

CP based Sequence Mining on the cloud using spark

Dissertation presented by
Cyril DE VOGELAERE

for obtaining the Master's degree in
Computer Science

Supervisor(s)
Pierre SCHAUS

Reader(s)
John AOGA, Guillaume Derval

Academic year 2016-2017

Contents

1	Introduction	4
2	Option for Scalability :	5
3	Solution Correctness - Verification Procedure	6
4	Performance Testing Procedure	6
4.1	Performance Testing Datasets	6
4.2	Performance Testing Procedure	7
4.2.1	Distribution Choice & Cluster Architecture	7
4.2.2	Program Parameters	7
4.2.3	Number of Partitions	8
4.2.4	Measurement Span	8
5	The PrefixSpan Algorithm	8
6	PPIC and Spark's Original Implementation	8
6.1	Spark's original Implementation	8
6.2	PPIC's Original Implementation	10
6.3	Performance Comparison	11
7	Improving the Performances	12
7.1	A First Implementation	12
7.2	Improving the Link	15
7.2.1	Adding Functionalities & Improving Middle-put Translation	15
7.3	Improving Spark	17
7.3.1	Automation of the choice between PPIC and Spark's local execution . . .	17
7.3.2	Additional Pre-processing Before the Local Execution	18
7.3.3	First/Last Position List	20
7.3.4	Priority scheduling in local execution	22
7.3.5	Specialize Spark's scalable stage for single-item problems	24
7.4	Improving PPIC	24
7.4.1	Pushing PPIC's Ideas farther	24
7.4.2	An Implementation Closer to Spark	26
8	Final Version	26
9	Scalability Tests	27
10	Performances Under Optimal Parameters	27
11	Conclusion	27
12	Annexes	27
12.1	Additional Performance Comparisons :	27

List of Figures

1	Performance comparison of Hadoop and Spark	5
2	The simple architecture used during the majority of our tests	7
3	An example of Spark's execution	10
4	Original performances of PPIC and Spark	12
5	A first implementation that makes PPIC scalable	13
6	The performance impact of adding new functionalities	17
7	Performance improvement - soft limit on the number of created sub-problem. . .	18
8	Automatic detection of item-sets type in dataset	19
9	Efficiency gain of cleaning the sequence database before any local execution . . .	20
10	Performance with first/last position lists	21
11	Spark's mapReduce	22
12	Naive priority scheduling	23
13	Performance improvement of sorting sub-problems during map stage.	25
14	PPIC with a map structure	26
15	PPIC's performances VS other specialized algorithm	27
16	Performance improvement of fixing Spark's pre-processing	29

List of Tables

1	Dataset features	6
2	Number of partitions used during performance tests	8

Abstract

TODO - Half a page to a page of content should be enough

1 Introduction

Frequent pattern mining is a widely studied problem concerned with discovering frequent subsequences in a given dataset of sequences, where each sequence is an ordered list of symbols. Those frequent pattern mining (FPM) problems can generally be further differentiated into two kinds of problems :

- Multi-item pattern mining problems, which consist in finding sequences of frequent sets of symbols. While time consuming, this type of FPM is frequently used in the industry. A practical example would be to find sequences of items often bought together, so that a promotion could be made and sales could be increased.
- Single-item pattern mining problems, which consist in finding sequences of frequent symbols. While also time consuming, it has been proven through various papers¹ that algorithms specialised for single-item pattern mining can widely outperform their multi-item counterparts on the same datasets. Motivating their common uses in applications such as web log mining and biological sequence analysis.

While many such specialized algorithms exist, constraint programming (CP) has often been proposed as a framework for frequent single-item pattern mining in recent years (### TODO : Check it has been proposed for multi-item pattern ###). The main benefit of CP based algorithms lying in their modularity, which allow the addition of various constraints with ease. Those constraints spanning over a wide range of possibilities, such as restriction on symbol occurrences, pattern length, respect of regular expressions, ...

However, it was feared that this increased flexibility came with a performance cost. Many of the developed algorithm staying uncompetitive with state-of-the-art specialized methods. Although many improvements had been made, notably by Kemmar et al², which had further extended this work by introducing one constraint (module) for both the pseudo-projection and the frequency pruning, only reasonable performances were obtained.

However, a recent paper³, shattered those fears by out-performing, not only other CP based implementation, but also state-of-the-art specialized methods (see Figure 15). This PrefixSpan based algorithm called PPIC, implemented on the open source CP solver Oscar⁴, achieve those performances by combining various ideas from pattern mining, such as the last-position lists of the LAPIN algorithm, as well as ideas from CP, such as trailing which allow faster computations and avoids unnecessary copying.

Their implementation is, however, limited in two aspects. First, it is limited to single-item pattern mining problems. Second, the developed implementation is not scalable, which is a serious flaw for the industry. As cloud computation is a must to solve large problems in a timely fashion.

Through this paper, we shall attack those two problems under various angles and prove that an efficient and scalable CP based FPM algorithms can be developed.

¹See references + TODO

²Link to Kemmar et al. paper TODO

³Link to PPIC's paper TODO + should I cite all three name ?

⁴<https://bitbucket.org/oscarlib/oscar/wiki/Home>

2 Option for Scalability :

Since the Oscala library currently doesn't allow scalable cloud computation, and that an efficient such implementation couldn't reasonably be made during the short time frame of a master thesis, we had to choose a cloud computing library to develop our scalable implementation.

Since, fortunately, FPM problems are embarrassingly parallel problem, we had the opportunity to choose from a selection of widely used open-source libraries. We, however, restricted ourselves to Scala compatible framework, as the CP library supporting PPIC was implemented in this language. If we wanted to compare performance and prove the implementation could be made scalable, our CP framework couldn't be changed. We thus considered two serious options :

1. Hadoop mapreduce
2. Spark

TODO : Should I make two subsections, one explaining Spark and its concept and one explaining Hadoop ?

Both of those libraries already disposing of a scalable PrefixSpan implementation. It was thus a matter of determining whose performances were better, and whether those implementations could be efficiently extended through CP technologies.

Fortunately for us, performance comparison had already been done in a widely recognised scientific paper on Spark's RDD[1]. The performances comparisons they obtained then, are presented in Figure 1

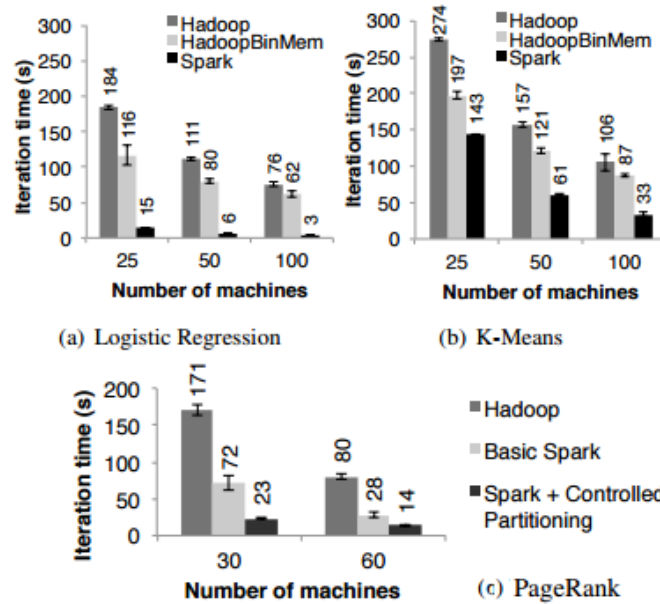


Figure 1: Performance comparison of Hadoop and Spark

As you can see, performance-wise, Spark vastly outperform Hadoop thanks to its ability to use both in memory and on disk computations. Allowing up to 100x speed-up under the right circumstances, as you can see in the logistic regression problem. According to the official website, Spark would also boast a 10x speed-up through on disk computation, but no performance benchmarks were provided.

In terms of extensions through CP technologies, we quickly realised that Hadoop would be far less practical. Although MapReduce can be used to execute the standard PrefixSpan algorithm, and could certainly be modified to introduce CP elements, Spark can support any coarse-grained transformation with its RDDs. Allowing more precise implementations where only required transformations would be made, instead of simple sequences of Map-Combine-Reduce.

We thus decided to use Spark to implement our scalable version of PPIC, as the advantages it provided far surpassed anything Hadoop could provide.

3 Solution Correctness - Verification Procedure

TODO ? (It feels like it would be too much, should I add that ?)

4 Performance Testing Procedure

To compare performances, we first need to discuss how they were obtained.

4.1 Performance Testing Datasets

For our performance tests, eight datasets were chosen for the different characteristic they displayed. The goal being to prove the efficiency of the developed algorithm in a wide range of situations. The chosen datasets and their characteristic are displayed in Table 1 :

	Dataset	#SDB	N	avg(#S)	avg(#Ns)	max(#S)	Sparsity	description
1.	BIBLE	36369	13905	21.64	17.85	100	1.18	text
	FIFA	20450	2990	36.24	34.74	100	1.19	web click stream
	Kosarak-70	69999	21144	7.98	7.98	796	1.0	web click stream
	LEVIATHAN	5834	9025	33.81	26.34	100	1.25	text
	PubMed	17237	19931	29.56	24.82	198	1.17	bio-medical text
	protein	103120	25	482.25	19.93	600	24.21	protein sequences
2.	slen1	50350	41911	13.24	13.24	60	1.0	generated dataset
	slen2	47555	62296	17.97	17.97	74	1.0	generated dataset

Table 1: Datasets features. The datasets of category 1 contain only single-item pattern, while the datasets of category 2 contain only multi-item pattern

- #SDB = number of sequences;
- N = Number of different symbols in the dataset;
- #S = Length of a sequence S;
- #Ns = Number of different symbols in a sequence S;
- Sparsity = $\frac{1}{\#SDB} * \sum \frac{\#S}{\#Ns}$

As you can see, our datasets vary largely in their characteristic, be it in their sparsity or in the size of their set of symbols. Although the selected datasets focus slightly more on single-item pattern problems, since it will be our focus through a major part of this paper, we also made sure to correctly represent multi-item problem.

4.2 Performance Testing Procedure

4.2.1 Distribution Choice & Cluster Architecture

Performance will be tested running different custom-distribution of Spark, compiled independently for each specific change implemented. The goal being to test the performances of an actual distribution containing our changes, rather than simply running the program on an existing distribution.

Although this will take us longer to obtain results, the goal behind this choice, is that our code could then easily be proposed as the standard for future distributions of Spark.

Additionally to that, unless specified otherwise, our algorithm will be executed on Spark's standalone cluster, in client deploy mode (not locally). More specifically, during our tests, we will run a simple architecture composed single master and a worker with four executor. A representation of this simple architecture is available in Figure 2.

Both the driver and the executors will also dispose of a large amount of memory (10G each) during our tests. As it serves no purpose to limit their abilities during those preliminary tests.

Later in our paper, scalability performances will also be tested in larger, memory-restricted, architecture. The architecture then used shall be specified in the same section.

