

Parallel and Distributed Computing

First Project

Bruno Oliveira - up202208700
Mansur Mustafin - up202102355
Rodrigo Silva - up202205188



Bachelor's in Informatics and Computing Engineering
Parallel and Distributed Computing

Professor: João Resende

March 2025

Contents

1	Introduction	2
2	Algorithms Explanation	2
2.1	Naive Matrix Multiplication	2
2.2	Line Matrix Multiplication	3
2.3	Block Matrix Multiplication	3
2.4	Parallel Matrix Multiplication	4
2.4.1	First Version	4
2.4.2	Second Version	4
3	Performance Metrics	5
4	Results and Analysis	5
4.1	Performance Evaluation of a Single Core Implementation	6
4.2	Performance Evaluation of a Multi-Core Implementation	6
5	Conclusion	7
6	Appendix	8
6.1	Comparison Charts	8

1 Introduction

This report presents the algorithms, performance results and respective analysis of the first project in the Parallel and Distributed Computing course. The objective of this project was to evaluate and compare the performance of different memory access patterns and parallelization techniques.

Matrix multiplication was chosen as the central problem due to its relative simplicity and its suitability for testing various performance-related aspects. In particular, we examine how operating systems and hardware handle memory access. Specifically, most CPUs employ techniques such as caching and prefetching, which retrieve data in contiguous blocks (cache lines). Consequently, when one element is accessed, the processor often fetches adjacent elements as well. To observe this effect, we implemented three single-core algorithms. To explore how different parallelization strategies can improve performance, we also implemented two multi-core algorithms (see Algorithms Explanation section).

By developing the code for these tasks, we became more familiar with modern and high performance C++ and Lua, since we chose it as our alternative implementation language. Moreover, we also had to research about matrix multiplication algorithms and in order to implement multi-core solutions, we explored and experimented with the OpenMP standard, allowing us to learn more about CPU multithreading.

2 Algorithms Explanation

To explore how cache locality and parallel computing affect matrix multiplication, we implemented five different algorithms. Although cache-related behavior is largely independent of the programming language, we additionally provide a Lua implementation to demonstrate that performance differences arise from algorithmic choices rather than language specifics. Although these algorithms perform the same operations and produce the same result, the order of operations can significantly affect memory access patterns and thus performance.

2.1 Naive Matrix Multiplication

The *naive* matrix multiplication algorithm is closely related to its hand calculation. This method for multiplication of two matrices: $A \times B = C$, where $A(m, n)$, $B(n, p)$ and $C(m, p)$, obtains the sum of a cell (i, j) in the solution, by linear combination of row i of A and column j of B .

Implementing this algorithm is fairly straightforward, as we only need to iterate through every row and column of the solution matrix, and, for each cell, perform the linear combination of the values:

As we can expect from the 3 nested loops, the overall time complexity for the matrix multiplication algorithm is (n^3) . This time complexity is maintained among

```
for (i = 0; i < m; i++)
    for (j = 0; j < p; j++)
        for (k = 0; k < n; k++)
            C[i][j] += A[i][k] * B[k][j];
```

Figure 1: Overview of Naive Matrix Multiplication Algorithm

all the implemented algorithms, therefore performance impacts are mainly caused by cache accesses and parallelism.

2.2 Line Matrix Multiplication

The line matrix multiplication algorithm is very similar to the previous algorithm. The only difference is that the second and third loops switch places:

```
for (i = 0; i < m; i++)
    for (k = 0; k < n; k++)
        for (j = 0; j < p; j++)
            C[i][j] += A[i][k] * B[k][j];
```

Figure 2: Overview of Line Matrix Multiplication Algorithm

In practice, the difference between this and the previous approach is that instead of calculating the result matrix cell by cell, we instead focus on calculating and storing some terms of the linear combination line by line. Because matrices are stored in row-major format, accessing data consecutively along a row can improve cache locality.

Although the correctness and overall time complexity of this algorithm remains the same.

