

Academic Year: 2023/2024

Cloud Computing project Letter frequency

Fabio Piras

Giacomo Volpi

Guillaume Quint

(Sam Sepiol)



UNIVERSITÀ DI PISA

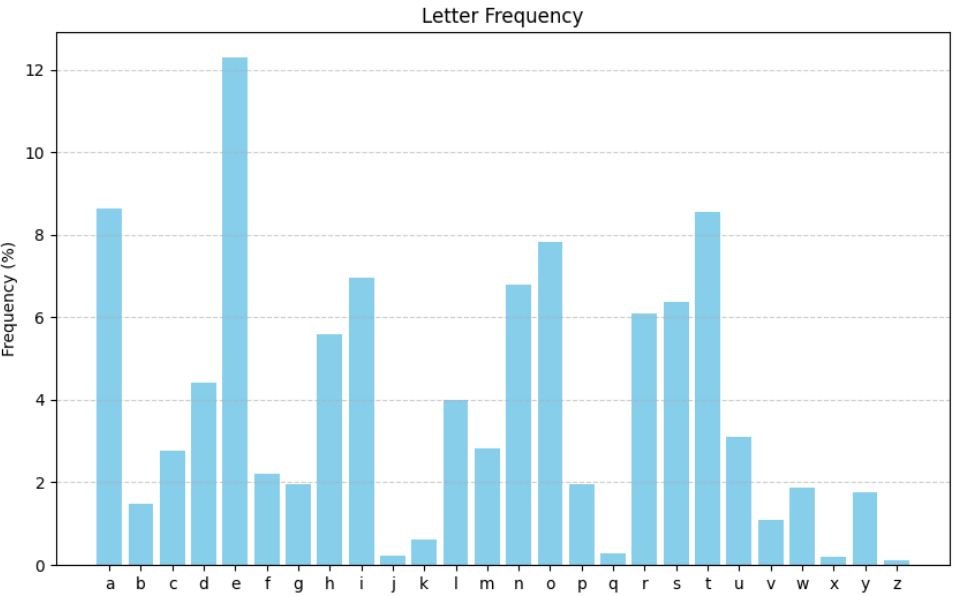
Problem and implementations

- The goal of this project was to analyze the frequency of letters in a text file using the Hadoop MapReduce.
- We implement 2 major versions of Hadoop MapReduce
- We also developed a pure Python, C and Spark implementation

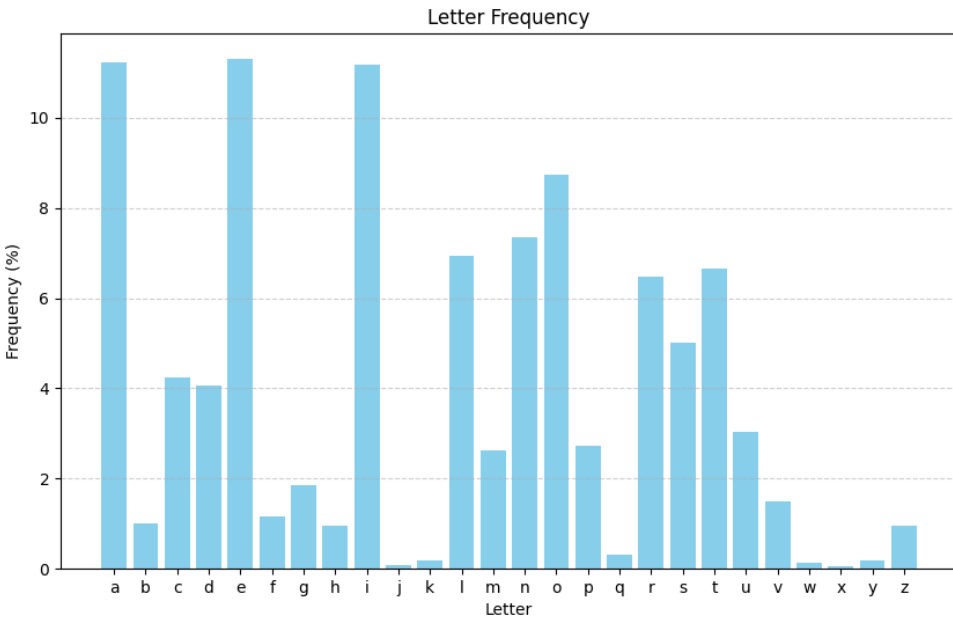
Dataset

Name	Language	Size
Project Gutenberg	English	1,2,5,10,20,50,100,654 MB
Wikipedia early June 2024 dump	Italian	3.63 GB
Base64 of OpenSSL random bytes	-	1GB and 6GB

