

Concurrency and Parallelism

Cilk⁺ Parallel Patterns Implementation

André Rosa
48043
af.rosa@campus.fct.unl.pt

João Geraldo
49543
j.geraldo@campus.fct.unl.pt

Rúben Silva
47134
rfc.silva@campus.fct.unl.pt

Abstract—In recent years, the parallel computation paradigm has been emerging, as a consequence to: the huge growth on the amount of data that needs to be processed and analyzed, and to the switching of processors' architecture evolution process from becoming smaller and having faster clock speeds, to integrate parallel functionalities, such as multiple cores, hardware threads and vector operations. Therefore, it is imperative to build algorithms that explore these functionalities to increase the efficiency on processing such huge amounts of data. However, building algorithms on top of these functionalities can sometimes be a difficult and complex task, due to them being low level primitives and sometimes platform dependent. Thus, parallel functions and libraries can provide an easy and generic way to better utilize these resources, when available. Therefore, we developed a library, that implements some of the most well known parallel patterns, and that can be easily integrated in any already existing sequential program, to improve its efficiency. We developed each algorithm to try to achieve the maximum possible parallel slack and ...

In order to be able to use it in a vast amount of use cases, these implementations are independent both from the data types they are manipulating and the parallel functionalities provided by each specific hardware platform.

What was our approach? What were the results? What did you learn?

Finally, we conducted a preliminary experimental evaluation on the performance of the different implemented alternatives, comparing them with their corresponding sequential version, that showed

Index Terms—Parallel Algorithms, Cilk⁺

1. Introduction

Nowadays, the parallel paradigm is a matter of great importance, and the demand for more scalable and efficient data processing algorithms has increased.

In the one hand, the amount of data produced every day grows exponentially [meter uma citação], requiring more computational power to process it at the same speeds.

On the other hand, due to single-core processors are reaching the physical limits of improvement through shrinking its components and increasing its clock speed, a revo-

lution happened in the processor's architecture design. Processors evolution switched to have parallel functionalities, such as multiple cores and hardware threads as well as vector operations, that allow the simultaneously execution of processes and parallel manipulation of multiple data [1].

For these reasons, it's crucial to take advantage of these features to build algorithms that allow the processing of such huge amounts of data, having good levels of scalability and latency. However, build such algorithms can be a very difficult and complex task, not only because they have to make use of the previous low level parallel primitives, but also because those primitives can be platform dependent and are highly heterogeneous, e.g., not all the processors provide vector operations and hardware threads and the amount of cores may vary greatly, even within the same manufacture, take Intel as an example - Intel's i7 lineup has chips with only 2 cores up to chips with a whopping 16 cores.

[Falar de que existem padrões paralelos, como tá no livro, para dizer que podem servir como building blocs de algoritmos e assim] [2]

On that account, we designed and implemented a library that contains some of the most well known parallel patterns.

This library provides an easy way to build parallel programs as well as to update already-built sequential programs to make use of the hardware parallel functionalities, allowing the application-level developer to integrate such patterns and give its application a boost of performance with little effort. We resorted to Intel's Cilk⁺ [3] to implement those patterns. In order to be able to use it in a vast amount of use cases, these implementations are independent both from the data types they are manipulating and the parallel functionalities provided by each specific hardware platform. We developed each algorithm to try to achieve the maximum possible parallel slack and ... [mudar um bocado isto]

The remaining of this report is structured as follows: at Section 2 we present the architecture of each pattern, at Section 3 we discuss the implementation details of our algorithms. The settings of the experiments along with their results are presented in Section 4. Lastly, we conclude this paper with the conclusions, at Section 5.

2. Architecture

In this Section, we describe the architecture of each of the implement patterns: Map, Reduce, Scan, Pack, Split, Gather, Scatter, Pipeline and Farm. [citações para os patterns]

2.1. Map

The map pattern represents the independent application of a function to every element of a collection, and thus each operation can be executed in parallel. This function needs to be pure, e.g., not have side effects in order to be parallelizable. However, the computational weight of the application of the function might be too small in comparison to the overhead of parallelizing the tasks, leading to using the sequential version being a better option. Thus, to overcome this problem, the elements can be grouped in batches, that are process in parallel, and, within each batch, the elements are processed sequentially. The size of each batch is called *Grain Size*. [Dizer mais o que?]