2.3 Block Matrix Multiplication

The last single-core algorithm implemented is the block matrix multiplication algorithm. This approach expands on top of the previous one and the same logic applied in calculating subresults line by line is now used to calculate the solution block by block. For this, we need 6 nested loops:

```
for (I = 0; I < m; I += block_size)
    for (K = 0; K < n; K += block_size)
        for (J = 0; J < p; J += block_size)
            for (i = I; i < min(I + block_size, m); i++)
                for (k = K; k < min(K + block_size, n); k++)
                    for (j = J; j < min(J + block_size, p); j++)
                        C[i][j] += A[i][k] * B[k][j];
```

Figure 3: Overview of Block Matrix Multiplication Algorithm

The main idea is that now, we have divided the matrices into smaller blocks and the operations are performed block by block. This approach can enhance cache locality by keeping calculations within a more localized area, potentially improving overall efficiency.

2.4 Parallel Matrix Multiplication

Besides the previous three single-core algorithms, we also implemented two variations of the line algorithm using multi-core processing. For this, we used the OpenMP API for parallel programming in C++. The next sections show an explanation and overview of the implementations.

2.4.1 First Version

The first approach aims to parallelize the outermost loop (i) of the line algorithm. There will be a division of i values between the different threads so that each one works on a separate row at the same time (they also get their own k and j variables). Due to working on separate rows, the threads do not need synchronization.

For this, we used OpenMP's `parallel for` and `private` variable pragmas:

```
#pragma omp parallel for private(i, k, j)
for (i = 0; i < m; i++)
    for (k = 0; k < n; k++)
        for (j = 0; j < p; j++)
```

Figure 4: Overview of First Parallel Matrix Multiplication Algorithm

2.4.2 Second Version

The second approach aims to parallelize the innermost loop (j) of the line algorithm. The thread pool will work on the same row at the same time but on different columns. In this case, the threads need to be synchronized after finishing the row operations which can lead to some overhead.

For this, we used OpenMP's `parallel region`, `for` and `private` variable pragmas:

```
#pragma omp parallel private(i, k)
for (i = 0; i < m; i++)
    for (k = 0; k < n; k++)
        #pragma omp for private(j)
        for (j = 0; j < p; j++)
```

Figure 5: Overview of Second Parallel Matrix Multiplication Algorithm

3 Performance Metrics

The performance of the algorithms was evaluated using the following metrics:

- **Time:** The time required to execute the algorithm. This is a crucial metric for comparing the speed and efficiency of multicore implementations.
- **Speedup:** Defined by the ratio of the sequential execution time T_{seq} to the parallel execution time T_{par} :

$$\text{Speedup} = \frac{T_{\text{seq}}}{T_{\text{par}}}.$$

This indicates how much faster the parallel algorithm runs compared to the sequential version. A higher speedup indicates more effective utilization of parallel resources.

- **Efficiency:** Defined by the ratio of the speedup to the number of threads p :

$$\text{Efficiency} = \frac{\text{Speedup}}{p}.$$

This measures how effectively the algorithm scales with an increasing number of threads. An efficiency of 1 (or 100%) indicates perfect scaling.

- **FLOPS:** The number of floating-point operations per second, given by the formula

$$\text{FLOPS} = \frac{2 \times n^3}{\text{time}},$$

where $2 \times n^3$ represents the total number of arithmetic operations (one multiplication and one addition per iteration) in the matrix multiplication, and time is the total execution time. Higher FLOPS values indicate better numerical performance.

- **Cache Misses:** To evaluate memory access patterns and their performance impact, we use number of cache miss, such as L1_DCM (Level 1 Data Cache Misses), L2_DCM (Level 2 Data Cache Misses), and L3_TCM (Level 3 Total Cache Misses). Each miss forces the processor to fetch data from a slower cache level or main memory, incurring extra clock cycles. Reducing cache misses can improve time execution by minimizing these wasted cycles.

4 Results and Analysis

To benchmark the implementations previously discussed, we executed the test set shown below 3 times, storing the results for analysis. All graphs produced are available in the Appendix section.

