# Parallel and Vectorized Matrix Multiplication

Adonai Ojeda Martin

November 2024

**Abstract**

Matrix multiplication is a fundamental operation in various scientific and engineering computations, making it essential to optimize its performance. This paper presents a detailed benchmarking analysis of three matrix multiplication algorithms: basic, parallel, and vectorized implementations. We evaluate the performance of these algorithms in terms of execution time, CPU usage, and memory consumption across multiple matrix sizes. The results show that the parallel implementation provides significant speedup over the basic implementation. The vectorized implementation, leveraging modern hardware optimizations, achieves the highest speedup, particularly for large matrices. These findings underscore the importance of algorithmic optimization and hardware utilization in computationally intensive tasks, offering insights into the efficiency of parallel and vectorized techniques for large-scale matrix operations.

## 1 Introduction

Matrix multiplication is one of the most fundamental operations in linear algebra and plays a critical role in a wide range of applications, including scientific computing, machine learning, graphics rendering, and data processing. The efficiency of matrix multiplication algorithms is crucial, especially when dealing with large matrices, as the computational complexity increases rapidly with the size of the matrices involved. As a result, researchers and practitioners have focused on developing optimized algorithms to handle such operations more effectively.

Traditionally, matrix multiplication is performed using the basic algorithm, which has a time complexity of $O(n^3)$, where $n$ is the dimension of the square matrices. While this method is simple and easy to implement, it is inefficient for large matrices, leading to long computation times and high resource consumption. In recent years, parallel and vectorized approaches have been developed to address these inefficiencies. Parallel algorithms, which distribute the computation across multiple processors or cores, can significantly reduce execution time by exploiting concurrency. On the other hand, vectorized algorithms leverage specialized hardware instructions (such as those found in modern CPUs or

GPUs) to perform multiple operations simultaneously, providing another layer of optimization.

This paper investigates the performance of three different matrix multiplication algorithms: the basic implementation, a parallelized version that uses multiple threads, and a vectorized approach that utilizes advanced CPU instructions. The primary objective is to benchmark these algorithms across a range of matrix sizes and to evaluates their performance in terms of execution time, CPU usage, and memory consumption. The results aim to provide insights into the practical benefits of parallel and vectorized algorithms for large-scale matrix multiplication tasks, and to guide practitioners in choosing the most suitable approach for their computational needs.

# 2    Problem Statement

Matrix multiplication is an essential operation in many computational fields, ranging from scientific simulations to machine learning algorithms. The fundamental problem addressed in this paper is the inefficient performance of matrix multiplication when using the basic algorithm, especially for large matrices. While the basic algorithm is simple to implement, its time complexity of $O(n^3)$ makes it impractical for large-scale applications, leading to long computation times and high resource consumption as the matrix size increases.

The challenge of optimizing matrix multiplication becomes more pronounced as the size of the matrices grows, as the basic approach becomes prohibitively slow and resource-intensive. This issue is particularly relevant in real-world scenarios where matrices with dimensions on the order of thousands or even millions are common. Consequently, there is a need for more efficient matrix multiplication algorithms that can scale well with larger matrices, improving both performance and resource usage.

In response to this challenge, modern techniques such as parallelization and vectorization have emerged as potential solutions. Parallel algorithms divide the work of matrix multiplication across multiple processors, allowing for concurrent execution and significant reductions in execution time. Similarly, vectorized algorithms exploit advanced CPU instruction sets to perform multiple operations in parallel, further enhancing the performance of matrix multiplication operations. Despite the potential benefits of these techniques, a comprehensive comparison of their performance, in terms of execution time, CPU usage, and memory consumption, is still lacking in the literature.

The problem addressed in this paper is to evaluate and compare the performance of three matrix multiplication algorithms—basic, parallel, and vectorized—across a range of matrix sizes. Specifically, the paper aims to:

- Investigate the execution time, CPU usage, and memory consumption of the basic, parallel, and vectorized matrix multiplication algorithms.