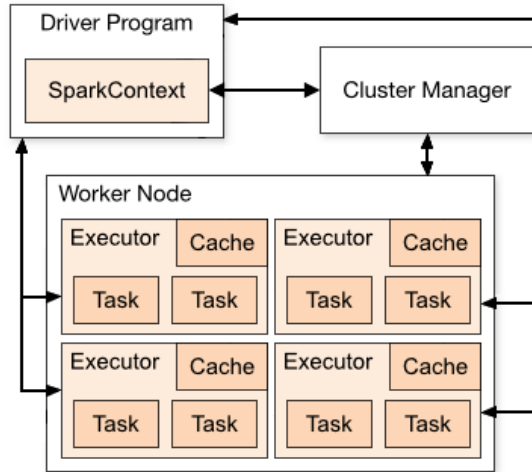


Figure 2: The simple architecture used during the majority of our tests

4.2.2 Program Parameters

At all time, unless specified otherwise, the program parameters will be kept at their default value, so that differences in parameters will never be reflected in our performances.

An exception to that rule is the 'maxPatternLength' parameter of Spark, since its default value of 10 doesn't fit our purpose. We will use `INTEGER.maxValue` instead, so that every solution pattern in our datasets may always be outputted. This will also be applied to future functionality we may add.

Later in our paper, we will also compare the performance of our implementation under default and optimal parameters.

4.2.3 Number of Partitions

As seen previously, Spark's performance may vary greatly depending on the number of partitions created from the Input dataset (see figure 1). Should too much be created, the algorithm will lose considerable amount of time switching its execution between partitions. Should it be too little, scalability will be affected as the partition will be too big. Spark's shuffler may even crash.

We will thus keep a constant number of partitions for each dataset. So that our performances can remain completely comparable. The number of partitions was specifically chosen for each dataset, as to guarantee good performances and the absence of unexpected crash from the shuffler.

The number of partitions created for each dataset are displayed in Table 2.

Dataset	File size (Ko)	Number of partitions
BIBLE	3065	250
FIFA	2594	300
Kosarak-70	2166	250
LEVIATHAN	713	100
PubMed	1646	200
protein	126046	5000
slen1	5400	500
slen2	6896	500

Table 2: Number of partitions used during performance tests

4.2.4 Measurement Span

For our performance tests, we will measure not only the running time of our algorithm, but also its pre-processing and dataset loading performance. The aim being to develop a new implementation that entirely surpasses the old one.

5 The PrefixSpan Algorithm

TODO ? Is it necessary

6 PPIC and Spark's Original Implementation

In this section, we will explain the original implementation of PPIC and Spark, and then compare their performance.

6.1 Spark's original Implementation

Spark's algorithm can be separated in four stages :

1. **Pre-processing** : The goal of this stage is to replace each item of the sequence database by an unique ID, to separate itemSets by a zero delimiter, and to clean the database from unfrequent items.
For example, the sequence $\langle (ABD)(ABC)A \rangle$ will become 0120123010, assuming only items A, B and C are frequent in the sequence database.
2. **Scalable execution** : The core of the algorithm. Its execution consist in extending large prefixes through a three sub-stage process, starting from the empty prefix.

- (a) First, a large prefix is projected on the database. If there is no prefix to project, the scalable execution comes to an end and the solutions will be returned.
 - (b) Then, from the set of supporting sequences, we discover symbols that can extend the current prefix. If no such extension exists, we try the next large prefix.
 - (c) Finally, for each possible symbol extension, we determine how long further expanding this prefix may take, by calculating the projected database size. Then, depending on the calculated projected size and of the value of a user defined parameter, we either further extend this prefix using another iteration of scalable execution, or store it for use in the local execution stage.
3. **Local execution :** The local execution is completely similar in its implementation to the scalable execution. Its only use is in improving performance by calculating all extensions from a Prefix at once, instead of doing so over multiple iteration. This stage is only launched once all large prefixes have been extended sufficiently. Depending on the parameters inputted by the user, this stage may even be skipped.
 4. **Post-processing :** During the post-processing step, we translate back the unique IDs into the item they each represented. Then we send back the collected and retranslated results to the user.

Since this implementation is specialized for multi-item pattern, one particular improvement has been created. During the prefix projection, for each item in the projected prefix, the algorithm will detect which itemsets in the sequence comply with what has already been projected and can still be extended in some way. Storing such items positions in a 'partial projection' list. That way, if we are projecting an itemset containing multiple symbol, until the end of that itemset, the following item projection will know where to search possible extensions, and extend from there without having to search the whole sequence. Also, should we be computing a single-item pattern mining problem, no partial start would be created and kept in memory.

If the itemset end, the full remains of the sequence, from the earliest partial position recorded, will have to be searched. During this search, new potential partial projection will be recorded, and the old one will be discarded.

During the prefix generation phase, the current partial projection will also be used to find extensions of the current itemset. The remainder of the database will also be searched, but only for extension that starts new itemsets.

An example of a fully scalable execution from this algorithm can be found in figure 3

TODO - Should I add some pseudo code ?

NB: During our first analysis of Spark's pre-processing stage, we noticed a small inefficiency in the cleaning of the database's sequences. When multiple itemsets were fully cleaned of their item, the algorithm had a tendency of creating sequences of zero delimiters, as the algorithm still delimited empty itemsets.

Although the internal representation was still correct and results weren't modified, those trailing zeroes substantially slowed the algorithm down. The performance improvement of this small correction is reflected in the annexes, figure 16.

Later, this small correction was proposed to Spark's community and quickly accepted into the default implementation. In the remains of this paper, we will thus consider this

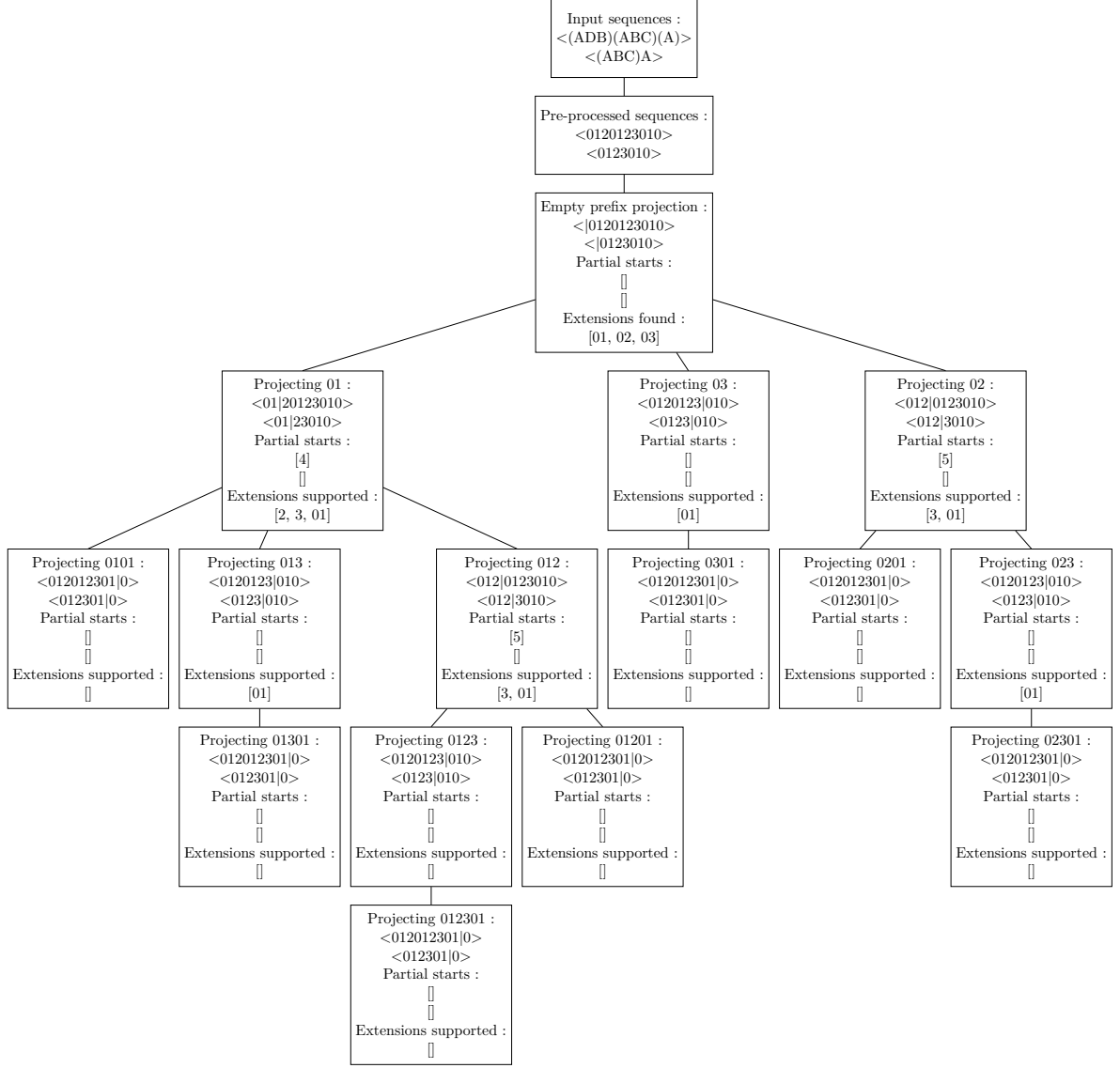


Figure 3: A simple execution of Spark's algorithm.

The solutions are the projected prefixes (except the empty prefix)

corrected version which has been approved by the community as the original algorithm for Spark, and compare our performance improvement with this corrected version as the basis.

6.2 PPIC's Original Implementation

PPIC's execution can be separated in two stages :

1. **Pre-processing** : In this first stage, we first clean the received sequences from unfrequent items, renaming them into unique ID's. Three matrices are then build from the sequence database :
 - (a) The 'first-position' matrix : A $\#SDB \times N$ sized matrix allowing $O(1)$ jumps to the first occurrence of a given item.
 - (b) The 'last-position' matrix : A $\#SDB \times N$ sized matrix allowing $O(1)$ check for the presence of a given item in the remains of a sequence.
 - (c) The 'interesting-position' matrix : A matrix with the same size as the original sequence database, but whose content are changed from the items forming those sequences, to

the positions of the next 'interesting' item. That is the next position where an item last appears in a sequence.

Although, at first glance, this matrix may seem redundant with the last-position matrix, its purpose appears when one realise that to achieve the same goal a whole columns of the last-position matrix would have to be checked. Similarly, keeping only this matrix would also be less efficient, since there would be no way to efficiently check if an item is present in the remains of a sequence. Both matrices are thus needed to achieve the greatest efficiency gain.

It is interesting to note that, since PPIC has only been implemented for single-item pattern mining problem. No delimiter are insert into the cleaned sequence database, making its representation more compact.

The pre-processing stage will also take care of adding multiple constraint, depending on the wishes of the user.