- Naive, Line and Parallel implementations (across both C++ and Lua) with input matrices from 600×600 to 3000×3000 with size increments of 400
- Line, Block and Parallel implementations (in C++) with input matrices from 4096×4096 to 10240×10240 with size increments of 2048, with block sizes 128, 256 and 512

4.1 Performance Evaluation of a Single Core Implementation

The first comparison in the single core implementations is between the naive and the line approach, which was compared both in C++ and Lua. The results show that, in both languages, the execution times are smaller for the latter, being even lower on the C++ Implementation when compared to Lua. Evidence for this can be seen in graphs 6 and 7.

The main reason why the line implementation is more performant than the naive approach is due to memory accesses, especially, cache misses. Due to the way memory is loaded in the line implementation, memory accesses are much more organized and sequential, leading to a higher rate of cache hits which reduces execution bottlenecks such as memory fetching.

Having established the advantages of the line implementation, we can try to further optimize the code to reduce cache misses during memory accesses. The best approach for this is to further subdivide the matrix calculation and perform it in small blocks. Theoretically, the best block size would be \sqrt{size} , however, we decided to test with different matrix and block sizes. The results of our benchmarks can be consulted in figures 9 and 10.

As we can see in graph 9, the execution times for the 4 algorithms start very similar for the smallest matrix size, but then, for the second measurement all block sizes run faster than the line algorithm.

We would expect this trend to continue on as we increase the matrix size, however, for matrix size of 8192×8192 , the block algorithms with block sizes of 256 and 512 suffer a spike in execution time (and consequent decrease in Mflops). We believe these results were caused by the number of cache misses that occur on the L3 cache. To verify this, we ran some tests while measuring the L3TCM and the results are shown in 13 (figures 11 and 12 depict the same measures for the other caches).

Finally, for the larger block sizes, the behavior aligns with the expectations and all the block algorithms are faster than the line algorithm.

4.2 Performance Evaluation of a Multi-Core Implementation

Besides the single core algorithms, we also implemented two multi-core approaches using OpenMP and based on the line algorithm. As expected, both multi-core algorithms have better performance when compared to the single core line approach.

By distributing the operation load across different processing units, these approaches are able to perform a larger number of floating point operations per second and, consequently, have smaller execution times, as we can see in graphs 14 and 15.

Nonetheless, the multi-core implementations follow different ideas, as discussed in the Algorithms section and have substantially different results. While still being faster than the sequential code, the first implementation is able to execute 3 to 4 times more Mflops than the second one.

While both take advantage of parallelism, the second approach has an associated synchronization overhead after the inner parallelized loop, which can lead to sizeable overheads, especially for larger input sizes.

As a result of this, we observe faster execution times for the first parallel approach which uses parallelism on the outermost loop, not requiring any synchronization whatsoever. The comparison between these two algorithms is evident when comparing their speedup and efficiency.

While the second approach has a speedup of slightly over $1\times$ for matrices over around 1000×1000 and efficiency below 0.2, the first implementation is at least $4\times$ faster (usually even around $6\times$) and has an efficiency over 0.5 when compared to the sequential version. These conclusions can be seen in the figures 16 and 17.

5 Conclusion

In conclusion, this project allowed us to deepen our understanding of fundamental programming concepts, such as data accesses, with emphasis on the impacts of cache misses, and parallel computing. By implementing multiple valid approaches to the same problem and comparing their performance metrics, we were able to verify the concepts taught in the theoretical classes.

Developing the code across two languages, in our case C++ and Lua, also provided us with better insights in those languages, particularly C++, where we explored the Performance API (PAPI) for benchmarking and OpenMP for the multi-core programming.

Overall, the project successfully met its objectives, giving us hands-on experience with the intricacies of parallel computing and demonstrating the performance capabilities of programs developed with memory layout and multi-core mechanisms in mind.