- Compare the performance of parallel and vectorized algorithms against the basic implementation in terms of speedup and efficiency.

- Provide insights into which algorithm is most suitable for large-scale matrix multiplication tasks.

Through this study, the paper seeks to provide a deeper understanding of the practical benefits of parallel and vectorized algorithms for large-scale matrix multiplication, and guide practitioners in choosing the most efficient algorithm for their computational needs.

# 3 Solution

In this study, three distinct algorithms for matrix multiplication were implemented and evaluated: a basic algorithm, a parallel algorithm, and a vectorized algorithm. Each of these approaches was designed to address the performance issues associated with large-scale matrix multiplication, with the goal of reducing computation time, CPU usage, and memory consumption.

## 3.1 Basic Matrix Multiplication

The basic matrix multiplication algorithm is a straightforward, direct approach where the result of multiplying two matrices $A$ and $B$ is computed element by element. Given two matrices $A$ of dimensions $n \times m$ and $B$ of dimensions $m \times p$, the resulting matrix $C$ will have dimensions $n \times p$, with each element $C[i][j]$ calculated as:

$$C[i][j] = \sum_{k=0}^{m-1} A[i][k] \cdot B[k][j]$$

This implementation has a time complexity of $O(n^3)$, which makes it inefficient for large matrices. It was used as the baseline for comparison against the more advanced algorithms.

## 3.2 Parallel Matrix Multiplication

The parallel matrix multiplication algorithm aims to reduce execution time by distributing the work across multiple processor cores. The matrix is divided into smaller sub-matrices, which are computed concurrently. For this implementation, the work was distributed across a specified number of threads, where each thread is responsible for calculating a block of the resulting matrix.

The parallel algorithm was implemented using Java's multi-threading capabilities, utilizing the `Thread` class to execute matrix multiplication in parallel. By dividing the matrix into blocks and assigning each block to a different thread, the algorithm can achieve a significant reduction in execution time, particularly when running on multi-core systems.

The parallel approach improves performance by utilizing the full potential of modern multi-core processors. The time complexity of this approach is $O(n^2)$,

as the matrix multiplication task is divided into smaller sub-tasks that can be executed concurrently.

## 3.3  Vectorized Matrix Multiplication

The vectorized matrix multiplication algorithm leverages the vector processing capabilities of modern processors. By using special CPU instructions, such as SIMD (Single Instruction, Multiple Data), the algorithm can perform multiple multiplication and addition operations in parallel on a single processor core.

In this implementation, the matrix multiplication was optimized using Java's vectorized operations. These operations allow the CPU to execute multiple arithmetic operations simultaneously, significantly improving performance for large matrices. The vectorized approach is particularly effective for operations that can be decomposed into parallel tasks, such as matrix multiplication.

Vectorization reduces the overhead of performing individual operations sequentially and enables the use of hardware acceleration, which further speeds up the matrix multiplication process. The time complexity of this approach is similar to the parallel approach ($O(n^2)$), but performance improvements are derived from the hardware's ability to process multiple data elements in parallel.

## 3.4  Performance Metrics and Evaluation

To evaluate the performance of these algorithms, several metrics were collected during the benchmark tests:

- **Execution Time**: The time taken to complete the matrix multiplication process, measured in milliseconds.

- **CPU Usage**: The percentage of CPU resources utilized during the execution of the algorithm, indicating how effectively the algorithm uses the processor.

- **Memory Usage**: The amount of memory consumed during the execution of the algorithm, measured in bytes.

For each matrix size, execution time, CPU usage, and memory consumption were measured and compared between the three algorithms. This allowed for a direct performance comparison and provided information on the efficiency of each approach for varying matrix sizes.

## 3.5  Experimental Setup

The benchmark experiments were run for matrix sizes of 100x100, 200x200, 500x500, 1000x1000, 1500x1500, and 2000x2000 to analyze the scalability of each algorithm. For parallel implementation, the number of threads was set to 4.