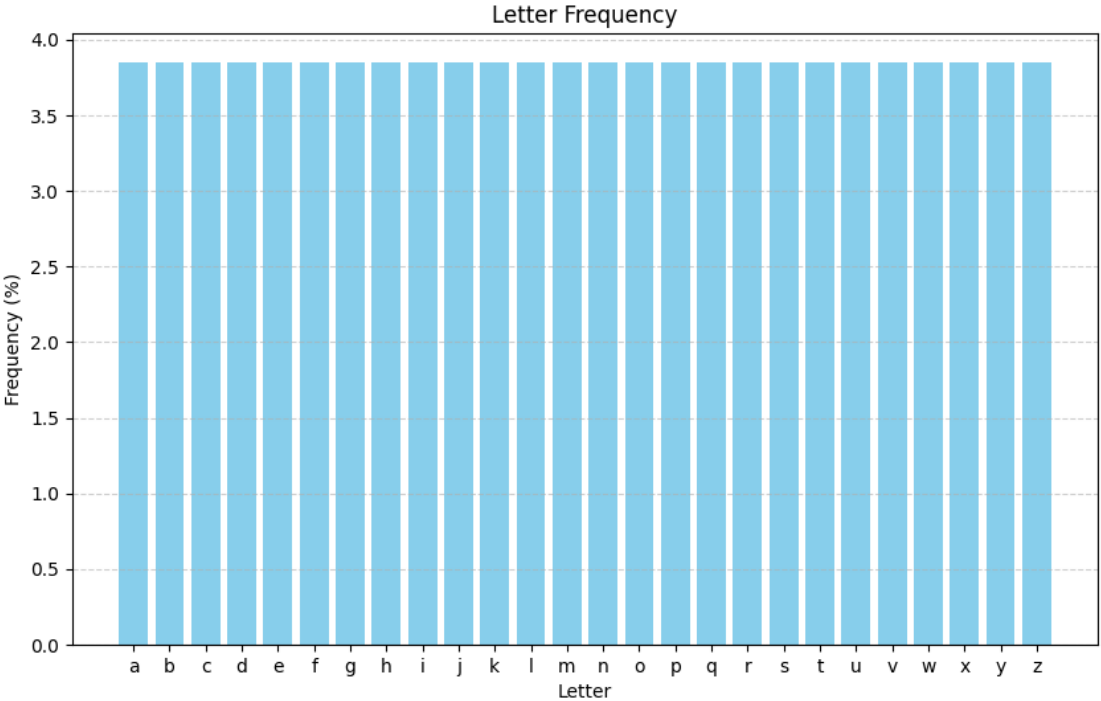
Frequency results



English



Italian



OpenSSL Rand

Hadoop MapReduce

First implementation

2 MapReduce Jobs, the first counts the total number of letters in the document, the second counts the occurrence of each letter and divides it with the output of the first job

```
1 function Mapper1(documentFragment):
2     for each character in documentFragment:
3         if character is a letter:
4             emit ("total", 1)
5
6 function Reducer1(key, values) :
7     sum = 0
8     for each value in values:
9         sum += value
10    emit (key, sum)
11
12 function Mapper2 (documentFragment):
13     for each character in documentFragment:
14         if character is a letter:
15             emit (character.lower(), 1)
16
17 function Reducer2 (letter, letter_count) :
18     total_letters = conf["total"] // This is the output from Reducer 1
19     sum = 0
20     for each value in letter_count:
21         sum += value
22     emit (letter, sum/total_letters)
```

In-Mapper combiner strategy. Only one job, use an HashMap to collect the occurrence of each letter

Hadoop MapReduce Second implementation

```
1 class MapperWithCombiner:
2     hashmap = {}
3
4     function setup():
5         hashmap.clear()
6
7     function map(documentFragment):
8         for each character in documentFragment:
9             if character is a letter:
10                 char_lower = character.lower()
11                 if char_lower in hashmap:
12                     hashmap[char_lower] += 1
13                 else:
14                     hashmap[char_lower] = 1
15                 hashmap['#'] += 1
16
17     function cleanup():
18         for each key, value in hashmap.items():
19             emit(key, [value, hashmap['#']]) // Emit letter count and total count
20
21 class Reducer
22     function reduce(letter, values):
23         if letter is '#':
24             return
25         sum = 0
26         sumTotal = 0
27         for each count, total in values:
28             sum += count
29             sumTotal += total
30         emit(key, (sum/sumTotal))
31
```

Hadoop MapReduce Optimizing Second implementation

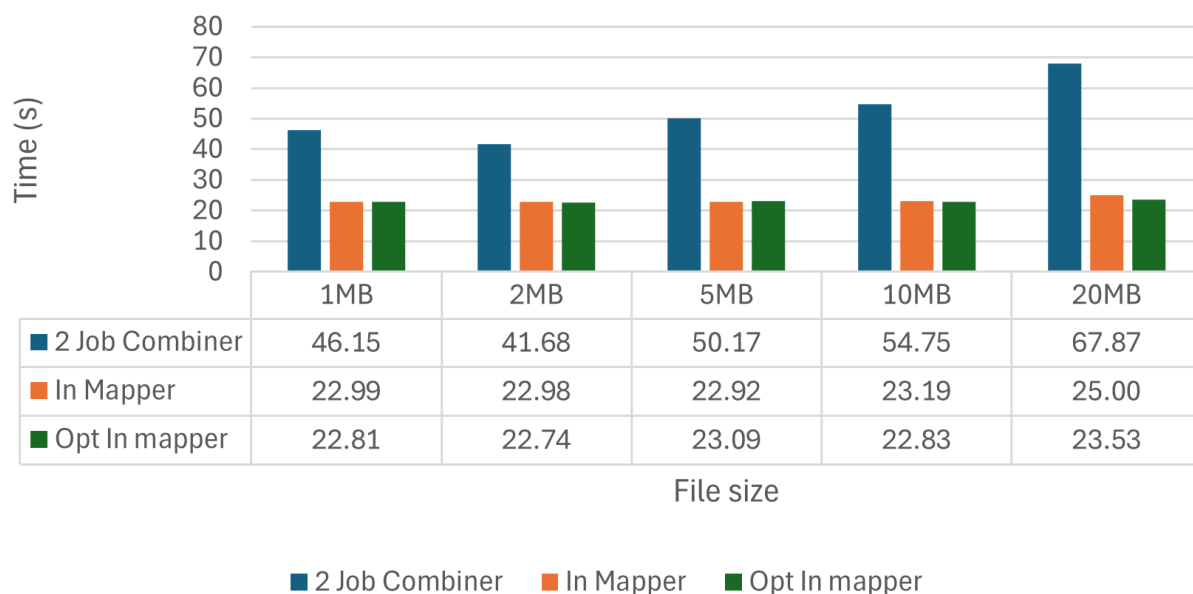
Use EnumMap instead of HashMap since keys are known a-priori

```
1 class MapperWithCombiner:
2     LetterEnum = a,b,c,d,e,f,g,h,i,j,k,l,m,n,o,p,q,r,s,t,u,v,w,x,y,z
3     enumMap = EnumMap<Letter,Long>
4     totalLetters=0
5
6     function setup():
7         for each letter in LetterEnum
8             enumMap.set(letter,0)
9
10    function map(documentFragment):
11        for each character in documentFragment:
12            if character is a letter:
13                char_lower = character.lower()
14                if char_lower in enumMap:
15                    enumMap[char_lower] += 1
16                else:
17                    enumMap[char_lower] = 1
18                totalLetters++
19
20    function cleanup():
21        for each key, value in enumMap.items():
22            emit(key, [value, totalLetters]) // Emit letter count and total count
23
24 class Reducer
25     function reduce(letter, values):|
26         sum = 0
27         sumTotal = 0
28         for each count,total in values:
29             sum += count
30             sumTotal += total
31         emit(key, (sum/sumTotal))
32
```