2. **Execution** : Once the pre-processing is finished, the algorithm will truly start to run. Using the three matrices, prefixes will then be projected and extended more efficiently. Thanks to trailing, memory consumption during the execution will also be minimal, since there will be no need to keep multiple copies of the data in storage.

Each time a valid solution is found. That is, a solution that satisfy all constraints. The execution will be momentarily interrupted. The solution will then be translated back to the items name corresponding the recorded ID's.

Once the solver determines that no further solution can be found under the specified constraints. The execution will terminate, and all resulting pattern will be returned.

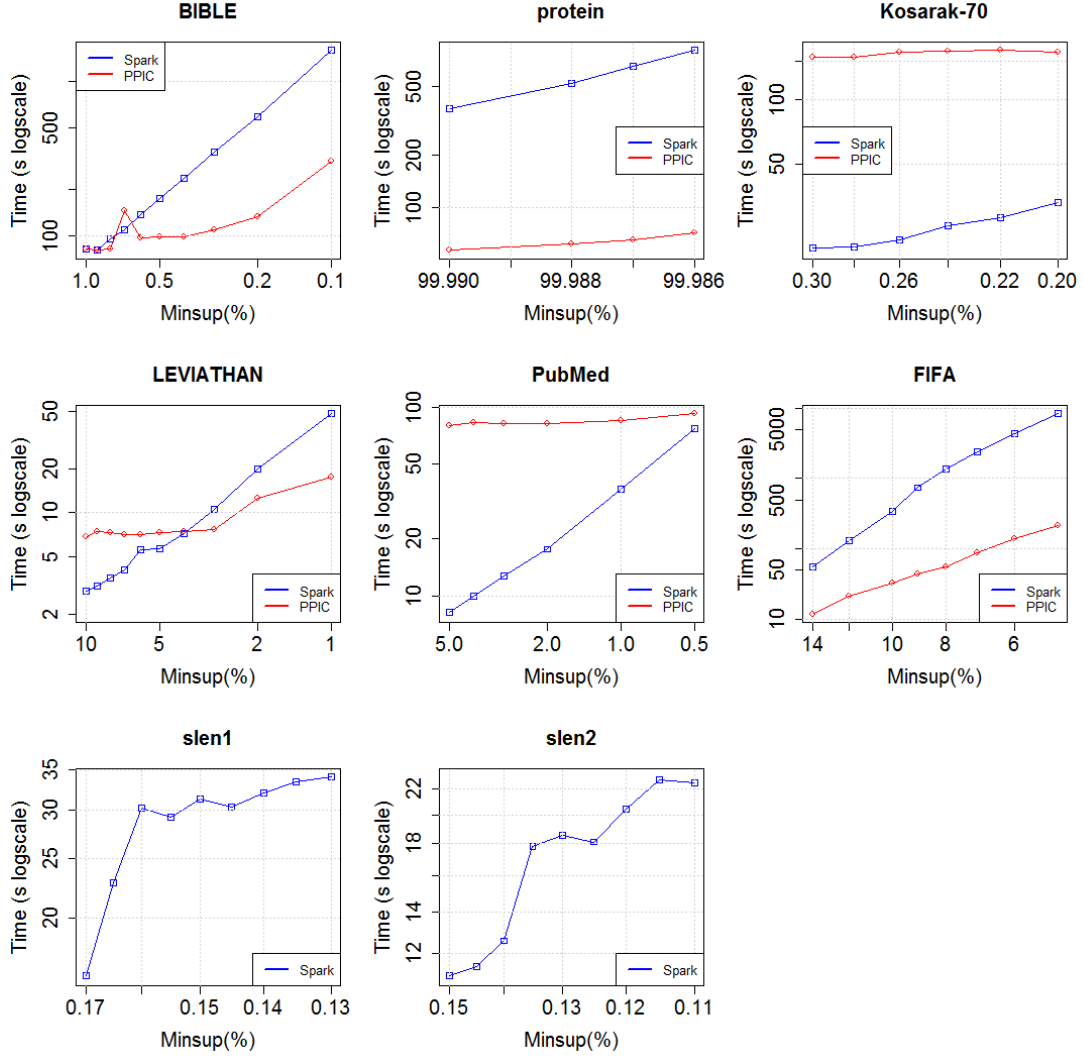
TODO - Speek more about how the execution is happening
TODO - Should I add some pseudo code ?

6.3 Performance Comparison

Now that we understand both algorithm better, the time has come to compare their performance under similar resources. Since PPIC's is neither scalable nor concurrent, we will exceptionally use and even simpler architecture for our test on Spark by restraining our worker to a single executor.

As you can see in Figure 4, given the same resources, PPIC greatly overcomes Spark on most datasets. But, although it completely outperforms Spark's implementation on sparser dataset like protein, it fails to do so on denser datasets like Kosarak-70. Especially when the minimal number of supporting sequence is large, and much of the dataset can be cleaned during the pre-processing stage.

We can also notice that PPIC's performance are far more stable, and that running additional tests for a smaller amount of minimal support may even show on inversion on denser datasets like PubMed.



⚠ Test realised on a simplified architecture ⚠

PPIC : 1 thread with 10G memory

Spark : 1 driver + 1 worker with 1 executor. 10G memory each

Figure 4: Original performances of PPIC and Spark

7 Improving the Performances

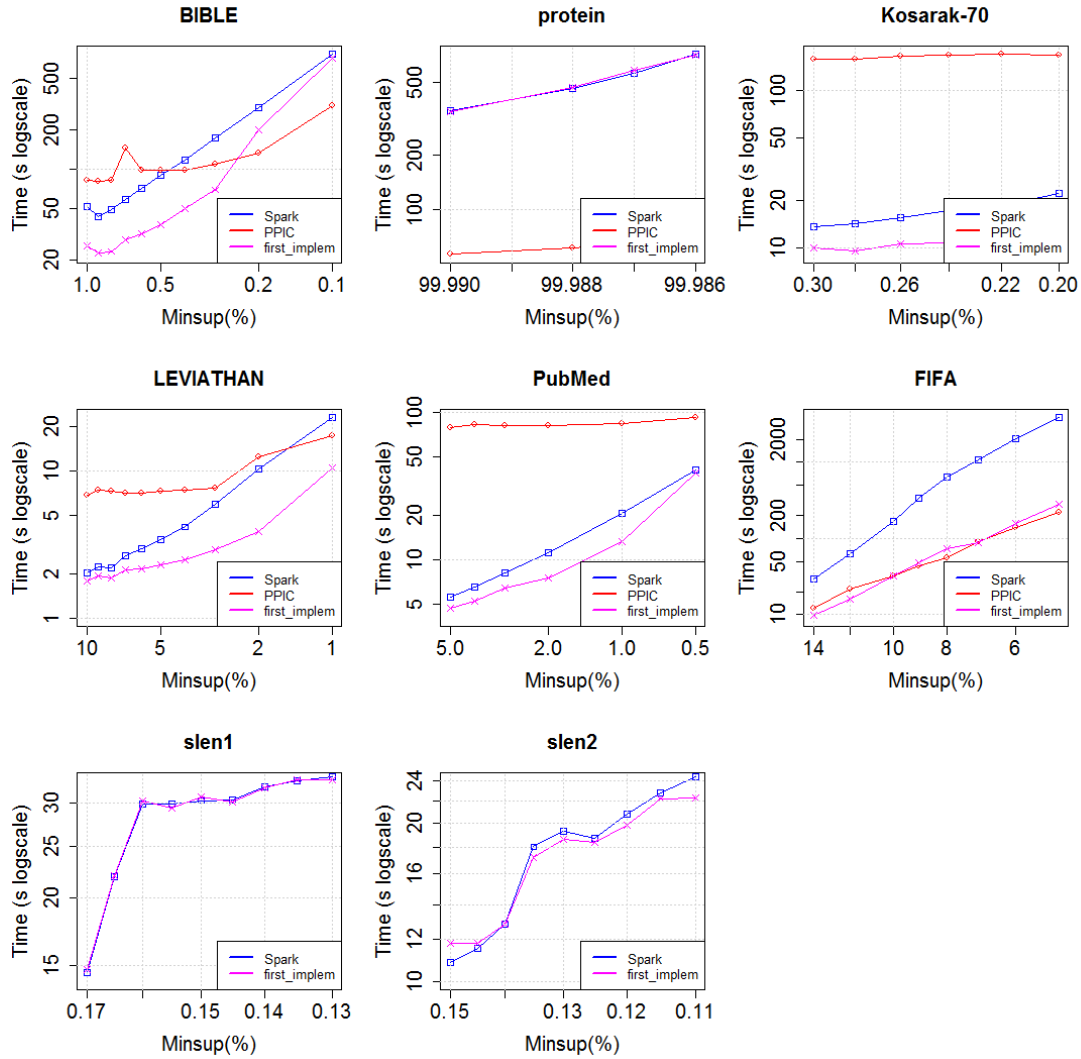
7.1 A First Implementation

Now that we observed the original performances of both algorithms, the time has come to design a first scalable implementation.

As mentioned earlier in this paper, Spark's algorithm was composed of two different execution stages. We thus analysed each stage independently to understand how CP technique could improve their implementation. Although the scalable could hardly be modified into an efficient implementation, since it would require us to maintain duplicates of CPvariables, which take more space to store due to their specificities. We realised the local execution stage could be significantly improved.

In fact, the entire local execution could be easily replaced by a CP based algorithm, for single-item pattern mining problems where PPIC is applicable, we could even use a nearly identical implementation. Only the pre-processing needed to be modified as to fit Spark’s middle-put.

We thus quickly implemented this change, also adding a simple boolean to specify whether PPIC could be used on the input dataset, and measured the difference in performance displayed in figure 5. To this graph, we also added PPIC’s performance as we judged this new implementation was closely linked, and that disposing of a comparison may be interesting, please remember PPIC’s performance were measured with lower resources than our other performances on this graph.



⚠ PPIC’s performance were still measured using a single thread and 10G of memory ⚠
 Other performances represented in this graph were measured using the 4-executor test architecture previously described.

Figure 5: A first implementation that makes PPIC scalable

As you can see, this first implementation which simple merges the two algorithm is already much more efficient on single-item pattern problems than Spark and PPIC original implementation were. At the exception of protein, which was a dataset on which PPIC thrived, which may appear disconcerting at first glance.

But, in truth, those performance can be perfectly explained. They result from three major factors :

1. First, the pre-processing of the original PPIC algorithm wasn't implemented efficiently, probably since only execution time were compared on PPIC's original paper. PPIC's true efficiency was thus brought down, while our new implementation wasn't, since this pre-processing was rewritten from scratch to fit Spark's middle-put. This time, with efficiency in mind. Explaining partly the differences between the red and pink lines.
2. Second, as explained earlier, three matrices are built from the original sequence database before the beginning of PPIC's execution. However, in PPIC's original implementation, those matrices could be very large. As their size depends not only on the number of sequences, but also on the number of unique item present in the input database.