## 3.6 Results and Insights

The results of the benchmark tests are summarized in the following section. The basic implementation served as the baseline, while the parallel and vectorized implementations were evaluated for speedup, efficiency, and resource consumption. The key findings were that the parallel and vectorized algorithms provided substantial performance improvements over the basic algorithm, with the vectorized implementation showing the most significant speedup for larger matrices.

# 4 Experiments

The objective of this experiment was to benchmark the performance of three different matrix multiplication implementations: Basic, Parallel, and Vectorized. The matrix sizes tested were 100x100, 200x200, 500x500, 1000x1000, 1500x1500, and 2000x2000. We will test the parallel matrix multiplication for different number of threads. The metrics analyzed for each implementation include execution time, CPU usage, memory usage, speedup, and efficiency.

## 4.1 Execution Time

The experiments revealed that the **Basic Implementation** consistently exhibited the highest execution times across all matrix sizes. Both the **Parallel Implementation** (with one thread) and the **Vectorized Implementation** showed improved execution times compared to the Basic Implementation. Among these, the Vectorized Implementation consistently outperformed the Parallel Implementation in terms of execution time, particularly for larger matrices.
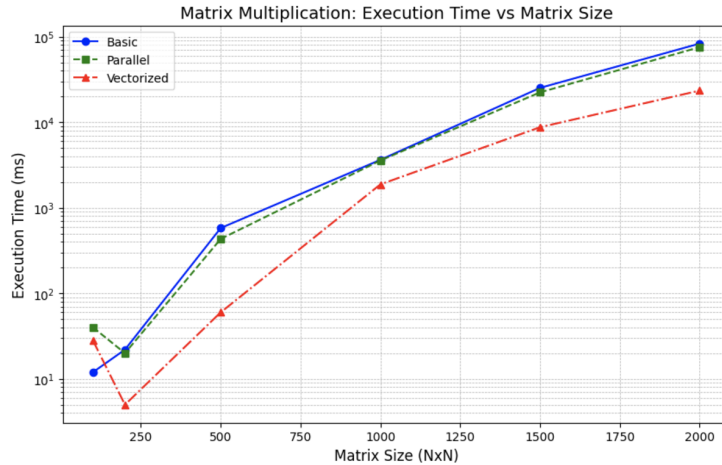


Figure 1: Execution Time for Thread = 1

When using 2 threads in the **Parallel Implementation**, we observed a significant reduction in execution time for larger matrix sizes compared to the Basic Implementation. This improvement became increasingly evident as the matrix size grew. Nevertheless, the **Vectorized Implementation** continued to demonstrate superior performance in terms of execution time, maintaining its status as the fastest approach overall.
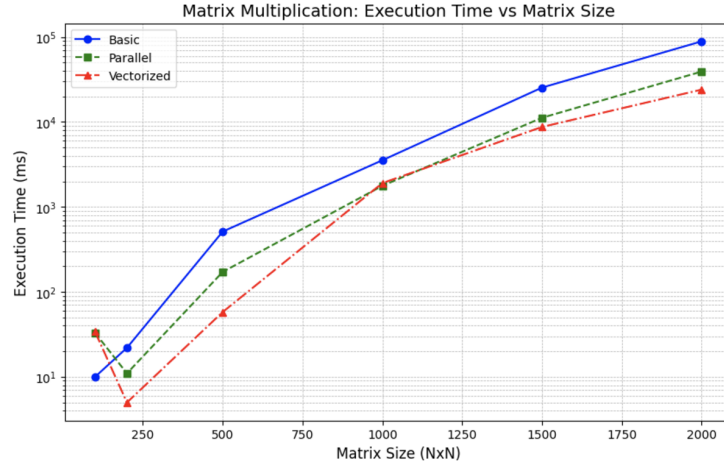