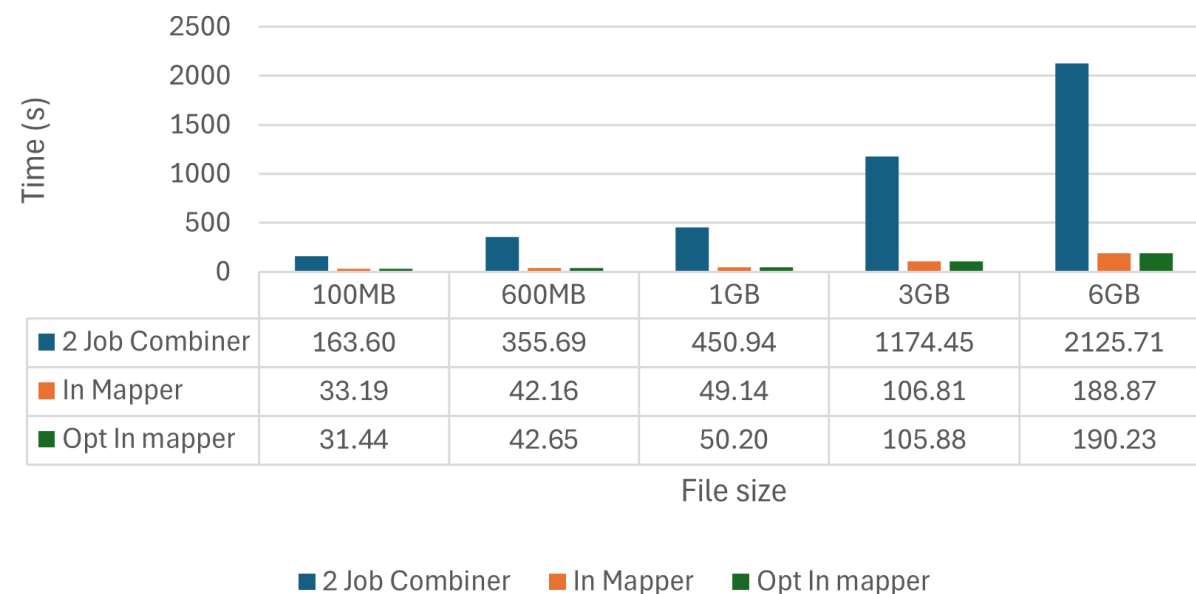
Times

- In order to obtain statistically valid results, we automated the launch of programs multiple times during the day