6 Appendix

6.1 Comparison Charts

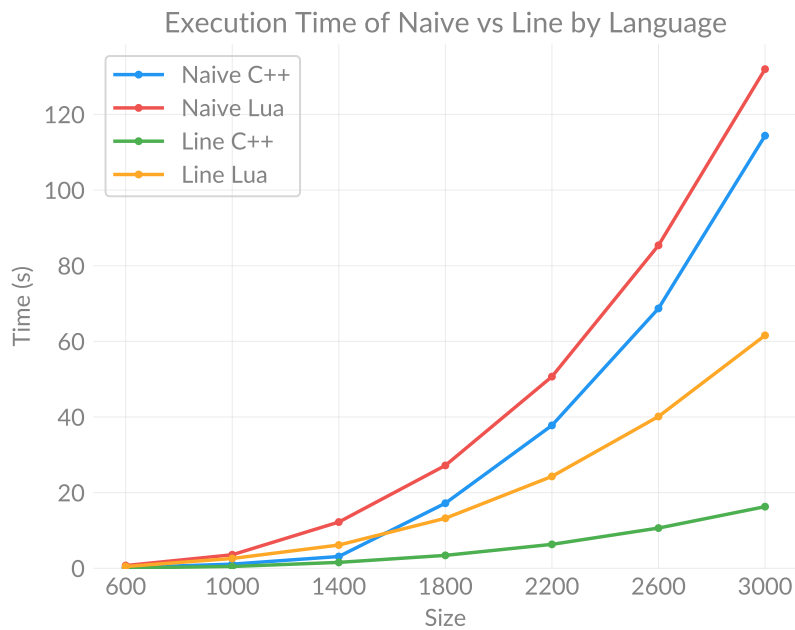


Figure 6: Time comparison between naive and line multiplication, in both C++ and Lua

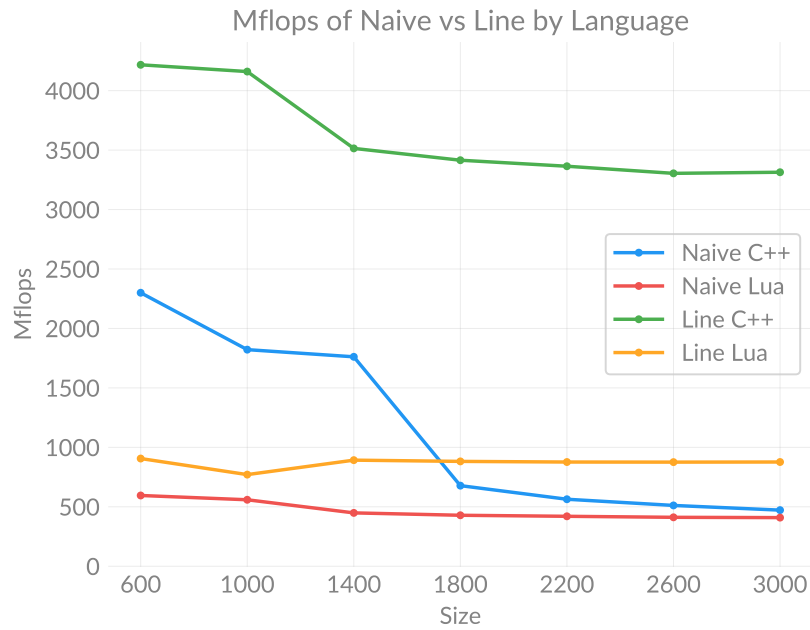


Figure 7: Mflops comparison between naive and line multiplication, in both C++ and Lua

