

Parallel and Vectorized Matrix Multiplication

David García Vera

November 2024

Source code in this [GitHub repository](#).

Abstract

Matrix multiplication is a fundamental operation in computational science, yet its high computational cost demands optimization for practical applications. This paper investigates two approaches to optimize matrix multiplication: parallelization and vectorization. The parallel implementation leverages Java's 'ForkJoinPool' framework to distribute computational workloads across multiple threads, while the vectorized implementation utilizes the 'jdk.incubator.vector' package to exploit SIMD (Single Instruction, Multiple Data) instructions for processing data in parallel. Both methods were benchmarked against a baseline sequential algorithm, focusing on metrics such as execution time, speedup, and resource utilization.

Results indicate that the parallel implementation significantly reduces execution time for large matrices, achieving substantial speedup through multi-threading. The vectorized approach demonstrated competitive performance for smaller matrices, offering efficient memory utilization and leveraging SIMD capabilities. The basic algorithm, while straightforward, was impractical for large datasets due to its sequential nature and high computational cost.

This study highlights the trade-offs between parallel and vectorized optimization strategies, emphasizing their suitability for different computational scenarios. By leveraging modern hardware capabilities, these methods provide valuable insights and practical tools for enhancing matrix multiplication performance in diverse applications.

1 Introduction

Matrix multiplication is a core operation in various fields such as numerical simulations, machine learning, and computer vision. The standard implementation of matrix multiplication, which involves triple nested loops, has a computational complexity of $O(n^3)$ for square matrices of size $n \times n$. While this approach is straightforward to implement, it becomes computationally prohibitive as the size of the matrices grows, making optimization a necessity.

Modern hardware architectures provide opportunities to optimize matrix multiplication by exploiting parallelism and vectorization. Parallelism leverages multiple cores to distribute the computational load, while vectorization uses SIMD (Single Instruction, Multiple Data) instructions to process multiple data elements simultaneously. These optimizations can significantly improve performance, but they also require careful algorithmic design to maximize the utilization of computational resources and minimize overheads such as memory latency.

In this paper, we implement and analyze two optimized approaches to matrix multiplication:

1. A **parallelized implementation** using the Fork/Join framework in Java, which divides the matrix into smaller tasks that can be executed concurrently across multiple threads.
2. A **vectorized implementation** using the `FloatVector` class from the `jdk.incubator.vector` package to exploit SIMD instructions, allowing operations on chunks of data in a single instruction cycle.

Both implementations are compared to the basic triple-loop algorithm, focusing on metrics such as speedup, parallel efficiency, and resource utilization. The results of this study contribute to understanding the trade-offs and practical considerations in optimizing matrix multiplication for modern hardware.

2 Problem Statement

Matrix multiplication poses a computational bottleneck due to its inherently high complexity. In the naïve approach, the process involves computing each element of the resulting matrix as the dot product of a row and a column. While conceptually simple, this approach fails to leverage the capabilities of modern multi-core and SIMD-enabled processors, leading to suboptimal performance.

To address these limitations, this paper explores two specific optimizations:

1. **Parallelization:** Using Java’s ‘ForkJoinPool’, matrix multiplication is divided into smaller, independent tasks that can be executed concurrently. The threshold for dividing tasks is dynamically adjusted to balance workload and minimize overhead. The matrix transposition is also utilized to improve cache efficiency during computation.
2. **Vectorization:** Leveraging the ‘FloatVector’ class from the ‘jdk.incubator.vector’ package, the vectorized implementation processes chunks of data simultaneously using SIMD instructions. This approach preprocesses matrices to align with the preferred vector size and handles tail cases to ensure correctness.

The study’s objectives include:

1. Implementing both parallelized and vectorized versions of matrix multiplication.
2. Comparing their performance against the baseline implementation in terms of speedup and efficiency.
3. Evaluating the scalability of each approach by varying matrix sizes and hardware resources.

These implementations illustrate how modern programming tools and hardware features can be combined to optimize computationally intensive tasks such as matrix multiplication, providing practical insights for real-world applications.

3 Methodology

The methodology for optimizing matrix multiplication is divided into two distinct approaches: parallelization and vectorization. These implementations were designed to leverage modern hardware capabilities effectively, including multi-core processors and SIMD (Single Instruction, Multiple Data) units. Both implementations were compared against a baseline sequential algorithm to evaluate their performance and scalability.