Small sized files Hadoop versions

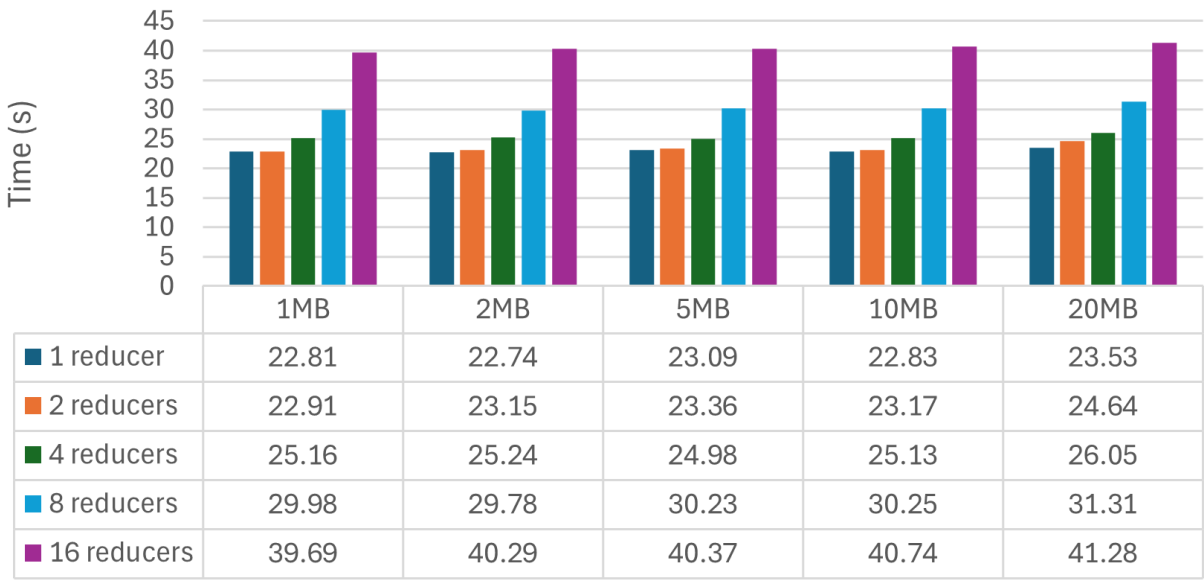


Big sized files Hadoop versions



Increasing the number of reducers

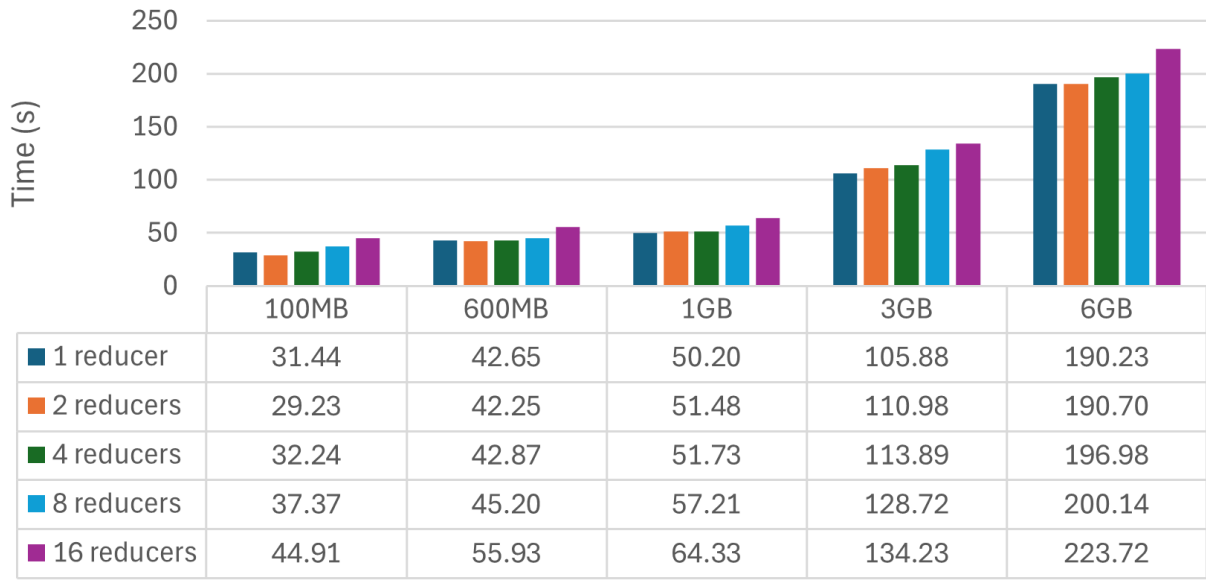
Small sized files #reducers



File size

1 reducer 2 reducers 4 reducers 8 reducers 16 reducers

Big sized files #reducers



File size

1 reducer 2 reducers 4 reducers 8 reducers 16 reducers



Alternative solutions

Python script: on a single machine


C program: on a single machine

Python spark: used on the cluster



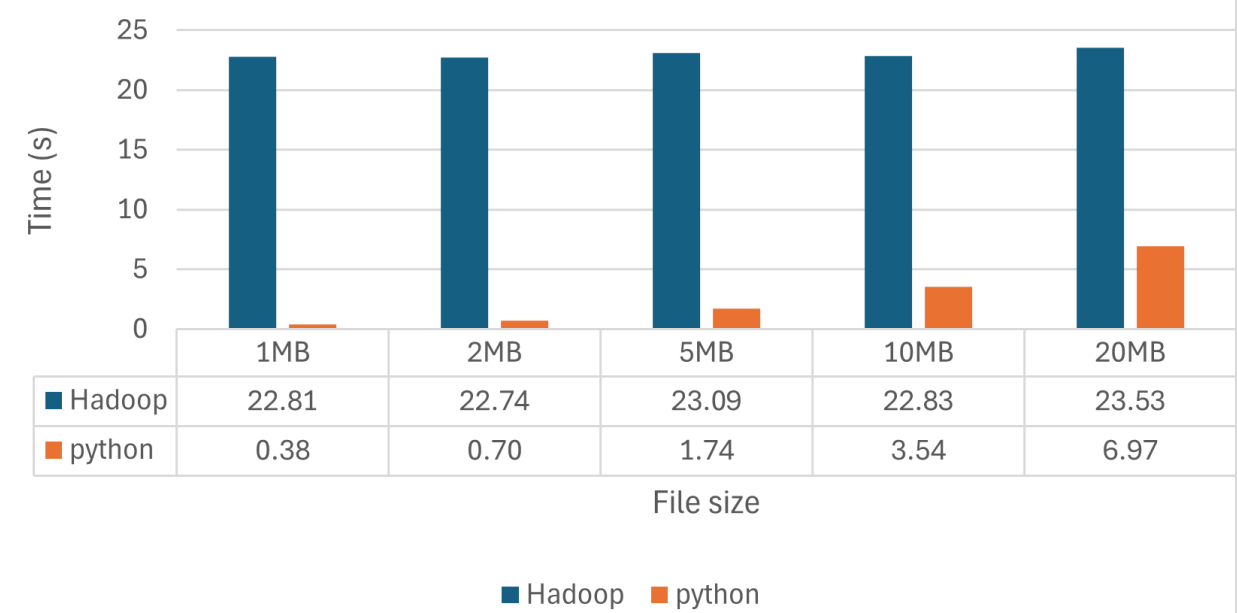
A large orange circle is positioned on the left side of the slide, partially cut off by the edge.