However, our first-implementation surprisingly solved this problem. Since only the interesting parts of the sequence database are kept before launching the local execution, many items and sequences which appeared in the initial problem don't appear during the local execution. For each sub-problem treated in the local execution, the constructed matrices will thus be much smaller in size. Which explain the large majority of our implementation's efficiency gain.

3. Finally, our third major factor is an inefficiency that comes as a by-product of scalability. It appear mostly on smaller datasets composed of a small number of repeated symbols, such as our protein dataset. In those cases, since the various sub-problem created before our local execution are actually very similar to each other, our new implementation will lose an important amount of time recreating the three input matrices before launching computations on those nearly identical sequence databases. While the original PPIC implementation would create those matrices once and use them through the reminder of the execution.

Sadly, this problem cannot be dynamically fixed without seriously affecting scalability or efficiency, as it would require us to compare the projection of multiple prefixes. This means that, either we would have to project the prefixes multiple time to obtain the results of those comparison, either we would project it once and keep multiple version of the database during comparison, which would be a disaster for memory consumption and scalability.

Fortunately, although a complete dynamic fix of the problem is unpractical for the implementation of a scalable and efficient solution. It is possible to give users the possibility to reduce or even negate the effects of this by-product inefficiency, through giving them control over the amount of sub-problem created. Since the less such sub-problems, the less those matrices will need to be recalculated.

In fact, Spark's implementation already contains a way to control the number of sub-problems created, thanks to the 'maxLocalProjDBSize' parameter. Stopping further sub-problem creation from problems that are already below the inputted parameter value in size. A more direct and precise way to control the amount of generated problems wouldn't be superfluous, but its addition will be studied later in this paper.

Also, we can notice that multi-item performance haven't changed much, if at all. Which is normal since the modification made only concerned single-item problems.

TODO : Should I put an emphasis on the important part below by putting it in its own section ?

Now that we understand where those performances come from, we need find how to improve them further. In that endeavour, we identified three options that need to be studied:

1. Improve the link between Spark’s middle-put and PPIC’s input. The easiest option, but also the most promising, since both algorithm have been separately optimised in terms of performances.
2. Improve Spark’s performances by incorporating more idea’s from pattern mining. Which would improve both the scalable execution and the original local execution component of Spark. Thus improving performances on both single and multi-item pattern problems.
3. Developing a new version of PPIC which can efficiently be applied to multi-item pattern. Making CP usable for every local execution opportunity, and improving performances.

Additionally to those performance improving options, we also decided to prove that PPIC’s modularity can be conserved through the addition of more functionalities.

7.2 Improving the Link

7.2.1 Adding Functionalities & Improving Middle-put Translation

From the first implementation presented previously, we set out to implement new functionalities and improve the link between Spark’s middle-put and PPIC’s input.

First, we implemented four new functionalities :

1. Unbounded max pattern length : Although Spark’s implementation already disposed of a way to control the maximum length of a pattern, no special value existed to allow unlimited max pattern length. We thus added this minor functionality.
2. Min pattern length : We added a functionality to specify the minimal length a pattern should have before being considered solution.
3. Limit on the maximal number of item per item-set : This functionality was added so that user could better control the outputted results. Supposing an hypothetical business would like to find all sequences of item bought in pair in a dataset, it would have no need for solutions where item-sets are larger than two. We can thus stop searching for further item-set extensions once the limit has already been reached, and de-facto, improve our algorithm performances for returning those specific solution.
4. Soft limit on the number of sub-problems created : By default inactivated, this parameters is enabling far better performance on protein like datasets where projected databases tend to be rather similar in sub-problems. Its implementation is a simple check at the beginning of each iteration of the scalable algorithm. If the number of sub-problem created is larger than the user inputted value, the implementation will forcefully put an end to the scalable stage, and switch to a local execution on each worker.

As you can see in Figure 7, performance improvement are only observed on the protein, Kosarak and slen2 datasets, but the increase in performance is significant. The problem being that the loss of performance on other datasets is generally even more significant. This loss in performance comes from large differences in sub-problems sizes appearing due to the forced local execution. Since the largest problems tend to be created and processed last, toward the end of the execution, only a few executor have problems to work on while the rests stays idle. Moreover, those remaining problems take a very long time to compute,

greatly slowing down this parameter performances.

As is, this parameter should be used with extreme care. However, we will consecrate a section to studying possible improvement of its performances later in this paper. [TODO link to that section.](#)

While these additions did not have any measurable impact on performances at their default value, they still allowed better control of the search space. Proving a certain level of modularity was easily kept from the original CP based implementation. In all the following graph, unless specified otherwise, those parameters will be kept to their default value, as to make sure they do not interfere with the observed efficiency gain/loss.

We then looked into improving our performances, adding two new functionalities.

First, during its pre-processing, Spark cleans unfrequent items from the database, finding frequent items in the process. In the original implementation, those frequent items are then discarded, and pattern searches are launched starting from an empty prefix. We thus decided to try keeping those frequent items, and passing them as starting prefixes to the various worker. Hereby removing the first and longest iteration of the algorithm. The performance obtained by implementing this change and the preceding 4 functionalities are reflected in figure 6.

As you can see, very small gain in time are generally made. At the exception of the Kosarak-70 and slen-2 datasets, which posses a lot of item. We believe those small loss of performance are due to the delay needed to transfer the already calculated frequent items being larger than the time needed to actually recalculate those items.

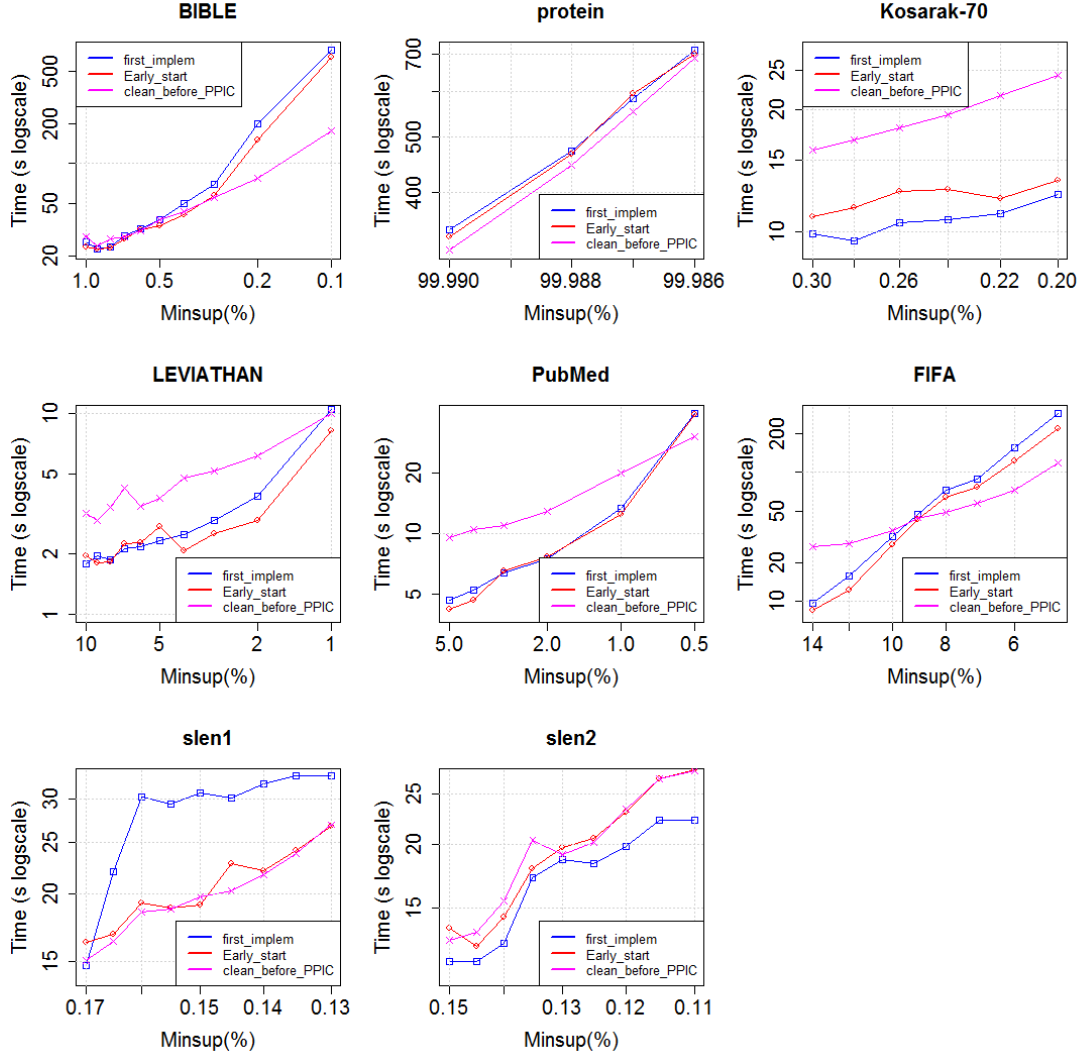
Despite the inefficiency on those two datasets, we decided to keep this change, as the efficiency gain can be relatively significant on larger datasets like Fifa or Bible.

We then implemented a second change, an additional pre-processing stage before the local execution of PPIC, so that unfrequent items would once again be detected and cleaned. The performance obtained, which can also be seen on figure 6, displayed great performance improvement on harder execution. This additional feature allowing a 3x performance speed-up on the Bible dataset and a 2x increase on Fifa. Lesser performance gain were also observed on protein and the latter stages of PubMed.

However, we can observe worsening performances in smaller datasets and in executions with larger minsup values. Theses loss however, pales in comparison to the performance gain observed. Although at first glance they may seem large on log-scale graph, they do not even exceed ten millisecond.

The reason behind those performances is fairly simple. The larger the dataset and the lower the minimal amount of support, the longer prefixes will need to be before being executed locally. The larger those prefixes are before local execution, the more items can be cleaned from the sequences, hastening the process of searching those very sequences for pattern. Finally, if less items are present in our sequences during our local execution of PPIC, the three matrices needed for PPIC's execution can be created faster. Thus, hastening performance another fold.

However, cleaning sequences comes with a variable cost depending on the database size. If, much like Kosarak, the database is very big but only contains small uncleanable sequences, performances will be worsened. Of course, since this change is only applied before PPIC, the performance measured on our two slen datasets weren't affected.



TODO -

Should I put another version of the graph without log-scale. The improvements seem downgraded due to that scale here.

Figure 6: The performance impact of adding new functionalities

In the end, since such large increase in performances were observed, and since so many other improvements needed to be tested, **we decided to use the performance obtained here as reference for the remainder of this paper**. To allow better comparison of performances gain/losses, all further improvement will thus be added to this version, and tested separately. Finally, by the end of this paper, a final version containing all changes providing performance improvements will be compiled into a single algorithm. Of whose the performances will be tested.

7.3 Improving Spark

In this section, we will discuss the improvement we tried to bring to Spark, and the results we obtained. For each improvement, we will first explain its nature and the trade-off's its implementation may encompass.