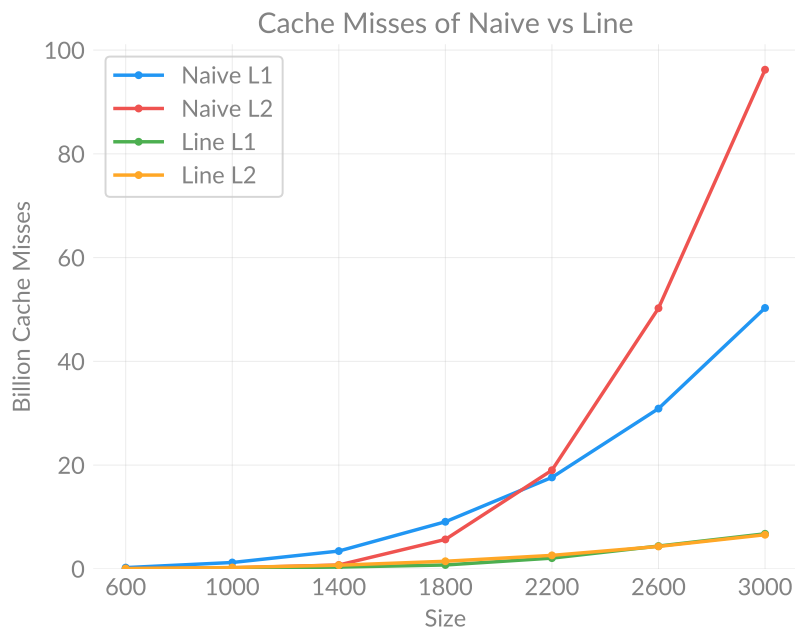


Figure 8: Cache misses comparison between naive and line multiplication, in C++

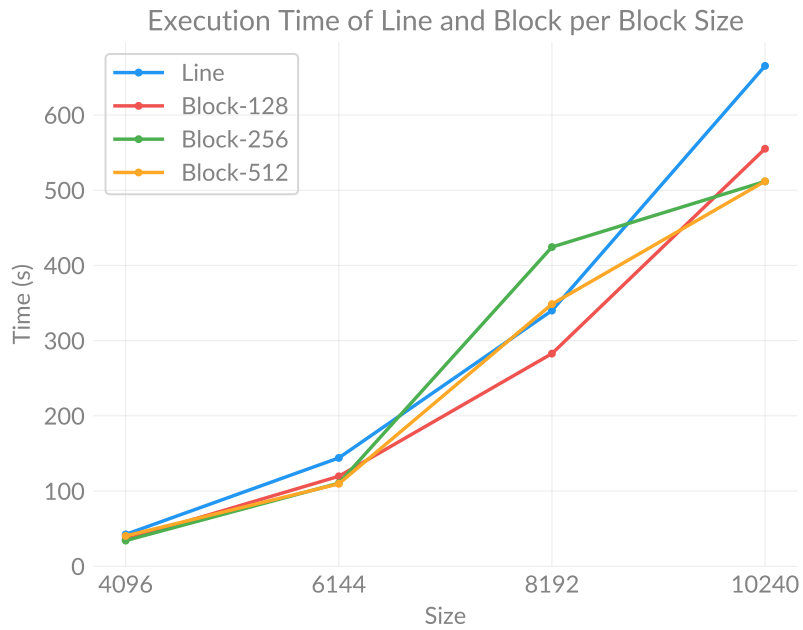


Figure 9: Time comparison between line and block multiplication, with different block sizes

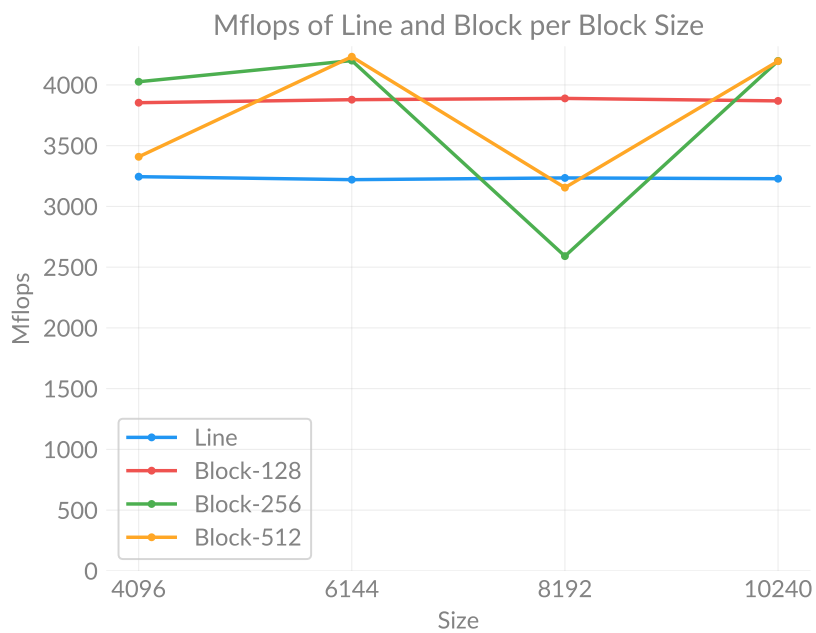


Figure 10: Mflops comparison between line and block multiplication, with different block sizes