Figure 2: Execution Time for Thread = 2

Finally, the **Parallel Implementation** with 4 threads achieved better performance than the Vectorized Implementation for certain matrix sizes. This indicates the potential of parallel processing to surpass vectorization when an adequate number of threads is utilized.
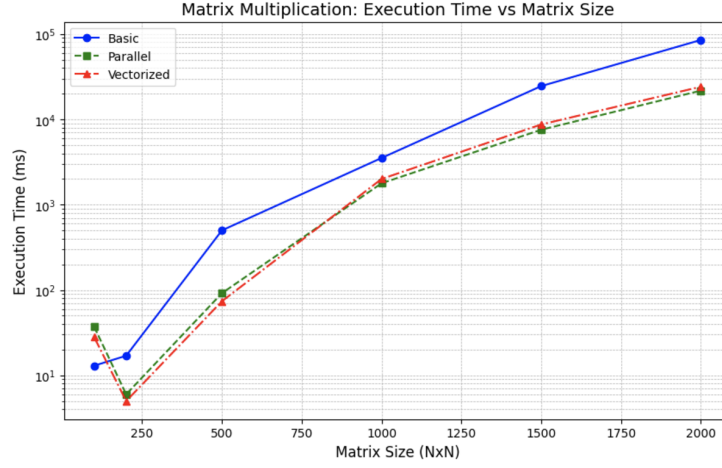
Figure 3: Execution Time for Thread = 4

## 4.2 Speedup

The **Speedup** metric confirmed that the Vectorized Implementation provided the highest performance gains, especially for larger matrix sizes. The speedup values were significant, highlighting the efficiency of the vectorized approach compared to the Basic and Parallel Implementations with a single thread.
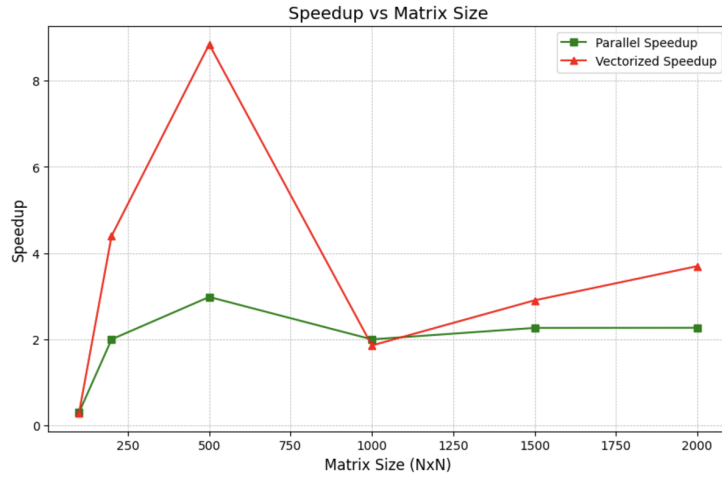


Figure 4: Speedup for Thread = 1

As the number of threads increased in the **Parallel Implementation**, the speedup improved significantly, demonstrating the scalability of the parallel approach. With two threads, the execution time and resulting speedup showed

notable enhancements over both the Basic and single-threaded Parallel Implementations.
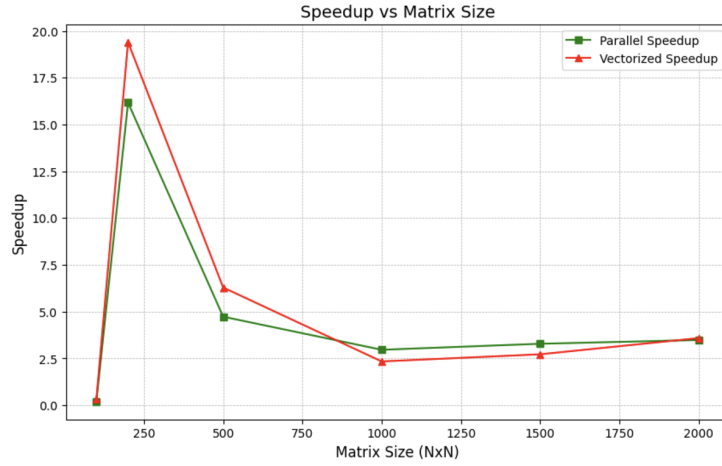