7.3.1 Automation of the choice between PPIC and Spark's local execution

The first implementation we decided to undertake, was to automatise the choice between PPIC and Spark's local execution. Using the former only for single-item patterns, and the latter for

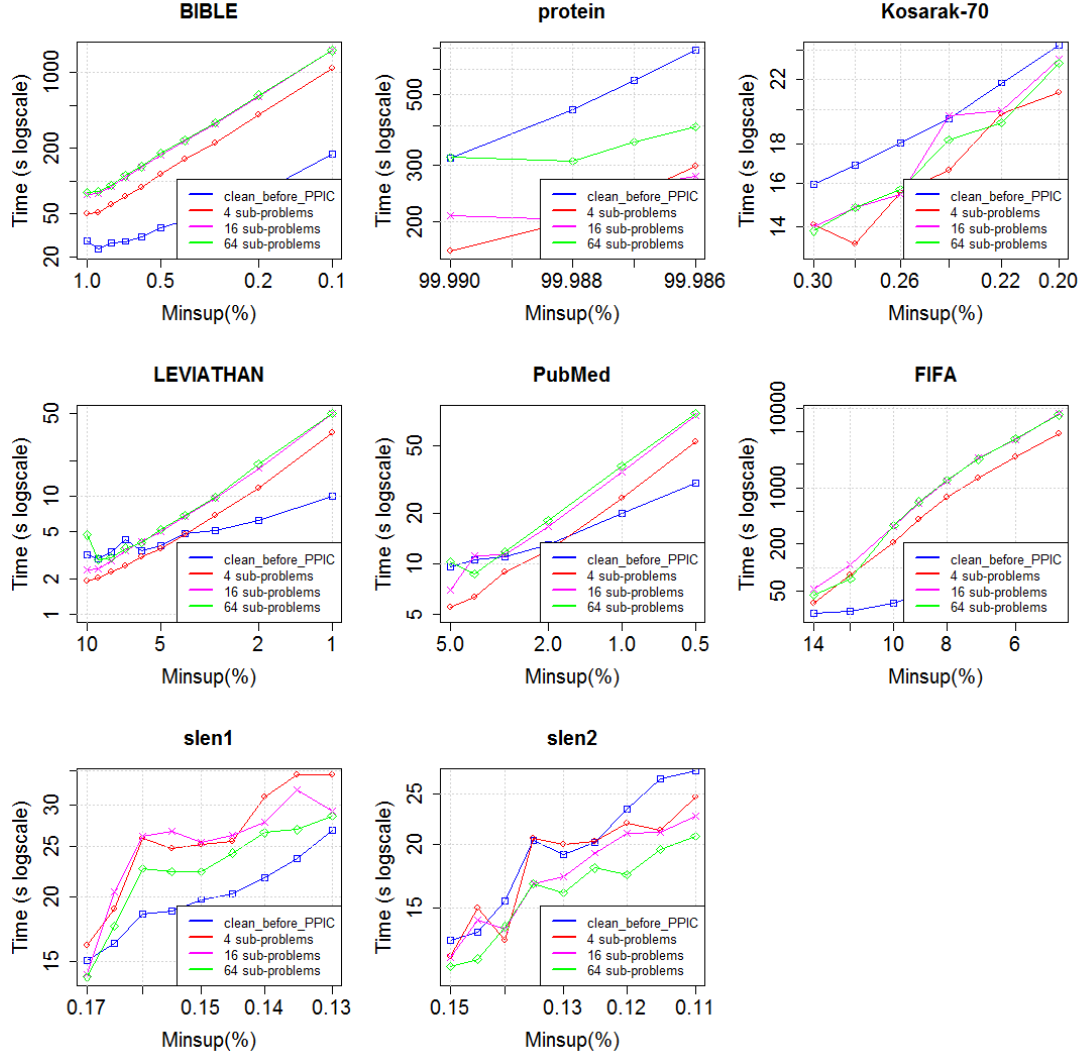


Figure 7: Performance improvement - soft limit on the number of created sub-problem.

multi-item patterns. Automatic detection being a feature we judged important to implement efficiently, since we judged it was something our user shouldn't constantly be aware of.

In that endeavour, our first go to was calculating the `maxItemPetItemSet` argument dynamically during the pre-processing stage. Since it allow us to determine whether we were dealing with single-item patterns, and thus, whether PPIC could be used.

As you can see on figure 8. While this improvement achieved its goal, slight performance degradation were also observed, due to the need for calculating this value.

However, we decided to discard this change as a more efficient way of doing that exist. As you will see in the next section, it can be implemented in a non performance damaging way.

7.3.2 Additional Pre-processing Before the Local Execution

Extending our previous idea of cleaning the sequence database before PPIC's execution, we wondered if cleaning the sequence database before Spark's original local execution could bring similar improvements.

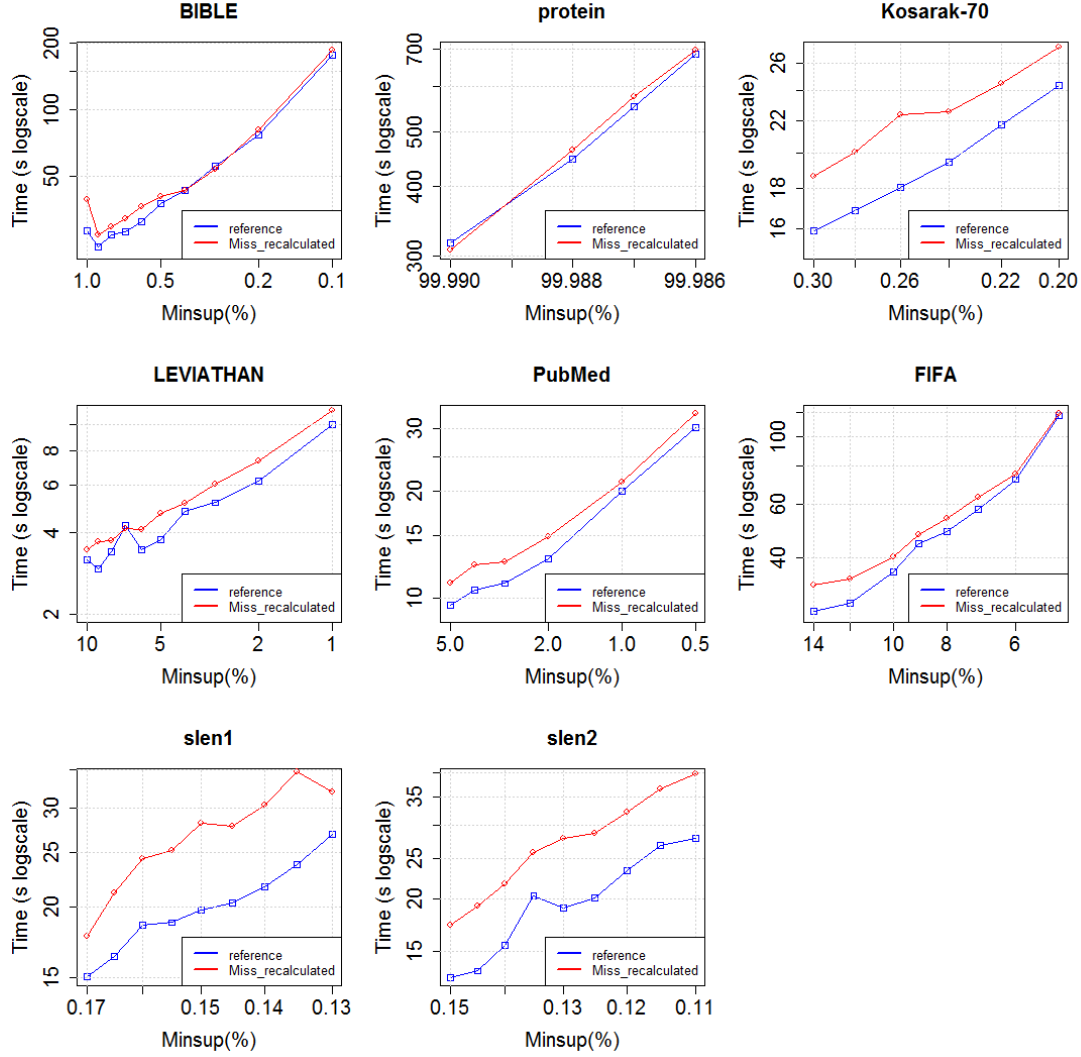


Figure 8: Performance improvement of recalculating `maxItemPerItemset` to avoid search for multi-item extension when unnecessary

We thus designed a new cleaning process suited both for Spark and PPIC’s local execution input. During this process, we also modified the algorithm so that it could automatically detect whether the current sub-problem concerned single-item or multi-item pattern mining, redirecting the execution to the appropriate local execution thereafter.

While creating the new cleaning algorithm, we also discovered a few inefficiencies on the old one. Such as the use of `ArrayBuffer` structure instead of the much more efficient `ArrayBuilder`. Our old implementation also considered the zero delimiters as items, counting them uselessly since they belonged in every sequences.

The new version of the algorithm was thus far more performant, as you can see in figure 9.

Another advantage of cleaning before the local execution, is that, should a sub-problem from a multi-item pattern dataset become a single-item pattern problem through cleaning, its local execution will be done through PPIC, which show improved performances.

For mixed datasets where some multi-item itemsets may rarely appear, this change may

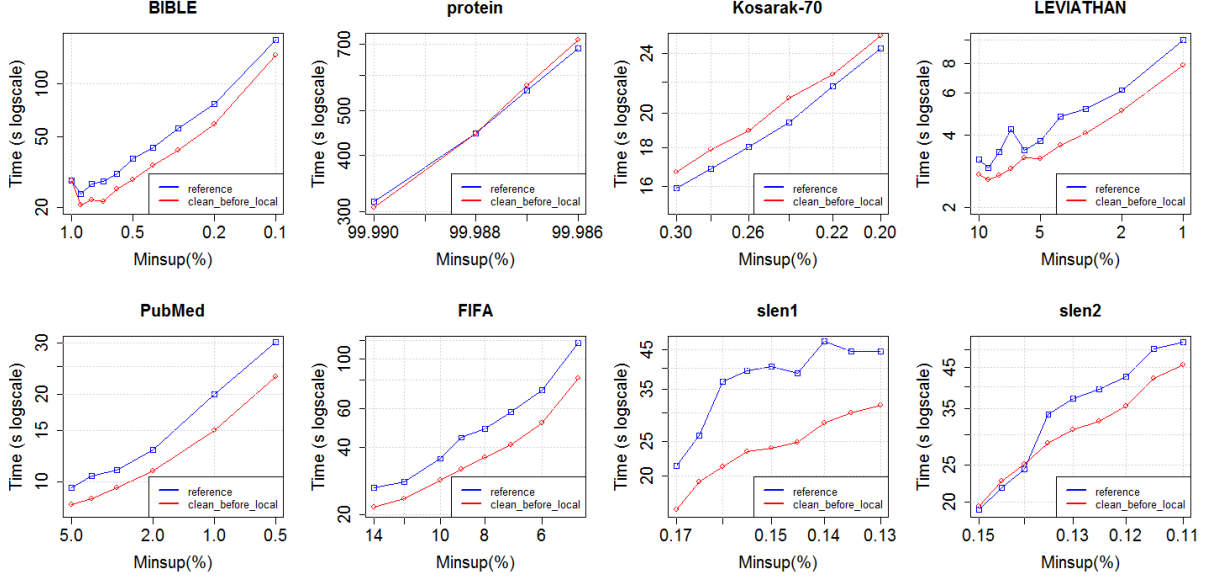


Figure 9: Efficiency gain of cleaning the sequence database before any local execution

provide great performance improvement. Sadly those kind of datasets are rarely, if at all, used in the industry. We thus couldn't get our hand on such a mixed dataset to run tests, and compare performances.