2.2. Reduce

The reduce pattern represents the application of a pairwise associative operation to all the elements of a collection, producing a single element. Since the operation is associate, multiple applications (to different elements) can be done in parallel. The execution of the parallel version of these pattern produces a binary tree, where the child nodes are the operands and the parent is the result of the operation.

In our library, we implemented two versions of reduce: (regular) reduce and tiled reduce. In regular reduce, the elements of the collection are paired and the worker (operator) is applied in parallel to each pair. This process is repeated to their results until there is only one element, the reduce result. However, the weight of the application of the operators may be so small compared to the parallel overhead, that it does not pay to use the parallel version. Thus, we decided to implement an alternative version for the reduce - tiled reduce, which groups several elements of the collection in a tile, instead of only two, and executes the sequential version for each in parallel.

2.3. Scan

The scan pattern

2.4. Pack

The Pack pattern is used to eliminate wasted space in a sparse collection and to handle variable-rate output from a map. From within map, each function activation is allowed to either keep or discard its outputs. The survivors are then packed together into a single collection. We consider make our own algorithm to solve pack but because it was not as efficient as the algorithm proposed by [those dudes] so we opted to implement that algorithm.

2.5. Split

The Split¹ pattern is a variant of pack which does not discard elements, but instead packs them to the top or bottom of the output collection. In the split pattern, the order of the elements within each output segment is the same as their relative position in the input (in other words, it is a "stable" reordering), and the total output size is always exactly the same size as the input. Split can obviously be implemented by running pack twice and merging the results. But because this wouldn't be the most efficient implementation we developed a more direct algorithm to improve latency and thus more time efficient.

2.6. Gather

2.7. Scatter

2.8. Pipeline

The Pipeline pattern represents the sequential application of multiple operators (or stages) at the same element. Due to having multiple operators, they can be applied in parallel at different elements. However, having only one element at each stage can be too slow and thus, we implemented an alternative parallel version of pipeline - pipeline farm, that have multiple elements being processed at each stage in parallel and thus increase the throughput of the algorithm. This alternative version results from the combination of the Pipeline and Farm parallel patterns.

2.9. Farm

The Farm pattern is very similar to Map with the constraint that the maximum amount of parallel operations is limited to the number of farms (parameter of the pattern). Thus, each job can be distributed evenly through all the farms and then processed sequentially in the respective farm.

3. Implementation

In this Section, we ...

No map deixamos a separação em batches para o runtime do Cilk, que faz a separação automaticamente. Tentamos alterar os valores default mas obtivemos piores resultados e, visto que não podemos assumir com que tipos de dados estamos a lidar, achamos melhor deixar esta separação para o Cilk.

Na nossa implementação do reduce, de forma a evitar dataraces e a manter os resultados intermédios contíguos em memória, resolvemos utilizar dois arrays auxiliares, um para escrita e outro para leitura, que vão sendo alternados a cada iteração do reduce.

1. Optional Extra Pattern

4. Experimental Evaluation

In this Section we present an experimental evaluation of the implemented parallel algorithms, comparing their performance against their sequential version.

4.1. Experimental Setting

To evaluate the performance of our algorithms, we built a tester application that, for each algorithm, runs each of its versions (parallel, sequential and alternative, when there is one) one after the other a parameterizable amount of times, in order to compute the average of their results and thus filter the noise introduced by the For each algorithm we measured its latency for multiple amount of jobs and its latency for a fixed amount of jobs, with variable computation work. In order to count the elapsed time between the start and finish of the algorithm, we resort to a system call [citar] that returns the amount of the process CPU time, and thus not counting the time when the process is interrupted by the OS's process scheduler, giving more accurate results for the latency. The units utilized are microseconds (us). We run the experiments on node9, a computer with 16 cores and

4.2. Experimental Results

The Figure 1 has the results obtained, for each algorithm, on the experiments.