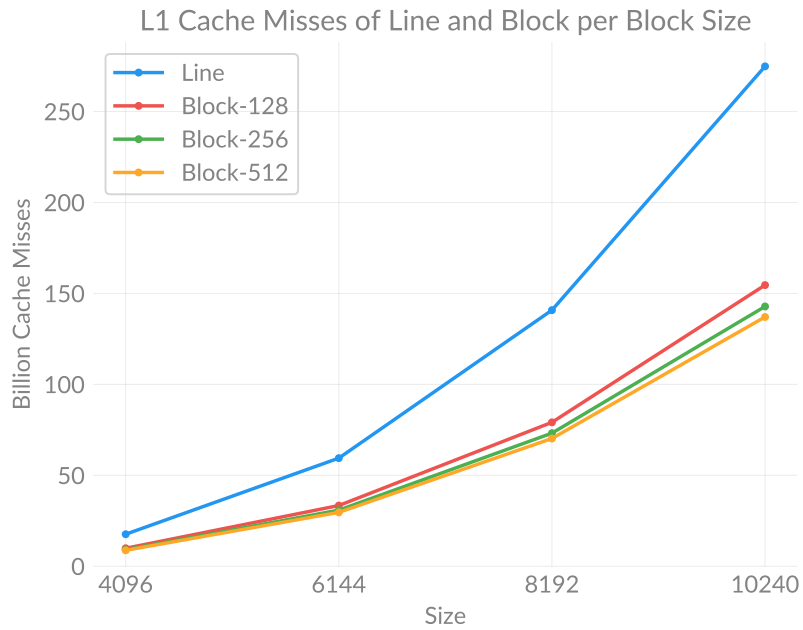


Figure 11: L1 cache miss comparison between line and block multiplication, with different block sizes

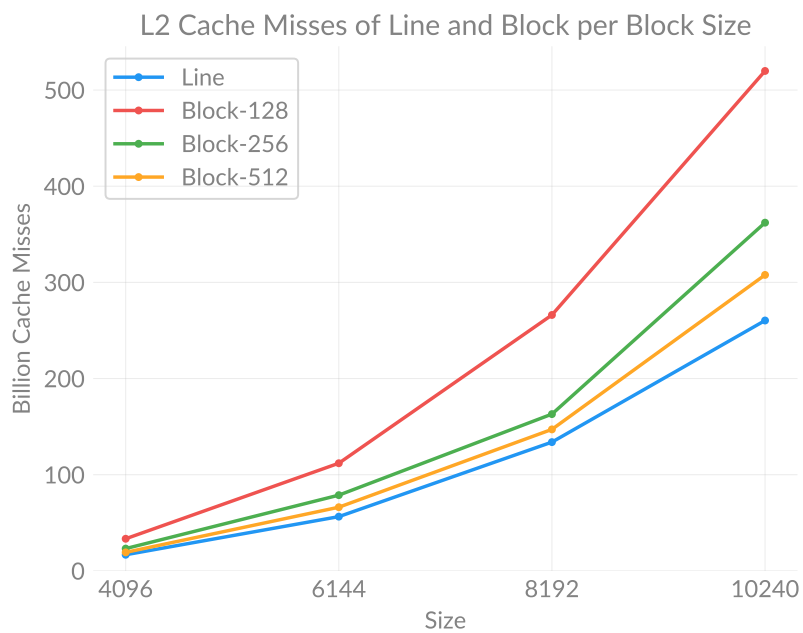


Figure 12: L2 cache miss comparison between line and block multiplication, with different block sizes