Pure Python implementations

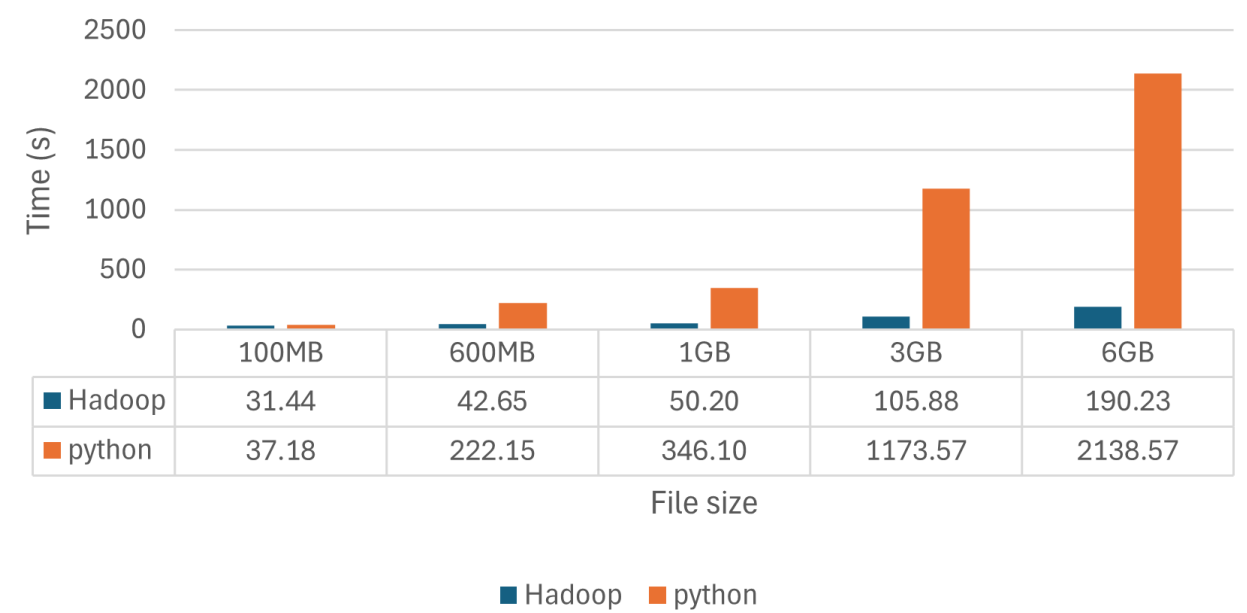
- **First implementation:** Simple, use `file.read()` and dictionary for letters count
 - **Second implementation:** Use `mmap` to map file in the virtual address space of the process, can process files bigger than RAM limits
- 
- A blue dashed line is located in the bottom right corner of the slide, consisting of several short, curved segments.

Hadoop vs Python times

Small sized files Hadoop v Python



Big sized file Hadoop v Python

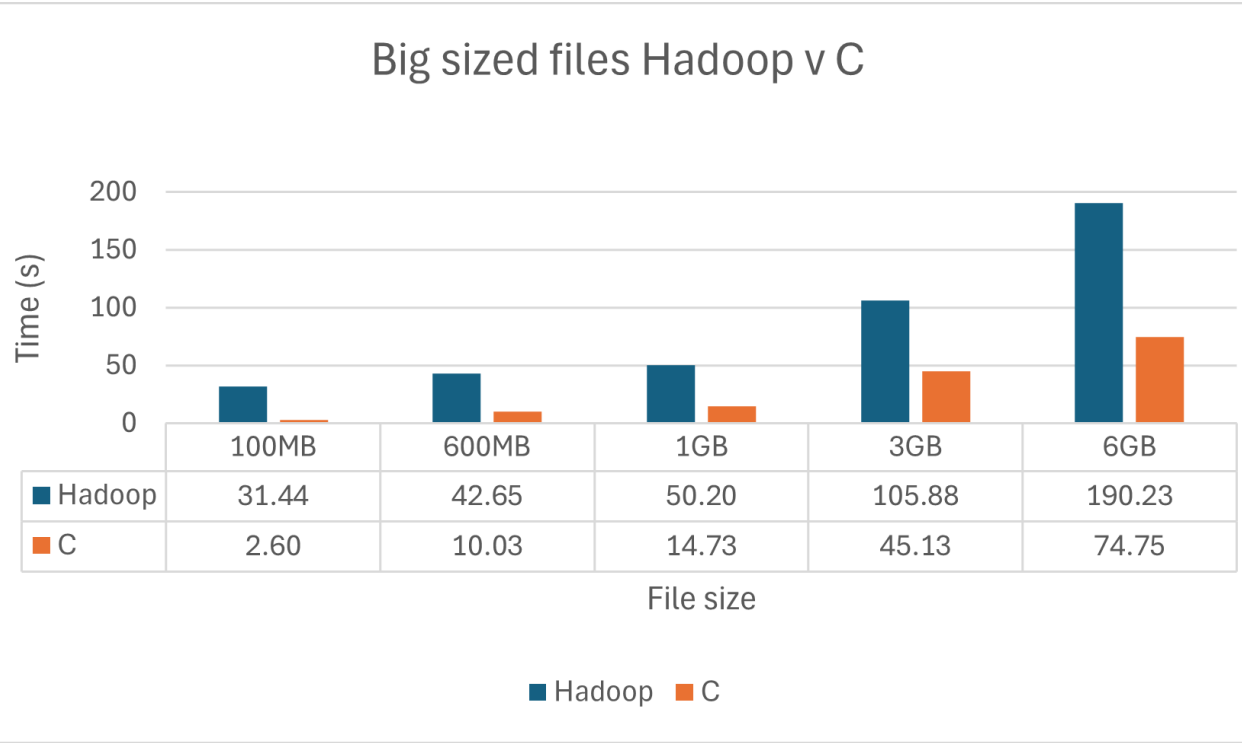
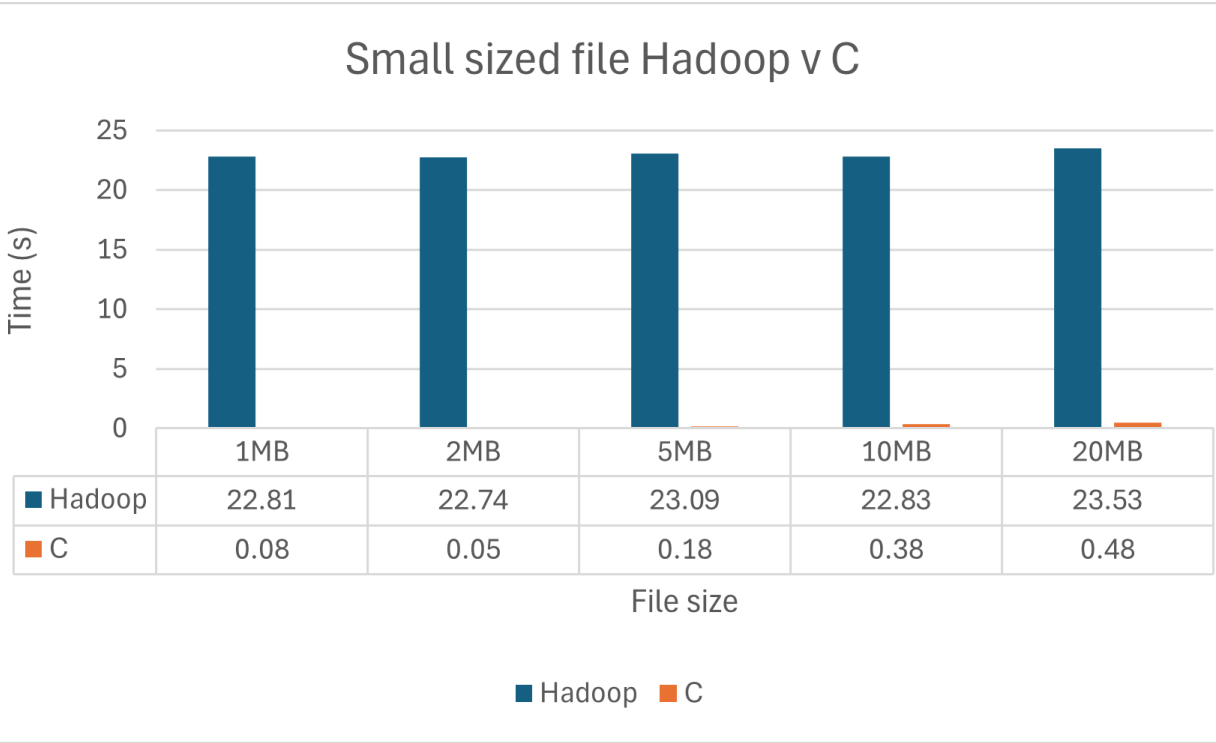