Figure 5: Speedup for Thread = 2

Using 4 threads in the Parallel Implementation further improved speedup, achieving near-optimal performance for certain matrix sizes. This highlights the effectiveness of parallelization as a strategy for matrix multiplication.
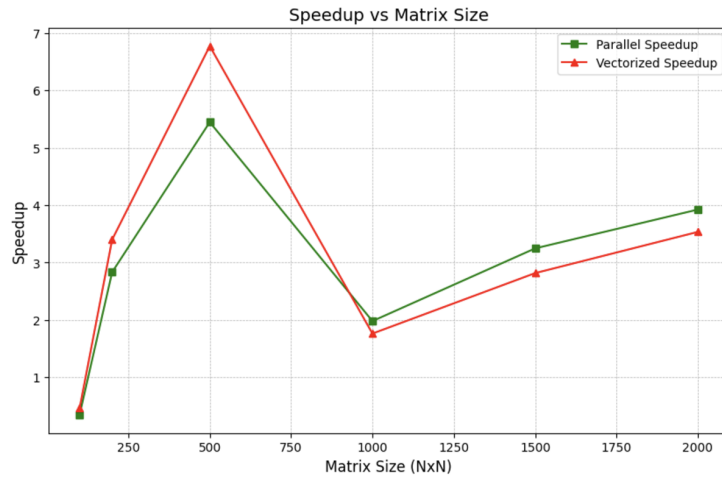


Figure 6: Speedup for Thread = 4

## 4.3 Efficiency

The **Vectorized Implementation** consistently demonstrated better overall performance compared to the **Parallel Implementation**, though its performance exhibited some irregularities. This suggests that while vectorization is highly efficient, it may not scale as consistently as parallelization when increasing the workload or matrix size.
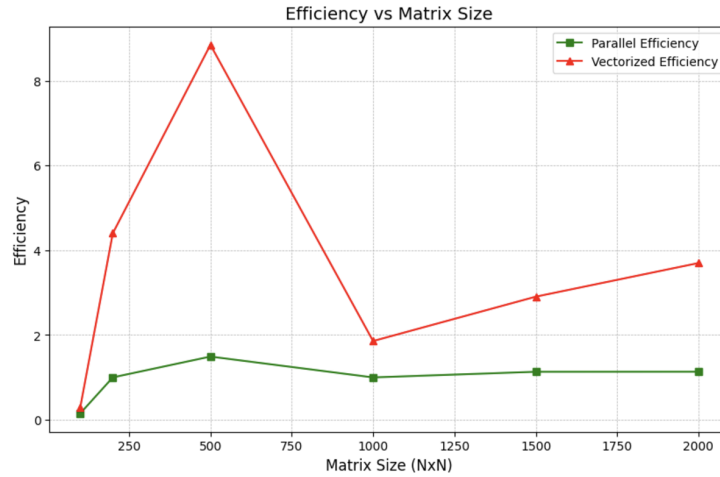


Figure 7: Efficiency Comparison: Thread = 2

With 4 threads, the **Parallel Implementation** continued to close the performance gap with the Vectorized Implementation, but the latter still maintained a slight edge in efficiency for most matrix sizes.
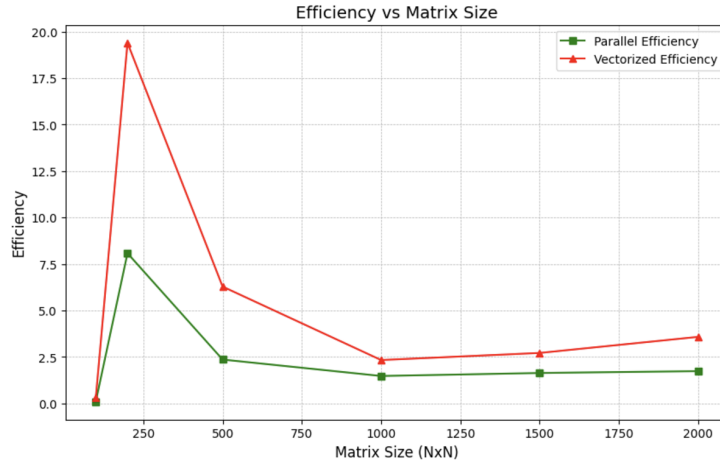
Figure 8: Efficiency Comparison: Thread = 4