7.3.3 First/Last Position List

The second idea we had to improve Spark's performances, was to use LAPIN's position list to our advantage.

In Spark's original implementation, the scalable and local stages of the algorithm perform their duty in three phases. First they receive an solution prefix which needs to be extended and project it on the whole database, allowing them to know which sequence support that prefix. Then, in the supporting sequence only, they search for symbols which may extend the prefix. Finally, If some symbols are found and they respect the constraint applied to the solutions, Spark's will save them as solution and try to extend them further.

While this implementation is very efficient in a scalable environment, we thought it could be improved through the addition of position list. More specifically, during the prefix projection phase, we determined it would improve performance if the algorithm knew earlier when a sequence couldn't possibly hold the currently projected pattern, or if the algorithm didn't have to analyse half of the sequence before starting to project this aforementioned pattern.

We thus modified our reference algorithm into three modified version. The first using only a last position list, the second using only a first position list, and the last using both together. We then tested their performance, obtaining the results of Figure 10

As you can see, the results weren't conclusive. On every dataset but protein, which is a very sparse dataset on which this type of techniques excel, the modified algorithm delivered worse performances than our current reference algorithm.

Those worsening performance appear because the benefit of those additional computation

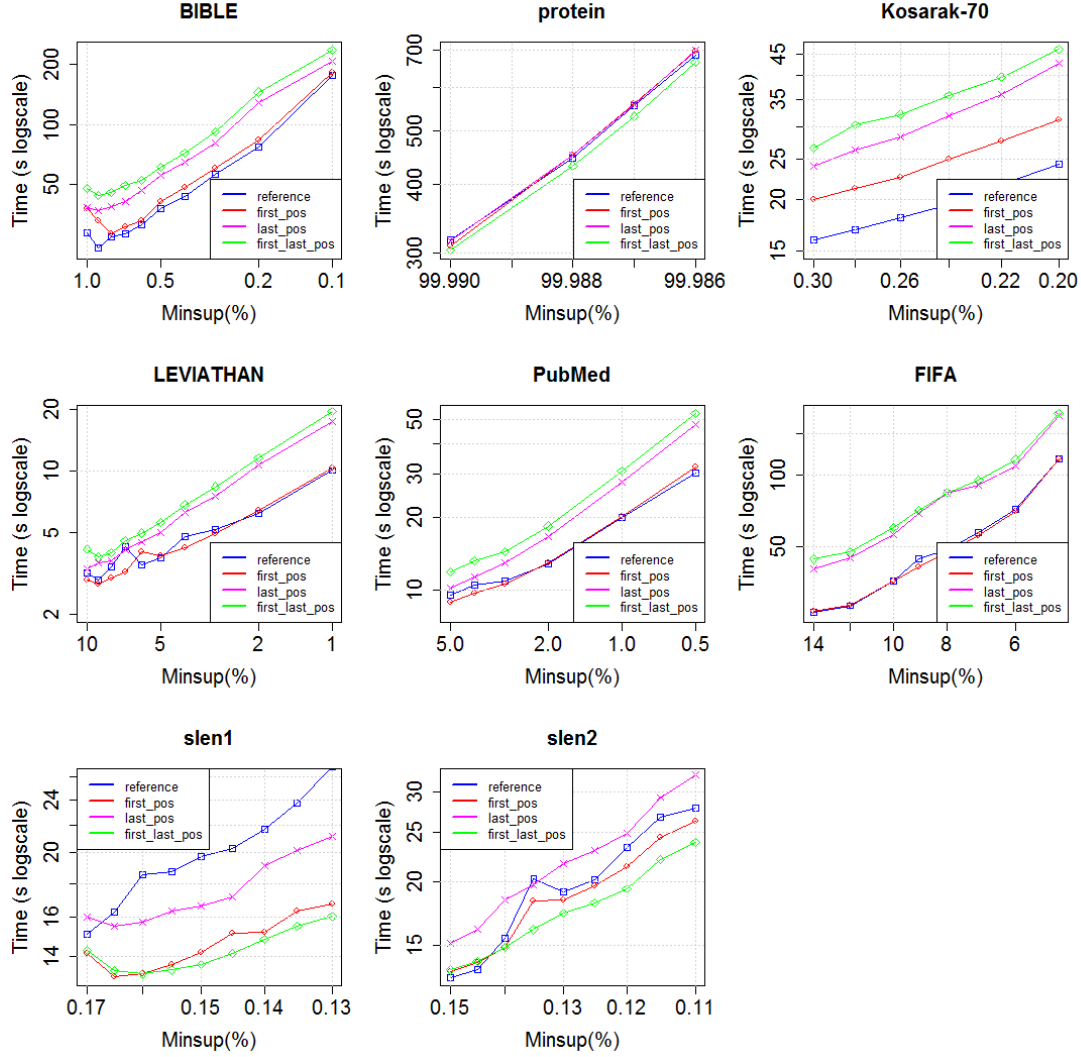


Figure 10: Performance with first/last position lists

simply do not have time to appear. After being calculated during the pre-processing, they only stay of use during the scalable execution of our algorithm. Since that, before the local execution, the sequence database will be compressed and the positions list will need to be recalculated.

Another reason for those worsening performance would be transfer time of the calculated positions list, which can be very large depending on the number of frequent items. The transfer time are thus far from negligible.

Finally, the last part explanation for those worsening performance come from our choice of multi-item pattern datasets. Since both of them are dense datasets on which the benefit of these positions list will have a hard time appearing, the calculation and transfer overhead are more strongly reflected in the performance. But those datasets were not chosen without thought either, since all multi-item datasets we managed to get our hand on displayed similar level of density. It thus seems like sparse multi-item datasets are rare oddities in the industry.

Since we saw no good reason to lower the performances of a vast majority of useful dataset, just to allow possible slight improvement on specific ones, the idea of using position list to improve performances was rejected.

7.3.4 Priority scheduling in local execution

The final idea we explored to improve Spark's performance, was to modify the order in which sub-problems are computed during the local execution stage.

In the original code, problem are decomposed until they become smaller than a size specified by the 'maxLocalProjDBSize' parameter. As mentioned before, in our reference algorithm, an extension of that idea, the subProblemLimit parameter, was also implemented to allow better control of the number of created sub-problem.

However, we have seen that a consequence of this new functionality is that sub-problems can largely vary in size, making some problem far harder to solve than others. Something would rarely appear in the original version. Coincidentally, we also realised that major drops in performance were experienced if those large problem were solved last, since some executor would be left with nothing to do while other would be stuck with atomic large workload.

The solution was clear, large problems needed to be solved first in the various executor, so that smaller workload could be shuffled around in the later stage.

7.3.4.1 Analysing sub-problem creation

The piece of code which created the various sub-problems was fundamentally a mapReduce process. The sequences from the original sequence database were one by one projected by different prefixes then mapped to some reducer, depending on the prefix's ID (see Figure 11).

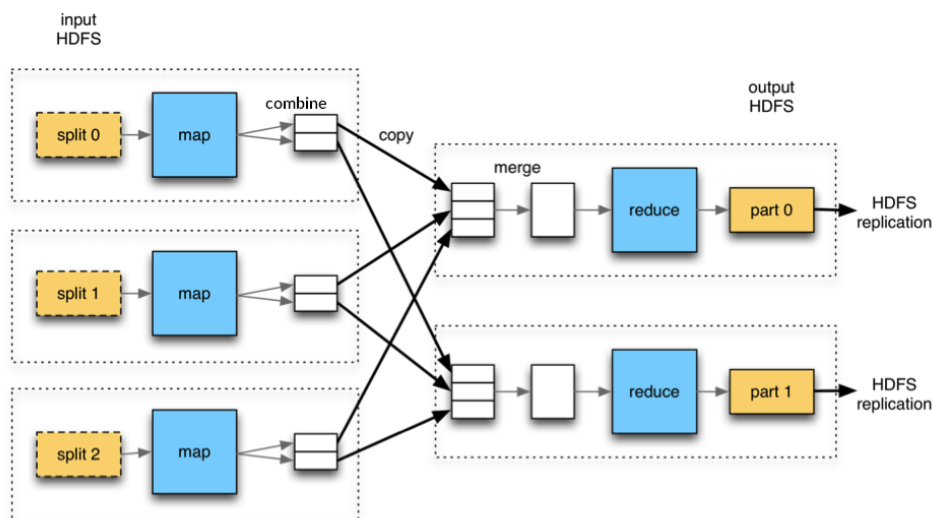


Figure 11: Spark's mapReduce

7.3.4.2 Sort sub-problem on the reducer

We thus concluded a simple solution could be implemented. According to the specification of Spark's sortBy function, the sorting stage will be executed locally on each reducer. We thus implemented this quick change by adding a simple sort between our map and reduce stage, obtaining the results of Figure 12.

