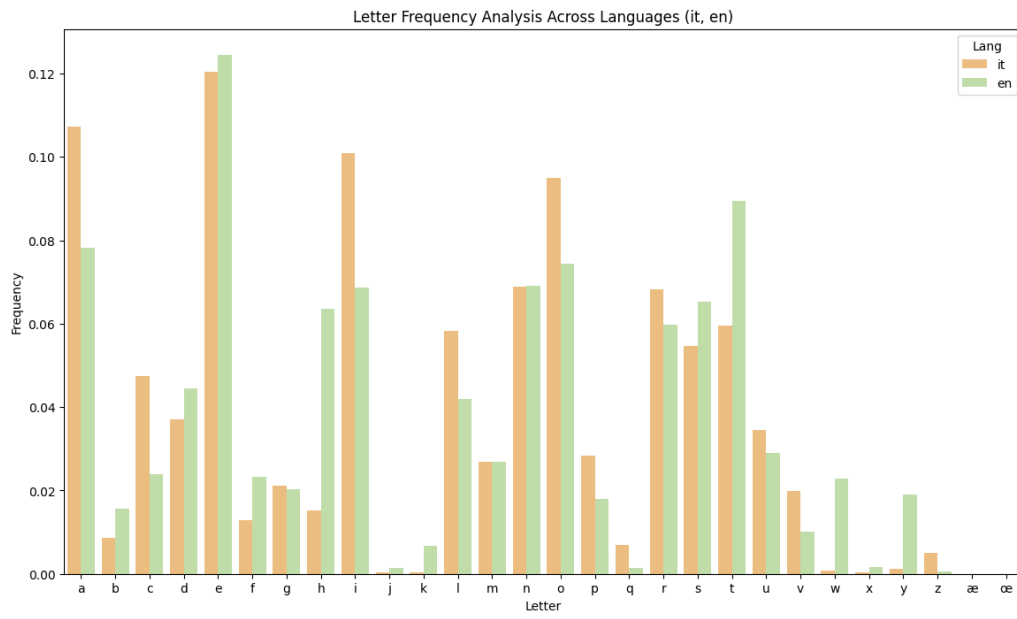


Figure 8: Letter Frequency Distribution for italian and english

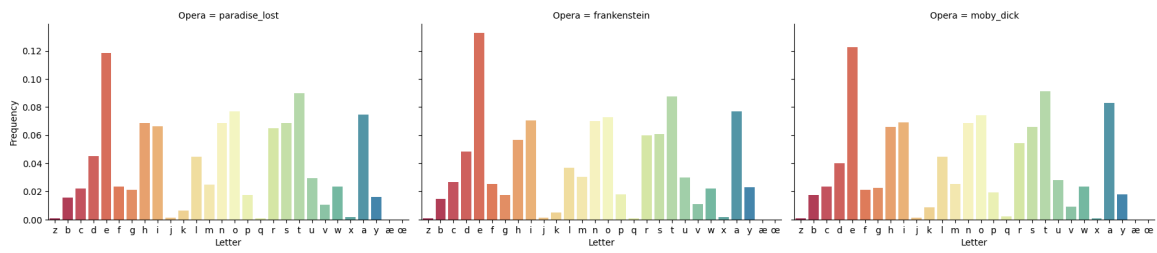


Figure 9: Letter Frequency Distribution for english

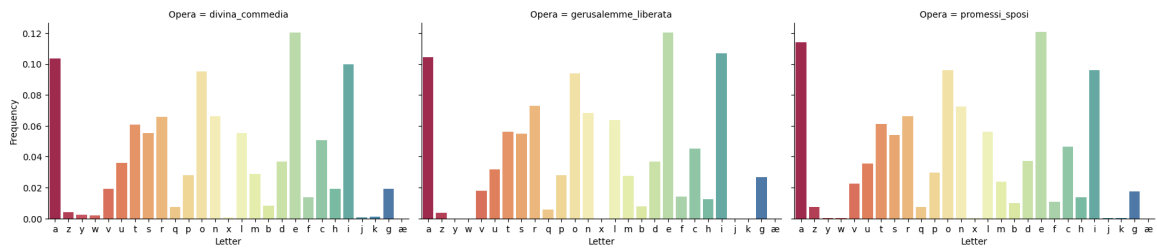


Figure 10: Letter Frequency Distribution for italian

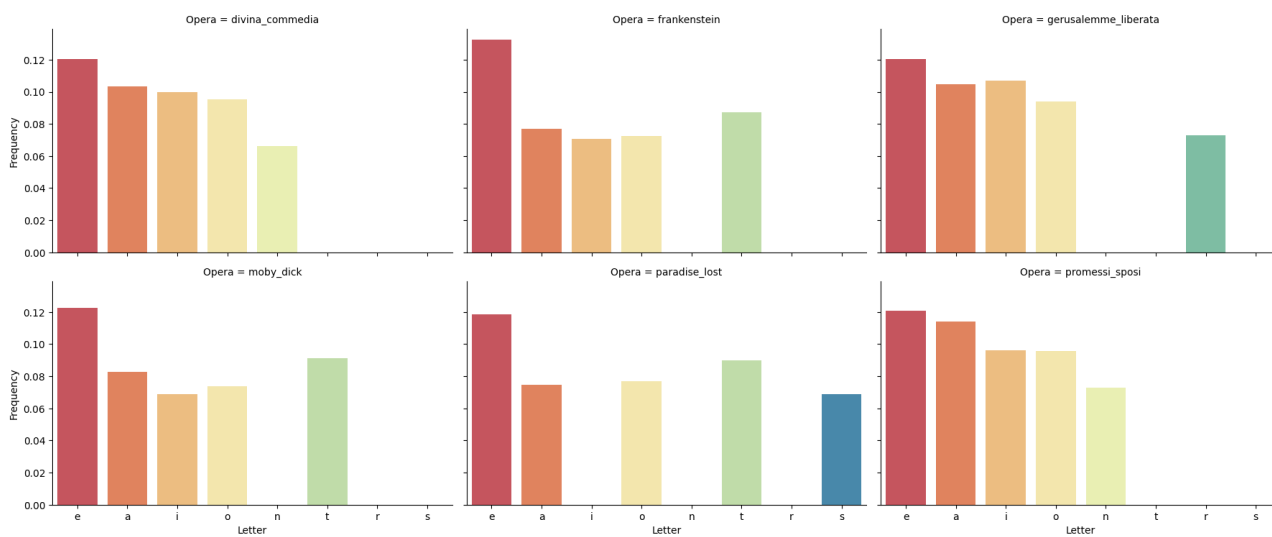


Figure 11: Top 5 letters

Conclusions

Performance

The performances of the letter frequency job and letter count one show the impact of input size and number of reducers. Moreover, the comparison between the Combiner and In-Mapper Combiner implementations highlights the better performances obtained by using the *In-Mapper Combining* design pattern.

- **Impact of Input Size:** larger input sizes, such as the 800 MB test file, grater the total processing time. This effect is evident across both the Combiner and In-Mapper Combiner implementations.
- **Number of Reducers:** it is shown that increasing the number of reducers do not lead to an overall significant improvement. This is due to the overhead of managing multiple reducers and the limited amount of data to process. In most cases, the performance is optimal with just two reducers.
- **Combiner vs In-Mapper Combiner:** the In-Mapper Combiner approach consistently outperforms the traditional Combiner method. The In-Mapper Combiner design pattern reduces the volume of intermediate data and minimizes data transfer overhead, leading to a more efficient MapReduce workflow.

Qualitative Analysis

The Letter Frequency Analysis provides significant insights into the variations and trends in letter usage across different texts and languages.