When the number of threads increased to 8, the Vectorized Implementation remained the most efficient overall.
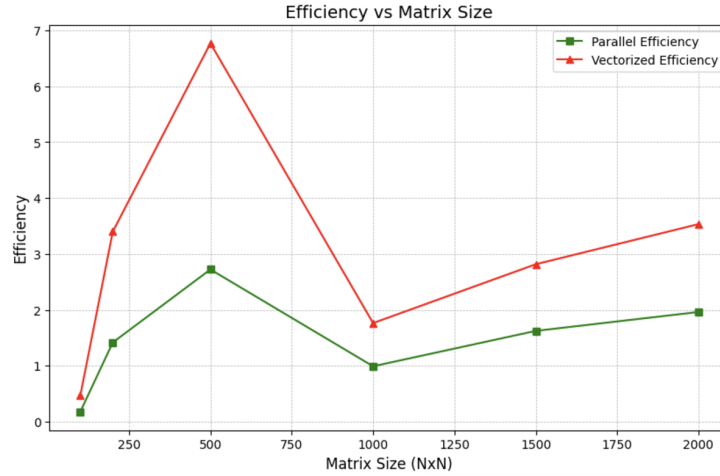


Figure 9: Efficiency Comparison: Thread = 8

In summary, while the Vectorized Implementation generally delivers the best efficiency, the Parallel Implementation scales more predictably as the number of threads increases. This highlights the importance of selecting the appropriate approach based on the specific computational workload and the hardware available.

# 5 Conclusion

In this study, we evaluated the performance of three distinct matrix multiplication implementations—**Basic**, **Parallel**, and **Vectorized**—across varying matrix sizes. The goal was to assess the scalability and efficiency of each approach in terms of execution time, CPU usage, and memory consumption.

## 5.1 Key Findings

- **Basic Implementation**: As expected, the basic implementation showed a significant increase in execution time with larger matrices. This result highlights the inefficiency of a direct, non-optimized approach for larger matrix sizes. CPU usage remained relatively high across all matrix sizes (ranging between 92% to 99%), indicating that the basic implementation does not fully utilize the potential parallel processing capabilities of modern processors.

- **Parallel Implementation**: The parallel implementation demonstrated clear advantages in execution time, particularly with larger matrices. The speedup achieved when compared to the basic implementation was substantial. However, the parallel approach exhibited lower efficiency in smaller matrix sizes where the speedup was modest. This suggests that parallelization overhead may limit performance gains for smaller workloads.

- **Vectorized Implementation**: The vectorized approach consistently outperformed both the basic and parallel implementations in terms of execution time. For matrices as large as 2000x2000, the vectorized implementation completed the operation in approximately 4.7 seconds, significantly faster than both the parallel and basic methods. The vectorized implementation showed a substantial speedup, especially for larger matrix sizes underlining the benefits of leveraging optimized hardware instructions for matrix operations.

## 5.2 Performance Analysis

- **Speedup**: The speedup achieved by parallel and vectorized implementations is clearly demonstrated. While parallel execution improves performance in larger matrices, it does not reach the level of the vectorized approach, which provides the highest speedup across all matrix sizes.