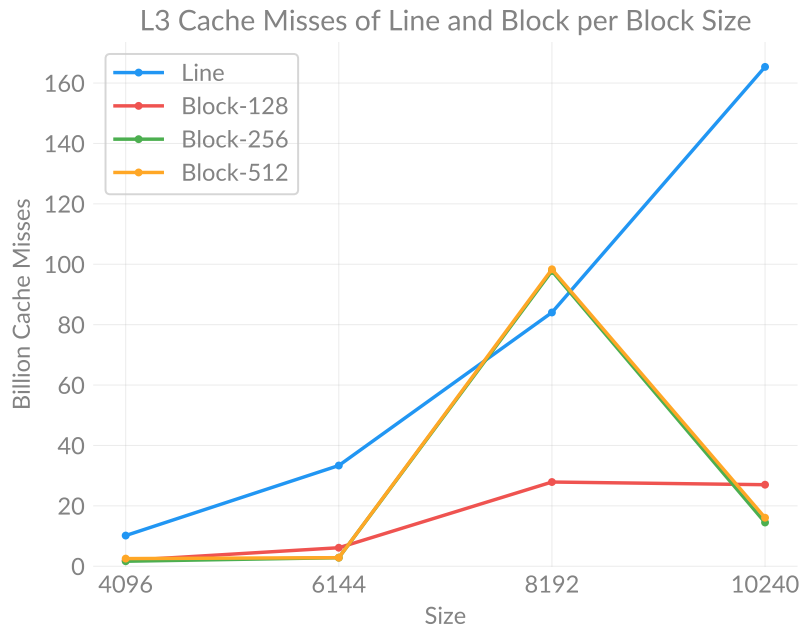


Figure 13: L3 cache miss comparison between line and block multiplication, with different block sizes