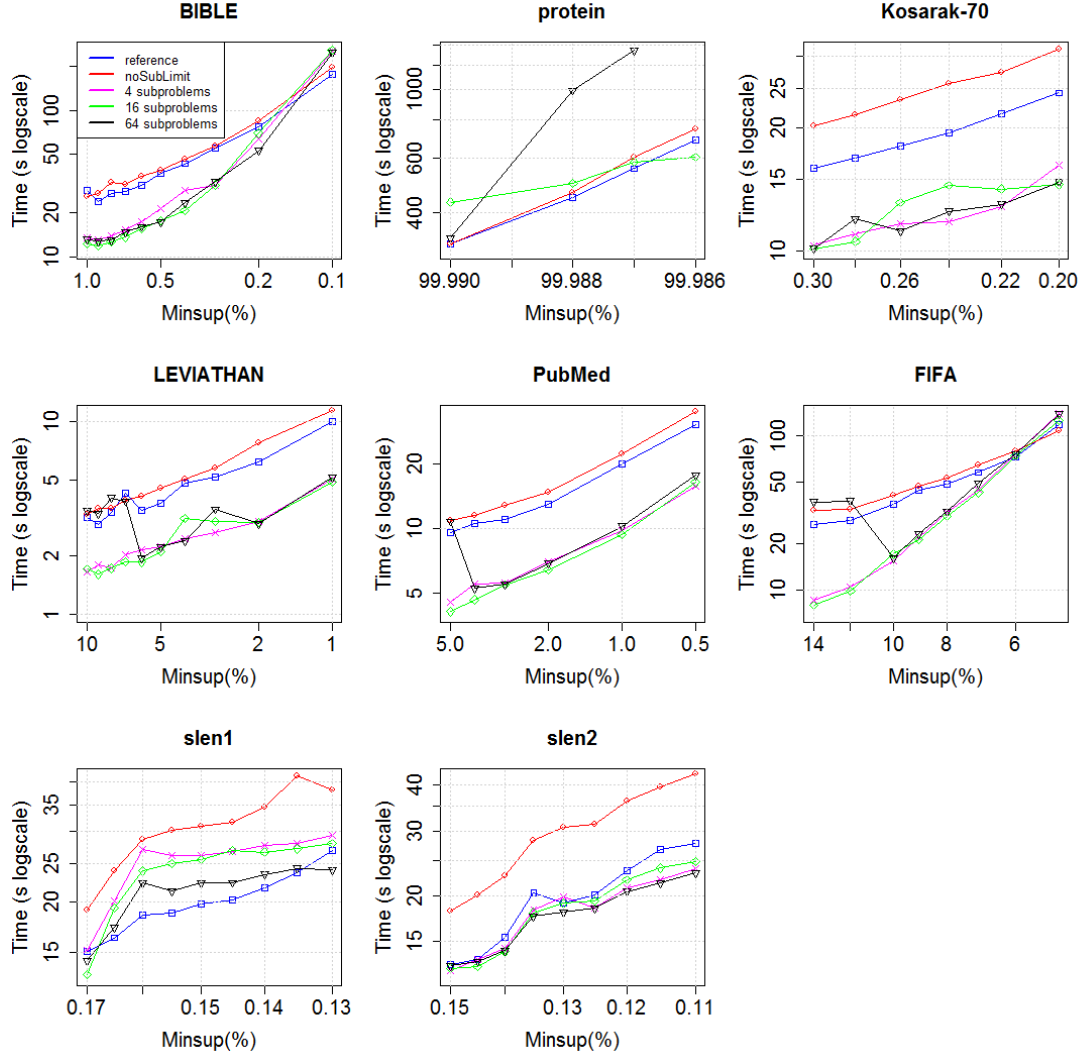


Figure 12: Performance of a 'biggest problems first' scheduling executed through the default sortBy() function of Spark.

As you can see, net performance improvements can be observed, especially for Kosarak, protein and slen2, and while BIBLE and FIFA suffer from slight performance loss in their later stages, the damage is extremely limited. Applying the sort function to the algorithm also isn't too damaging even when sub-problems aren't limited, as we can see with the red line.

At first glance, this thus looks like a nice little improvement. However, a huge problem remain, as you can see the protein measurement lack its fourth point. The reason being a constant crash of the algorithm due to a lack of memory !

In fact, to sort the various sub-problems depending on their size, Spark's need to calculate all sub-problems assigned to an executor and sort them while keeping those sub-problems either on disk or memory.

Since we forbid our algorithm from accessing the disk during our tests, asserting the 10G of memory given to each executor would be enough, we easily concluded our new implementation consumed more than 10G memory on at least one executor. Something which shouldn't happen in a scalable algorithm which should work on even the smallest configuration !

We thus had to come with another solution to sort our sub-problems. One that did not involve computing multiple sub-problem and keeping them in memory..

7.3.4.3 Sort sub-problem on the mapper

The answer was obvious, if something could be done, it would be during the mapping phase of the mapReduce. But, despite the multiple attempt at modifying Spark's map function to sort its results, we failed to accomplish our goal without using huge amount of memory. No matter what, the sorting function needed all sub-problems computing before it could work.

That is, until we realised the painfully obvious. The mapping function of Spark created sub-problem one by one, following the mapping code and sent them directly to the reducer through the groupBy function, which delivered them in the very same order they were sent.

The solution was thus to change the order in which we created our sub-problem, those would then be mapped to the reducer and be computed in the same order that they were mapped.

Meaning that, should we map the harder problem first, they would also be executed first.

We thus modified our code to sort our prefixes in descending order through their projected database size, a size which was already computed during scalable stage of the algorithm and could simply be stored in the prefix for this sort.

As it turns out, sorting a few prefixes uses nearly insignificant amount of memory in comparison to sorting huge RDD containing complete projected database, while producing equivalent performances. This implementation performances were thus a huge success, as you can see on Figure 13.

Not only to we obtain our previous boost in performance, we also cling very closely to the performance of the reference algorithm when we use the default parameters in this new extension.

This implementation advantages should thus be kept for the final version of our algorithm, since allows far better performances when using the correct set of parameters and nearly identical performance for default parameters.

7.3.5 Specialize Spark's scalable stage for single-item problems

TODO - Didn't think about it before, but it would also reduce representation size, since there is no need for delimiters.

TODO in priority, and add to final version if good perf detected.

7.4 Improving PPIC

Our final implementation attempts, were to create a CP-based implementation to solve multi-item pattern mining problems. Of course, to replace the original local execution, this implementation would need to be more performant than it's predecessor.

7.4.1 Pushing PPIC's Ideas farther

Our first attempt at creating such an implementation, was to try pushing PPIC's ideas further.

We decided to forgo all three pre-computed matrices, and to change the structure of our sequence database to fill those matrices purposes more efficiently. Instead of an array, we changed

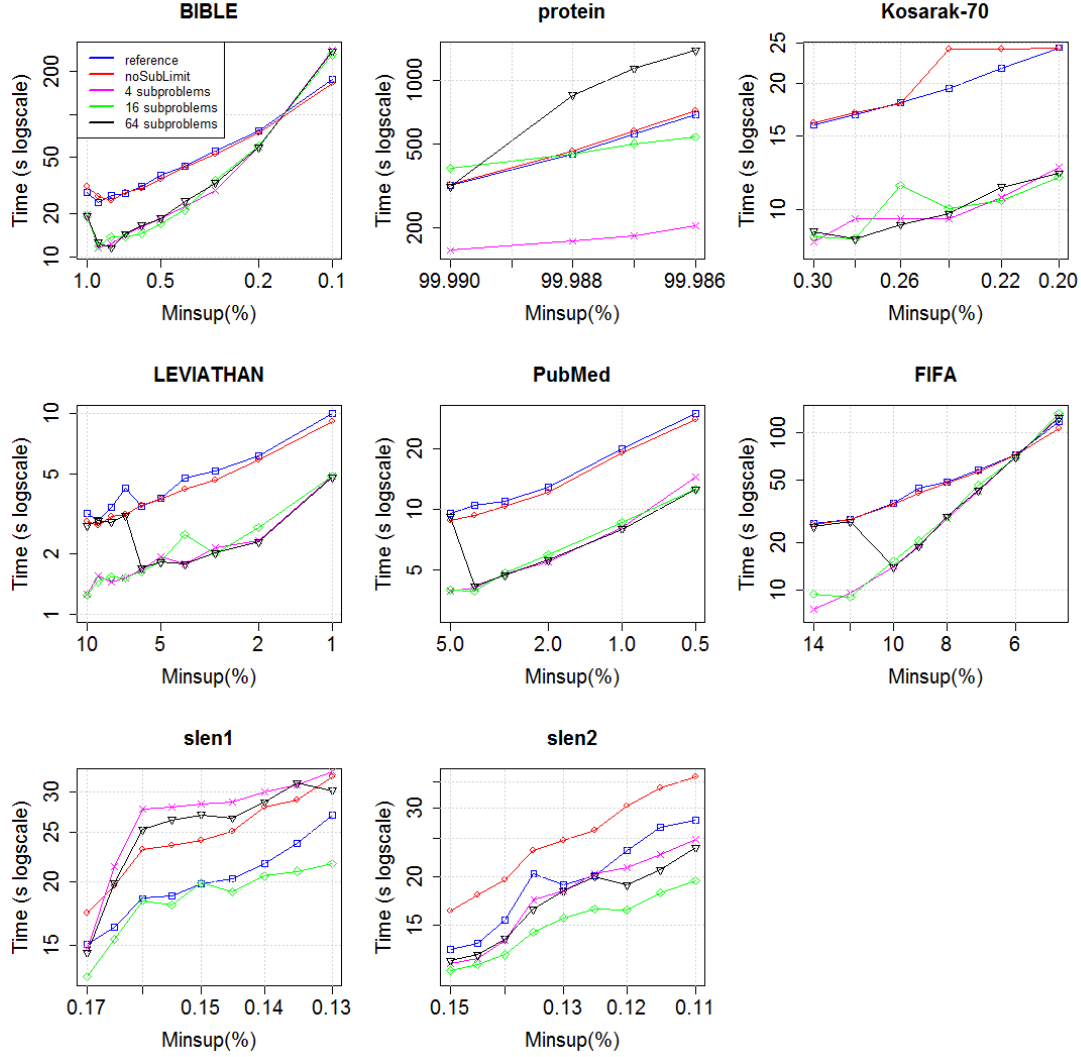


Figure 13: Performance improvement of sorting sub-problems during map stage.

each sequence into a map containing its symbols as key and the symbols positions as value. Additionally, we represented the sequence's symbols positions (including zero delimiter position) by a ReversibleArrayStack structure, thus allowing efficient backtracking through trailing.

This new structure's purpose was to allow us to make distant jumps and checks more efficiently, at the cost of an higher memory consumption. Theoretically, finding the next position of a given symbol would be $O(1)$, be this the last or first position of the symbol. Should the symbol's position list be empty, or should the last recorder be smaller than our current position in the sequence, we would also know that no more such symbol were contained in the remains of the sequence.

We thus had good faith our performance could be improved, until we saw the results of our performance tests. As you can see in Figure ??, the performance were terrible. So terrible that some measurement even surpassed our huge performance test time-out of 15000 seconds.

While we couldn't see exactly why at first, profiling our algorithm revealed it came from two majors factor. First, to attain the same purpose as the three matrices with our map structure, we needed many more trailing points which had to be backtracked at each step of the search.

Regularly generating significant period of time were our algorithm did nothing but backtrack the various sequence of its database.

The second factor was garbage collection. Since this map structure used significantly more memory than PPIC's implementation did, they were far more frequent. Damaging our performances further.