C implementations

- **First implementation:** Use multithread approach with mmap
- **Improving first implementation:**
Remove the check for ascii characters from the thread code. This results in more computation (counting non-ascii chars too) but reduces chances of branch misprediction

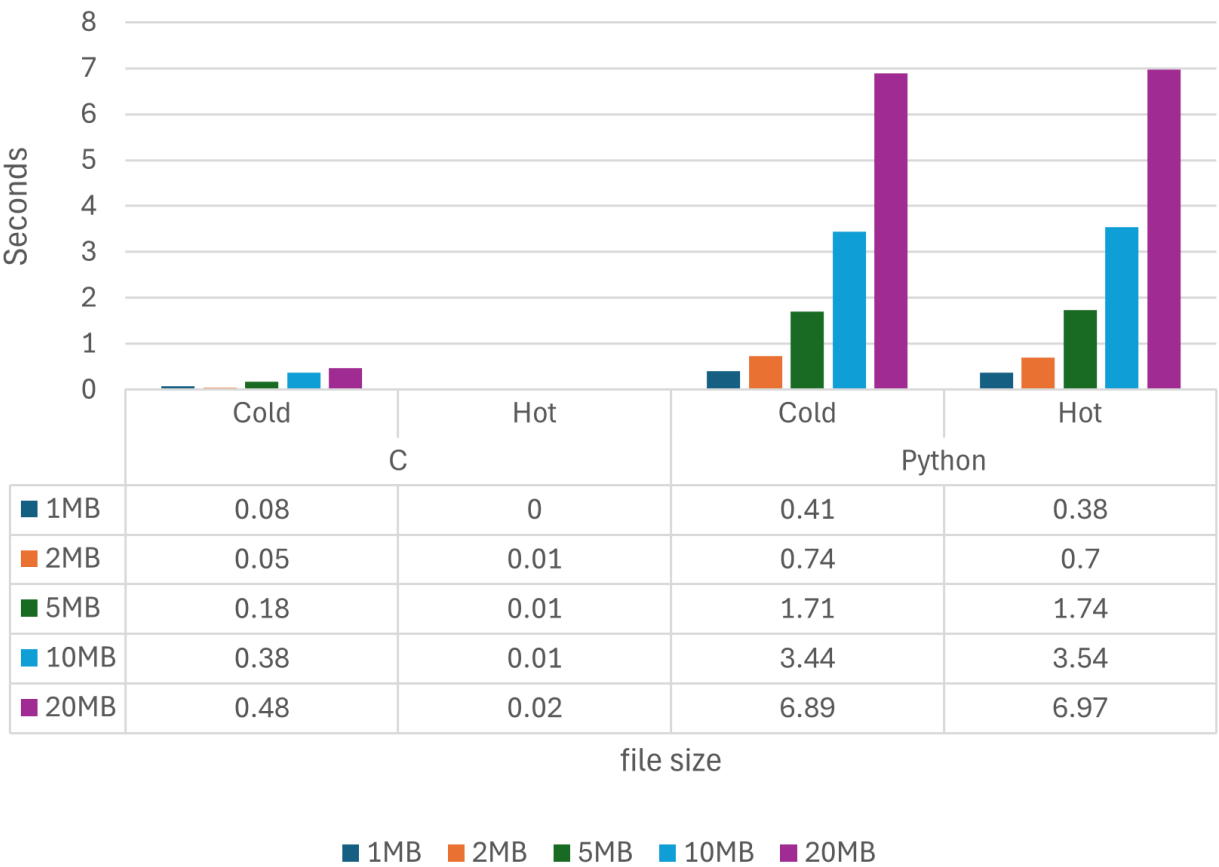
All solutions were compiled with enabled optimizations (gcc -O2 flag)