- **Letter Frequency Distribution (Figure 6):** the general letter frequency distribution graph highlights the overall trends in letter usage across the analyzed documents. Common letters like /e/, /a/, /t/, /i/, and /o/ show high frequencies, consistent with typical linguistic patterns in English and other Latin-based languages.
- **Comparison Across Different Operas (Figure 7):** when comparing letter frequency distributions across different operas, noticeable variations can be seen. These differences may reflect the unique stylistic and linguistic choices of the authors, as well as the historical and cultural contexts in which these operas were written. For instance, certain operas may exhibit a higher frequency of specific letters due to their thematic or lexical peculiarities.
- **Italian vs. English (Figures 8-10):** the comparison between Italian and English texts reveals distinct frequency patterns. Italian texts, for instance, tend to have higher frequencies of vowels such as /a/ and /o/, which are more prevalent in Italian phonology. Conversely, English texts show a higher occurrence of consonants like /t/ and /h/. These differences underscore the phonetic and orthographic characteristics inherent to each language.
- **Top 5 Letters (Figure 11):** the analysis of the 5 most frequent letters across the documents provides further insights into common linguistic elements. In English, letters like /e/, /t/, and /a/ dominate due to their essential roles in constructing common words and grammatical structures. In Italian, vowels play a more prominent role, reflecting the language's phonetic structure.

In conclusion, the qualitative analysis of letter frequency data across different texts and languages highlights both universal and language-specific patterns. These findings not only reveal underlying linguistic trends but also enhance our understanding of how language and stylistic elements vary across different cultural and historical contexts.