Análise comparatória da performance dizendo porque correu bem e porque correu mal para cada um. Análise geral, comparando a performance dos vários uns cons os outros. Ex o map é o que aprenseta melhores resultados comparativamente à versão sequencial, bla bla Podíamos medir a difenreça entre a paralela e a sequencial dividdo e assim tínhamos um valor que ea comparável com os outro salgoritmos

5. Conclusion

The conclusion goes here.

Acknowledgments

The authors would like to thank...

Comments

References

- [1] A. Fasiku, O. Oyinloye, S. Falaki, and O. Adewale, "Performance evaluation of multicore processors," *International Journal of Engineering and Technology*, vol. 4, no. 1, 2014.
- [2] M. D. McCool, "Structured parallel programming with deterministic patterns," in *Proceedings of the 2nd USENIX conference on Hot topics in parallelism*. USENIX Association, 2010, pp. 5–5.
- [3] A. D. Robison, "Cilk plus: Language support for thread and vector parallelism," *Talk at HP-CAST*, vol. 18, p. 25, 2012.

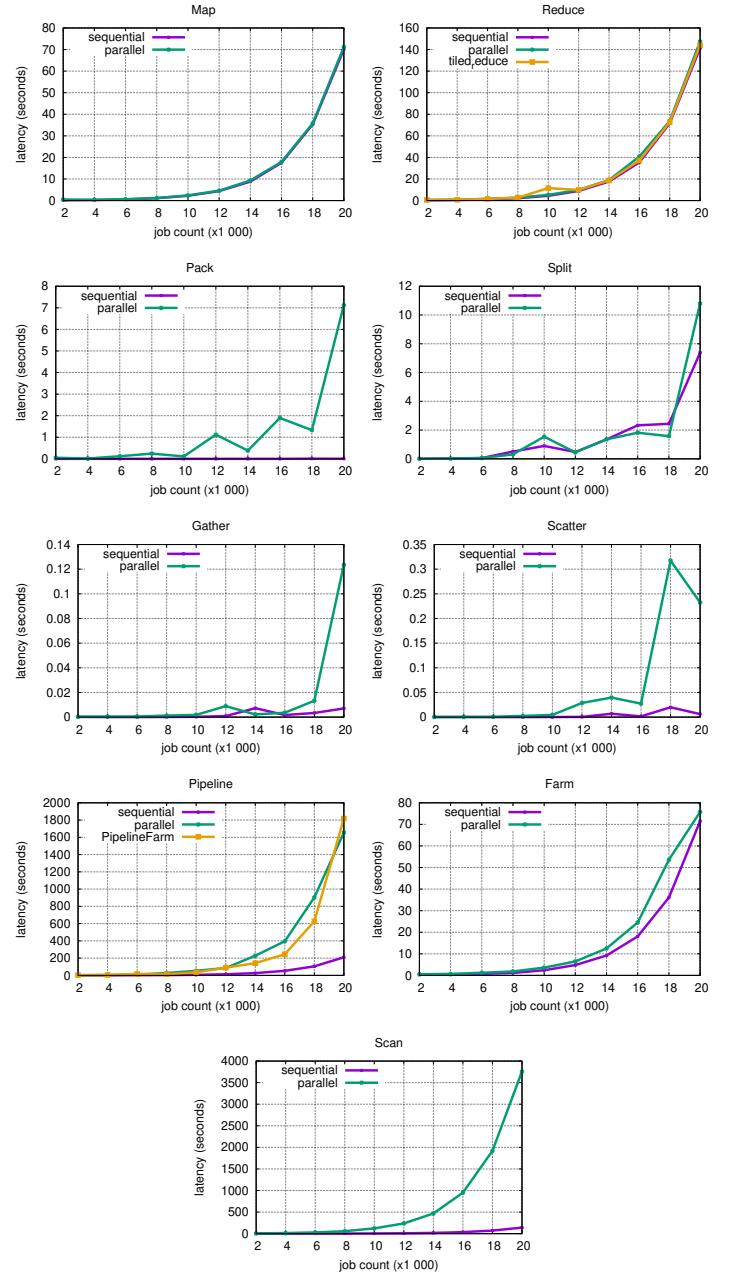


Figure 1. Latency of the algorithms.