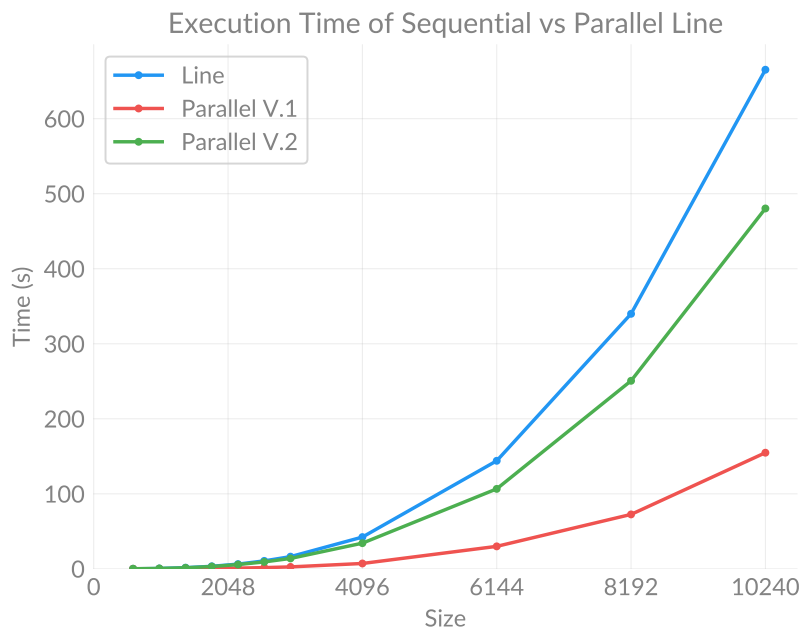


Figure 14: Time comparison between sequential and parallel line multiplication

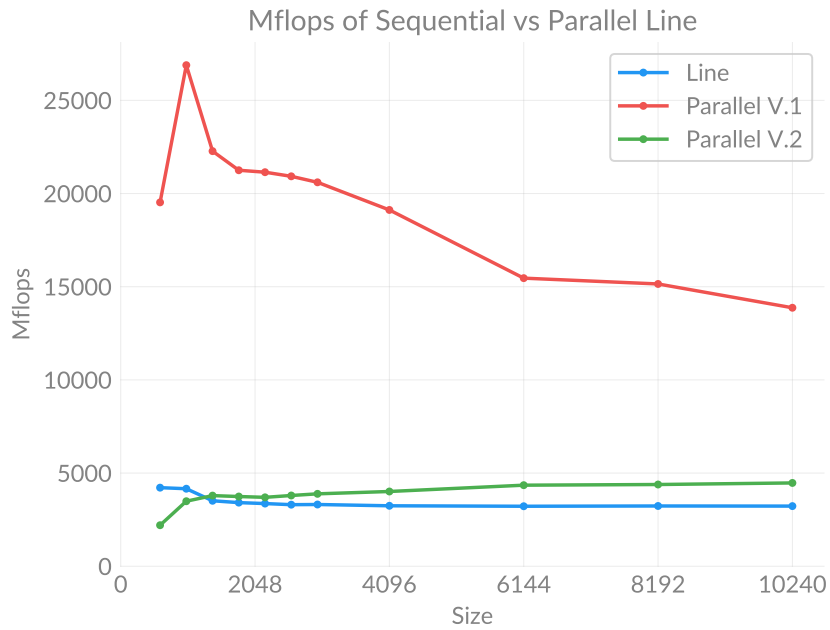


Figure 15: Mflops comparison between sequential and parallel line multiplication

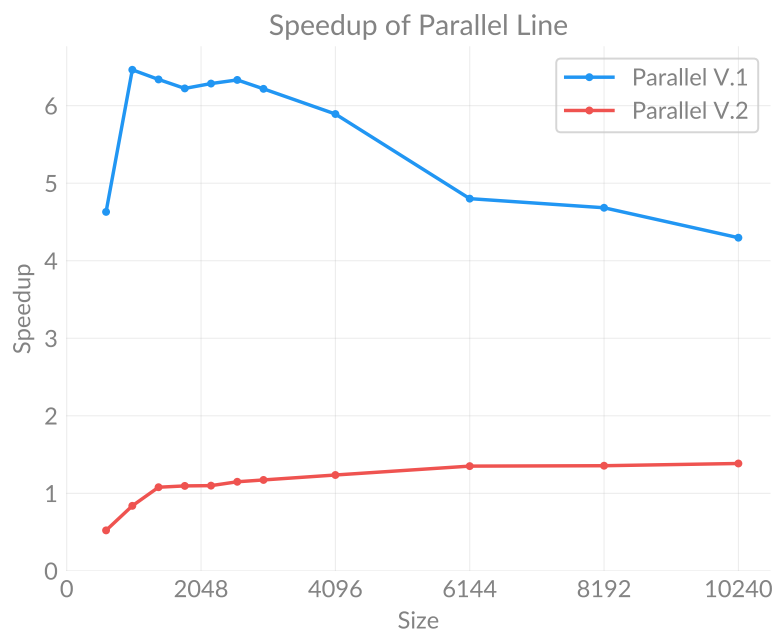


Figure 16: Speedup comparison between versions of parallel line multiplication, relative to the sequential version

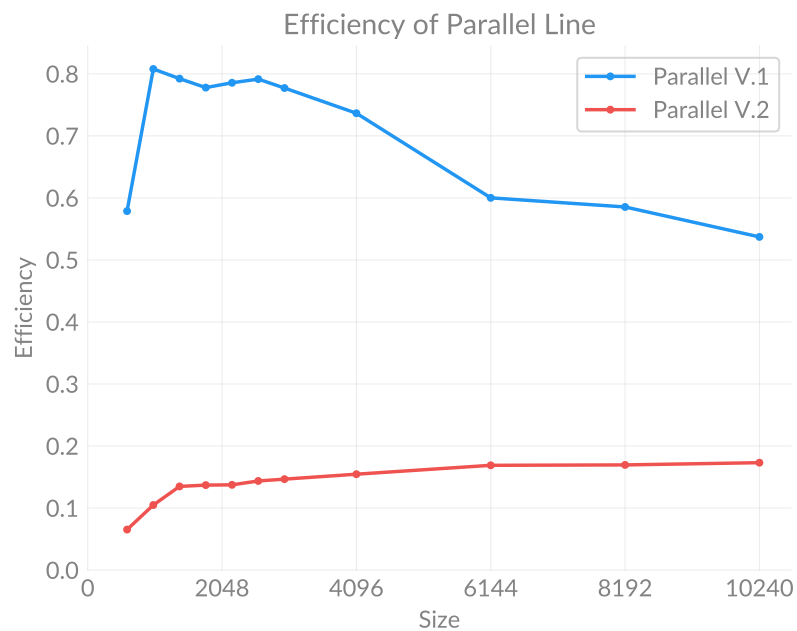


Figure 17: Efficiency comparison between versions of parallel line multiplication, relative to the sequential version