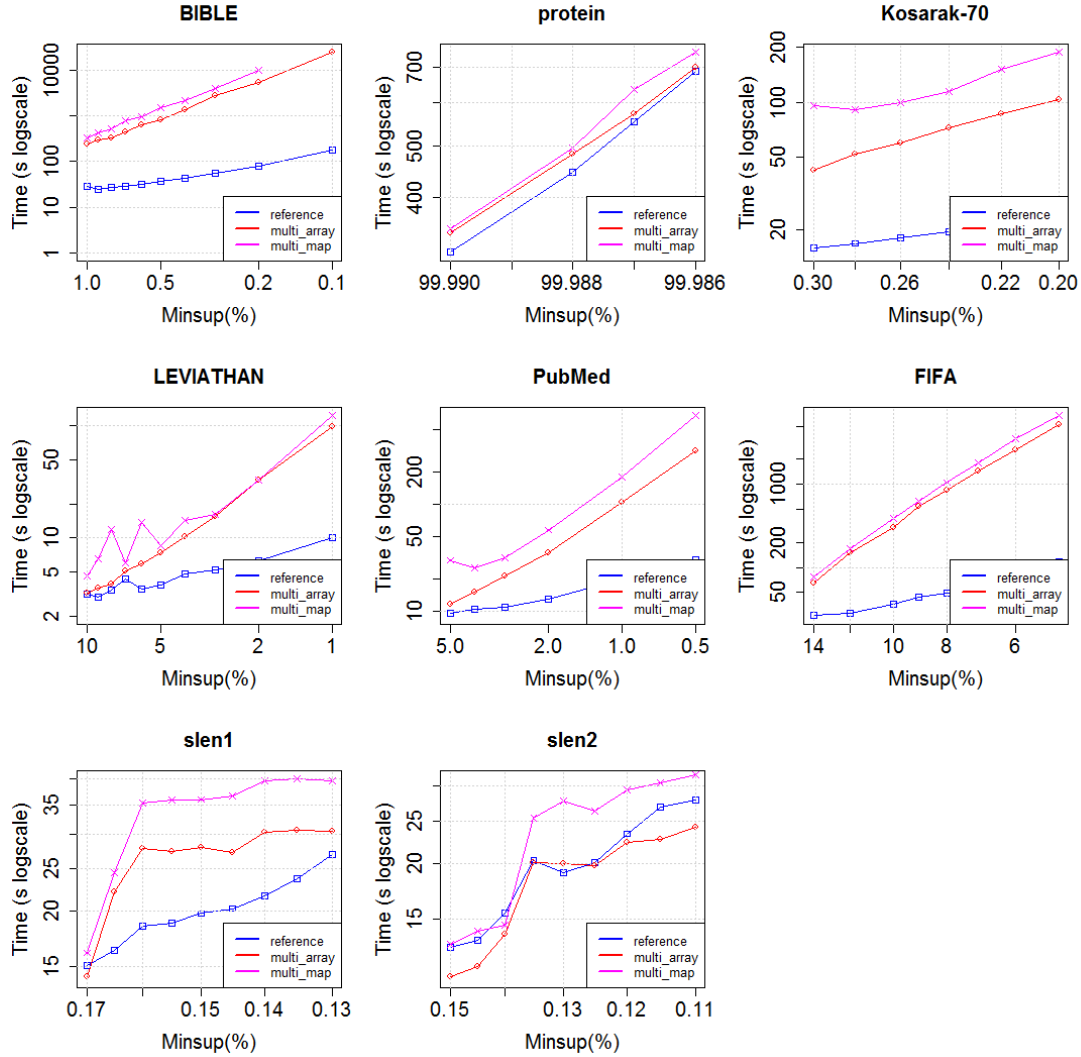


Figure 14: Performance of a CP based algorithm replacing all three pre-processed matrices from PPIC by a Map structure

7.4.2 An Implementation Closer to Spark

TODO - Solve multi-item bug + add partial-start support

8 Final Version

TODO - Redo a final version after all other is finished

9 Scalability Tests

TODO - On original simple + final implem

10 Performances Under Optimal Parameters

TODO - Only on final version and if result relevant.

11 Conclusion

TODO - ~ A page of content

12 Annexes

TODO : Add code of all implementation ? Only main ones ?

TODO : Table des abréviations + symbols !

TODO : Remerciement/Acknowledgement + pas oublier de remercier les serveurs CECI

12.1 Additional Performance Comparisons :

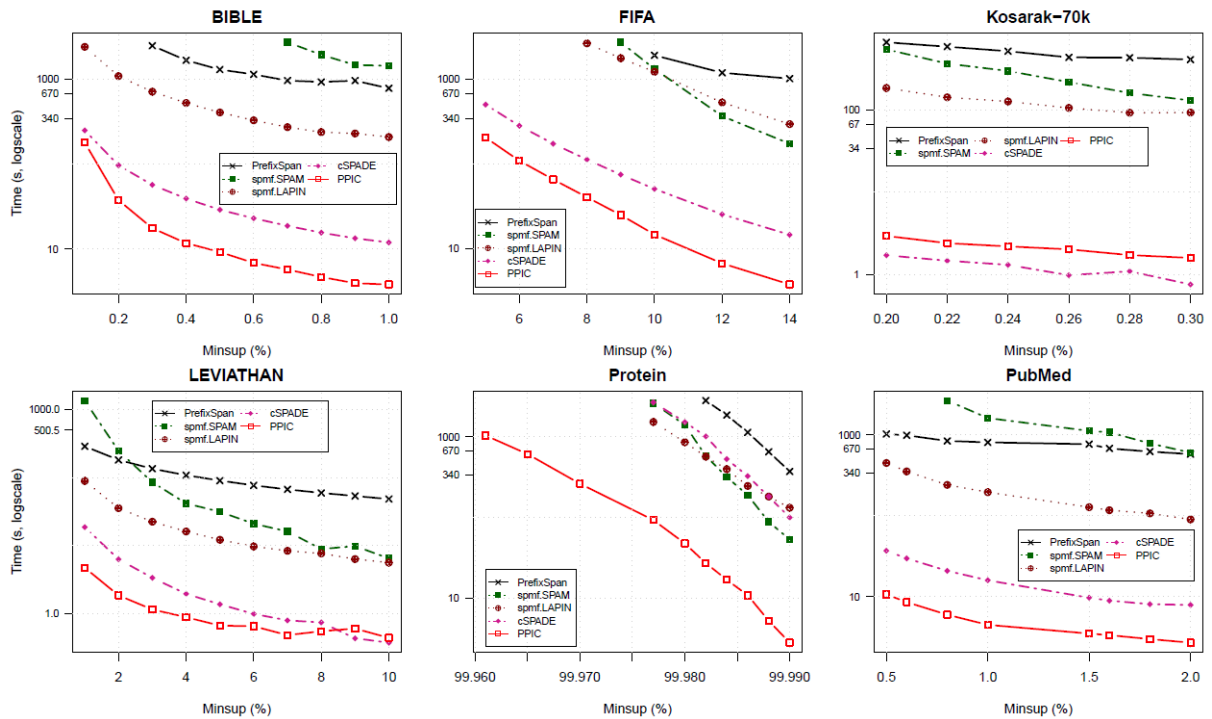


Figure 15: PPIC's performances VS other specialized algorithm

References

- [1] Matei Zaharia , Mosharaf Chowdhury , Tathagata Das , Ankur Dave , Justin Ma , Murphy McCauley , Michael J. Franklin , Scott Shenker , Ion Stoica. *Resilient distributed datasets: a fault-tolerant abstraction for in-memory cluster computing*. Proceedings of the 9th USENIX

conference on Networked Systems Design and Implementation, April 25-27, 2012, San Jose, CA

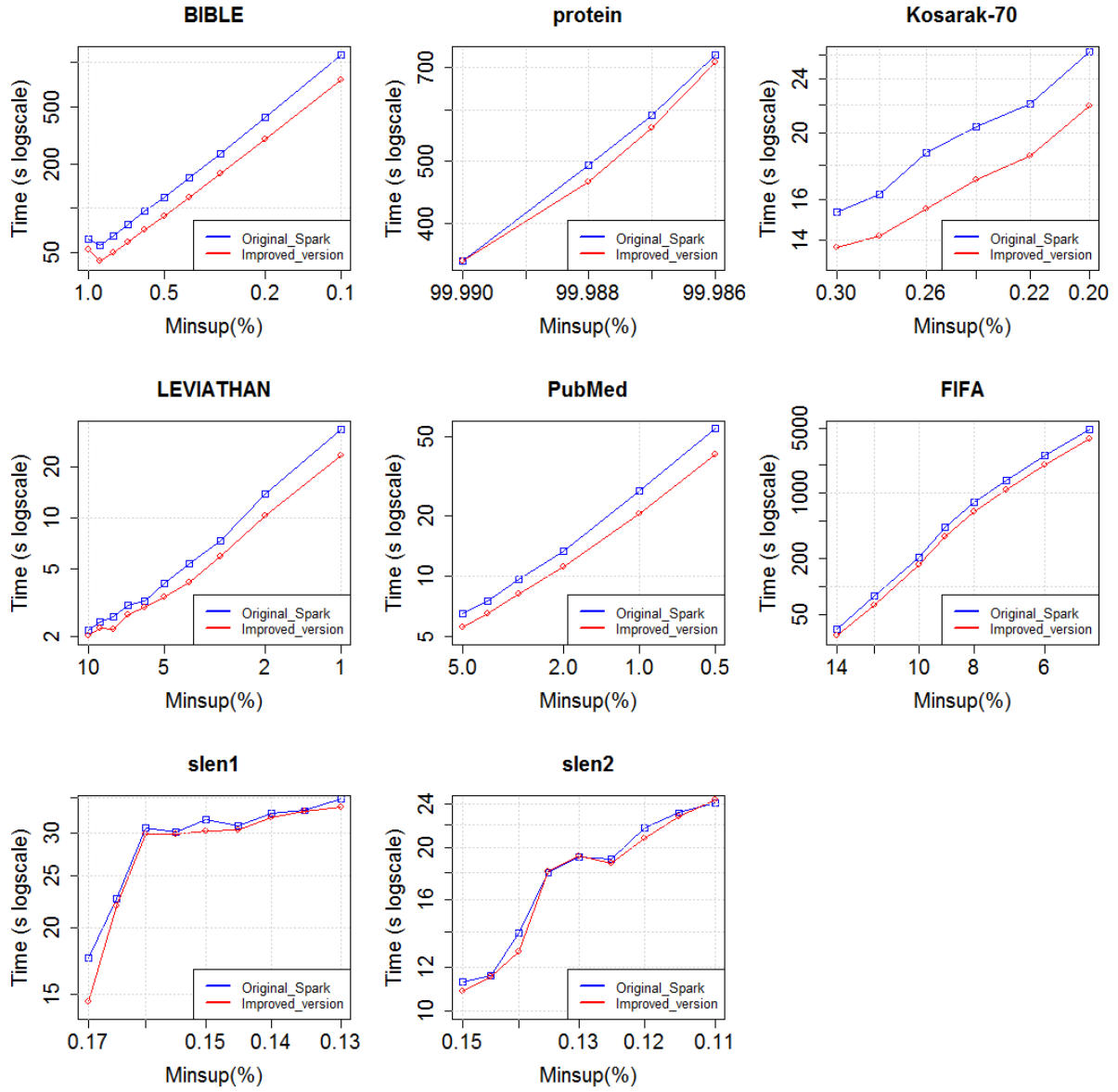


Figure 16: Performance improvement of fixing Spark's pre-processing