Hadoop MapReduce vs C times

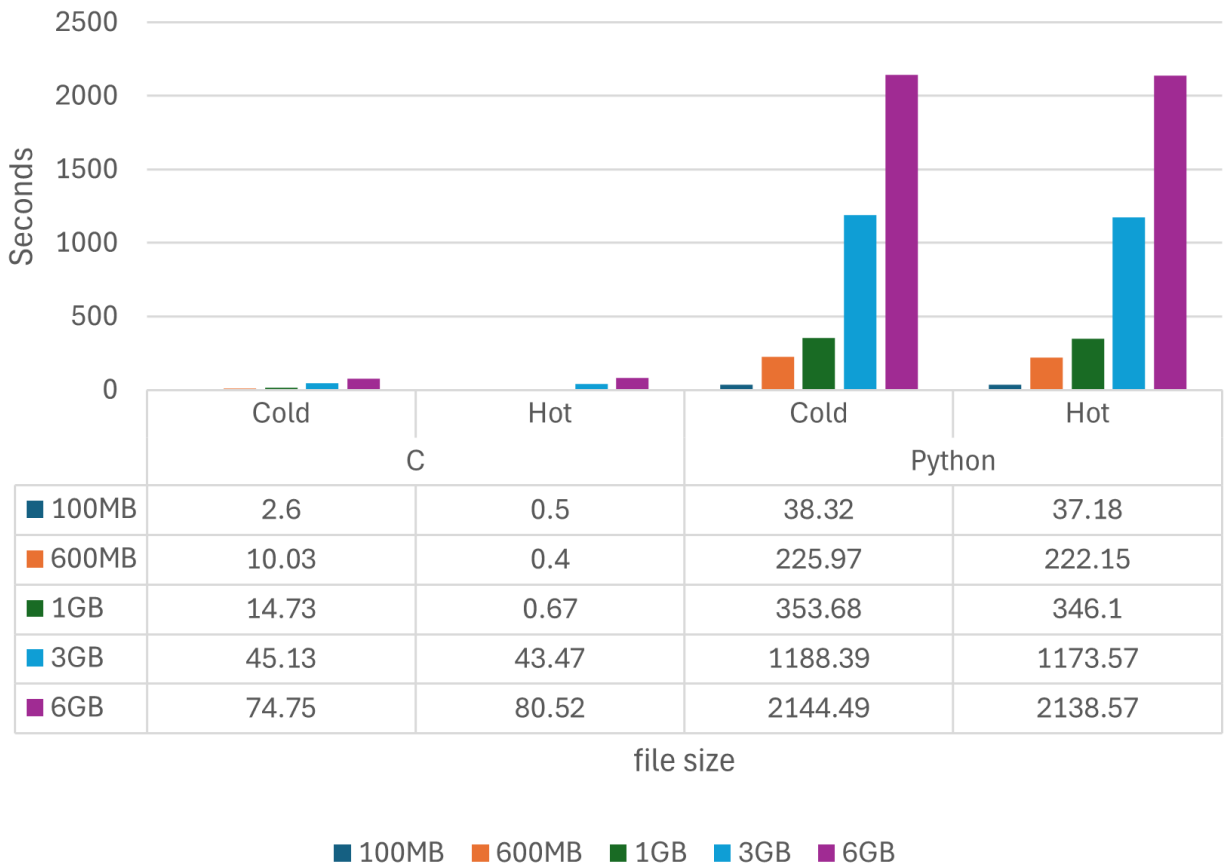


C vs Python times

Small sizes files C v Python



Big sized files C v Python



A large orange circle is positioned on the left side of the slide, partially cut off by the edge.

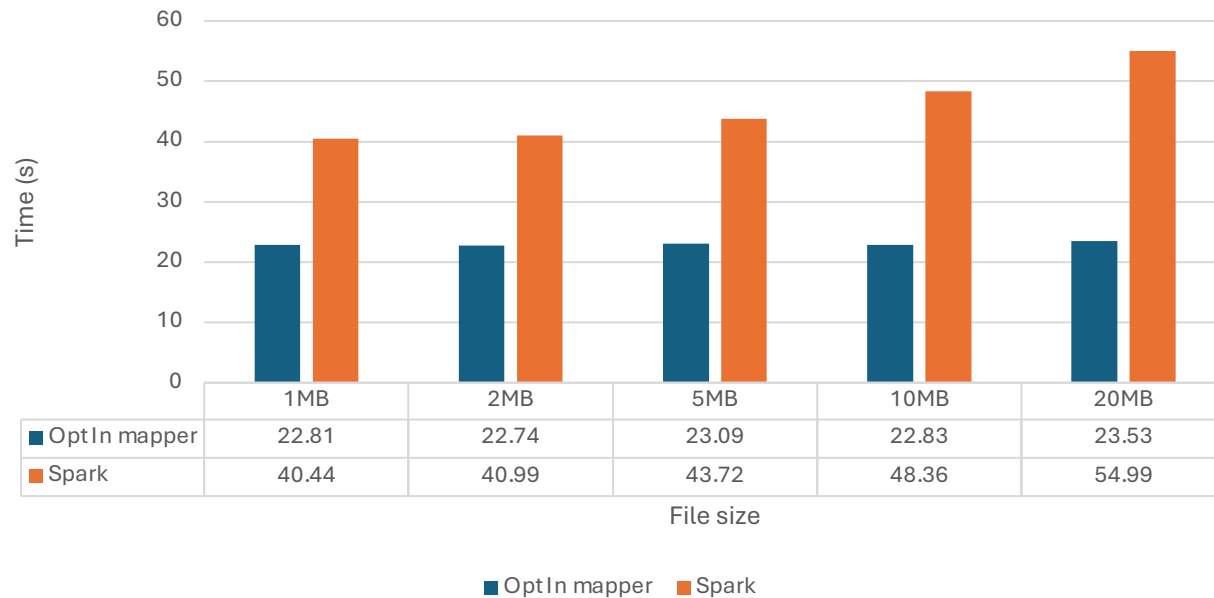
Spark implementation

- Use **RDD** and lineage tracking which allows for failure handling
- The usage of **RAM** should improve performances