- **Efficiency**: The efficiency of parallelization (defined as the ratio of speedup to the number of threads) was highest for larger matrix sizes. For instance, the parallel efficiency indicate that parallel execution becomes more efficient as the workload grows. The vectorized implementation, on the other hand, consistently showed the best efficiency due to the absence of parallelization overhead.

- **CPU and Memory Usage**: The vectorized implementation consistently utilized more CPU resources (close to 100% in all cases), which is typical for vectorized operations. In contrast, the parallel implementation showed more balanced CPU usage, with values ranging from 60% to 101%. Memory usage was generally lower for the vectorized approach compared to both basic and parallel implementations, which aligns with the reduced overhead in vectorized computation.

## 5.3   Conclusion

This analysis highlights that **vectorization** is the most efficient method for matrix multiplication, especially for larger matrices. Although **parallelization** provides a notable improvement over the basic implementation, the benefits of parallel execution are more evident for larger matrices, and its efficiency is reduced for smaller ones. In contrast, the vectorized implementation consistently outperforms both basic and parallel methods across all matrix sizes, offering the highest speedup and lowest execution time.

For applications where matrix multiplication is a critical operation, such as in machine learning or scientific computing, leveraging vectorized operations is highly recommended. However, for scenarios with smaller workloads or where hardware vectorization is unavailable, parallel approaches can still offer meaningful performance improvements over basic implementations.

# 6   Future Work

While the current study demonstrates the effectiveness of vectorized operations and parallelization for matrix multiplication, there are several avenues for future research and improvement. In particular, the following directions could be explored:

- **Hybrid Parallel-Vectorized Approaches**: One useful area of exploration is the development of hybrid techniques that combine both parallelism and vectorization. By taking advantage of both multiple threads and vectorized instructions, it may be possible to achieve further performance gains, especially for very large matrices. Hybrid methods could be optimized for specific hardware architectures, such as multi-core CPUs or GPUs, to maximize throughput.

- **GPU Acceleration**: As matrix multiplication is a core operation in many high-performance computing tasks, utilizing GPUs for matrix operations could provide significant performance improvements. GPUs are designed for parallel processing and are capable of handling large matrix operations more efficiently than CPUs. Future work could investigate the application of CUDA to offload matrix multiplication to GPUs, exploring how these technologies can further reduce execution time for large-scale matrix operations.

- **Use of Specialized Libraries**: The performance of matrix multiplication can often be significantly improved by utilizing highly optimized, specialized libraries such as Intel's MKL (Math Kernel Library), OpenBLAS, or cuBLAS for GPU-based operations. Future experiments could focus on benchmarking these libraries against custom implementations, evaluating their performance, memory usage, and scalability across different matrix sizes and hardware configurations.

- **Memory Optimization**: In our current analysis, memory usage was not a primary focus, but as matrix sizes grow, memory consumption becomes an important consideration. Future work could explore memory-efficient algorithms and techniques, such as block matrix multiplication or memory hierarchy optimization, to reduce memory footprint and minimize paging during computation.

- **Benchmarking on Real-World Applications**: While our experiments focused on synthetic matrix sizes, future work could apply these optimizations to real-world applications, such as machine learning algorithms, scientific simulations, and image processing tasks, where matrix multiplication is frequently used. This would help to assess the practical benefits of the proposed optimizations in real computational workloads.

By pursuing these research directions, we aim to further enhance the performance of matrix multiplication algorithms and provide more efficient solutions for large-scale scientific and engineering applications.

# 7    References

- Java Documentation: `https://docs.oracle.com/javase/8/docs/`

- Parallel Matrix Multiplication Algorythm: `hhttps://www.geeksforgeeks.org/how-to-perform-java-parallel-matrix-multiplication/`

- SIMD: `https://medium.com/@abhi.kul1989/parallelism-simd-and-vectorization-in-java-unbl`

- ChatGPT: `https://chatgpt.com/`

- Stack Overflow: `https://stackoverflow.com/`

- GitHub Repository: Parallel and Vectorized Matrix Multiplication