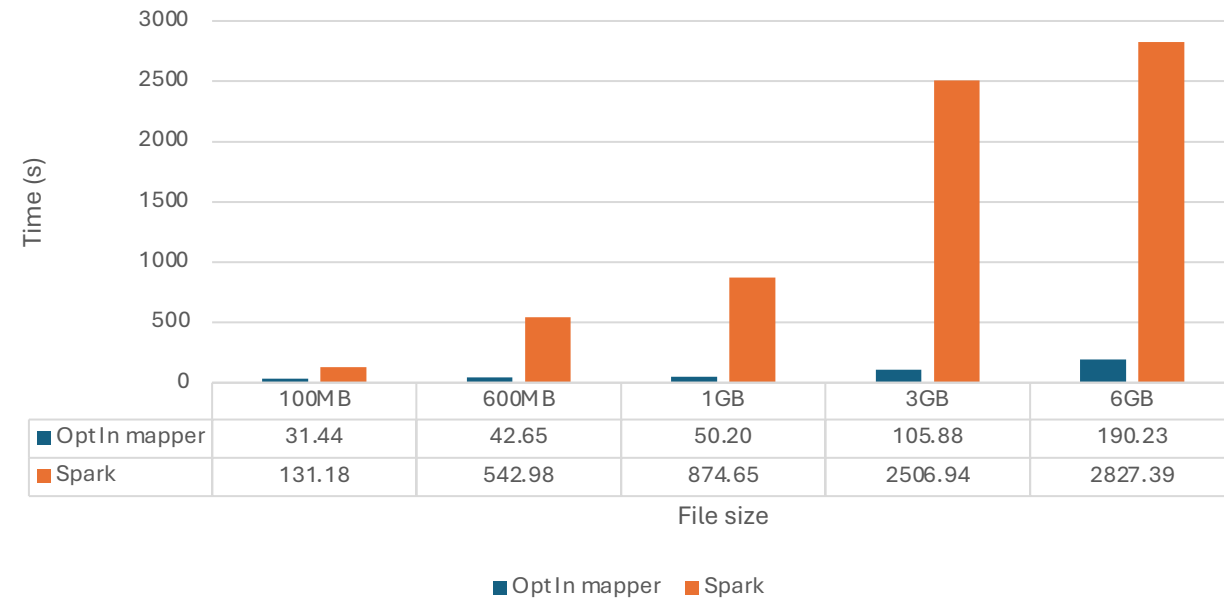


Hadoop vs Spark times

Small sized file Hadoop v Spark



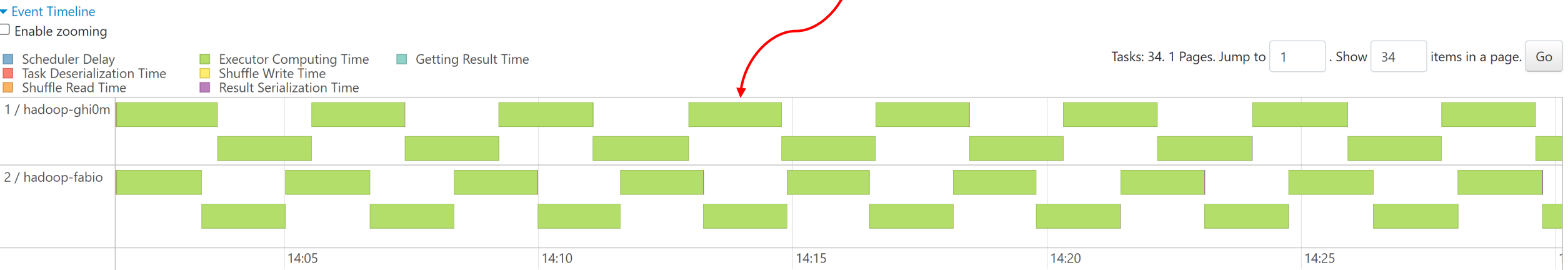
Big sized file Hadoop v Spark



On smaller files Spark can have a bigger overhead, but on larger ones it should not scale this bad, why?

Hadoop vs Spark times

~100 seconds of computing time for each block (6Gb file)



Unfortunately, the Spark deployment is CPU bounded on the cluster causing long delays

ID	User	Name	Application Type	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU Vcores	Allocated Memory MB	Reserved CPU Vcores	Reserved Memory MB	% of Queue	% of Cluster	Progress	Tracking UI	Blacklisted Nodes
application_1718967685527_0017	hadoop	letter frequency	MAPREDUCE	default	0	Fri Jun 21 16:48:27 +0200 2024	Fri Jun 21 16:48:28 +0200 2024	N/A	RUNNING	UNDEFINED	8	8	2304	0	0	50.0	50.0	<div></div>	ApplicationMaster	0

ID	User	Name	Application Type	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU Vcores	Allocated Memory MB	Reserved CPU Vcores	Reserved Memory MB	% of Queue	% of Cluster	Progress	Tracking UI	Blacklisted Nodes
application_1718967685527_0018	hadoop	LetterFrequency	SPARK	default	0	Fri Jun 21 17:13:20 +0200 2024	Fri Jun 21 17:13:20 +0200 2024	N/A	RUNNING	UNDEFINED	2	2	3072	0	0	100.0	100.0	<div></div>	ApplicationMaster	0

Key points

Hadoop MapReduce

- Only need to redefine few functions
- Distributed file system, can scale, fault tolerant
- Overhead of the JVM is visible with smaller files

Python

- Easy to implement, few lines of code
- Overhead of VM is visible w.r.t. C, MapReduce is faster with large files
- Doesn't fully exploit mmap'd files
- No failure handling, files are limited by the disk size

C

- Faster, no virtual run-time
- Hard to implement a *fast version*
- No failure handling, files limited by the disk size