3.1 Baseline Implementation

The baseline implementation follows the traditional triple-loop algorithm, where the resulting matrix is computed as:

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

for all $i, j \in [0, n - 1]$. This approach, while simple and intuitive, has a computational complexity of $O(n^3)$ and does not utilize parallelism or SIMD instructions. It was implemented in Java to serve as a reference point for performance comparisons.

3.2 Parallel Implementation

The parallel implementation was developed using Java’s ‘ForkJoinPool’ framework. This approach recursively divides the computation into smaller sub-tasks, which are executed concurrently by multiple threads. The primary steps are as follows:

1. **Matrix Transposition:** To improve cache efficiency during computation, the second matrix (B) is transposed before multiplication. This ensures that matrix elements are accessed sequentially in memory, reducing cache misses.
2. **Recursive Task Division:** The computation is divided into smaller sub-problems based on row ranges of the result matrix (C). If the number of rows in the sub-problem exceeds a predefined threshold (e.g., 200), the task is split further.
3. **Sequential Computation for Small Tasks:** For sub-problems below the threshold, the matrix multiplication is performed sequentially to minimize the overhead of task management.
4. **Execution in ‘ForkJoinPool’:** The tasks are submitted to a shared thread pool, allowing efficient utilization of available cores.

The parallel implementation effectively reduces computation time by distributing the workload across multiple threads. However, its performance is influenced by factors such as the threshold value, matrix size, and the number of available cores.

3.3 Vectorized Implementation

The vectorized implementation utilizes the ‘jdk.incubator.vector’ package, which provides access to SIMD instructions through the ‘FloatVector’ class. The algorithm follows these steps:

1. **Matrix Preparation:** Similar to the parallel implementation, matrix B is transposed to align memory access patterns with SIMD operations.
2. **Chunk Processing:** The computation for each element of the result matrix ($C[i][j]$) is divided into chunks of size equal to the preferred vector species (e.g., 128-bit or 256-bit, depending on the hardware). This ensures maximum utilization of SIMD registers.
3. **Vector Operations:** For each chunk, the corresponding rows of A and columns of B are loaded as ‘FloatVector’ objects. The dot product of these vectors is computed using element-wise multiplication followed by reduction.
4. **Tail Handling:** Elements that do not fit into a full SIMD vector are processed sequentially to ensure correctness.

This approach accelerates computation by processing multiple elements in a single instruction cycle. Its performance depends on the alignment of data with the SIMD vector size and the efficiency of the underlying hardware.

3.4 Experimental Configuration

The implementations were tested on large square matrices with dimensions ranging from 512×512 to 4096×4096 . The following hardware and software configurations were used:

- **Processor:** Multi-core CPU with support for SIMD instructions.
- **Memory:** Sufficient RAM to handle matrix sizes without disk swapping.
- **JVM:** Java 21, with access to the ‘jdk.incubator.vector’ module.

Performance metrics such as execution time, speedup, and parallel efficiency were measured. Each experiment was repeated multiple times, and the average performance was recorded to account for variability in runtime conditions.

3.5 Comparison and Analysis

To provide a comprehensive evaluation, the following comparisons were performed:

1. **Baseline vs. Parallel Implementation:** To measure the speedup achieved by multi-threading.
2. **Baseline vs. Vectorized Implementation:** To evaluate the benefit of SIMD optimizations.
3. **Parallel vs. Vectorized Implementations:** To identify the strengths and weaknesses of each approach.

These comparisons highlight the trade-offs between parallelization and vectorization, including their impact on computation time, resource utilization, and scalability.

4 Experiments

This section presents the results of benchmarking the three implementations: the basic sequential algorithm, the parallelized implementation, and the vectorized implementation. The experiments were conducted on square matrices of various sizes ranging from 100×100 to 3000×3000 . Each implementation was evaluated in terms of average execution time and memory usage.

Each test was repeated 10 times, and the average execution time and memory usage were recorded. Confidence intervals (99.9%) were calculated to quantify variability.

4.1 Results

The results for average execution time (ms/op) and memory usage (MB) for the three implementations are summarized in Tables 1 and 2, respectively.

4.1.1 Average Execution Time

Table 1 lists the average execution times for matrix sizes $n = 100, 500, 1000, 2000$, and 3000 . Figure 1 shows the corresponding graphical representation.

Method	$n = 100$	$n = 500$	$n = 1000$	$n = 2000$	$n = 3000$
Basic	0.970	191.193	3856.535	56091.469	235173.618
Parallel	0.842	31.574	149.939	1318.025	4693.941
Vectorized	0.385	35.479	259.552	2695.634	9836.863

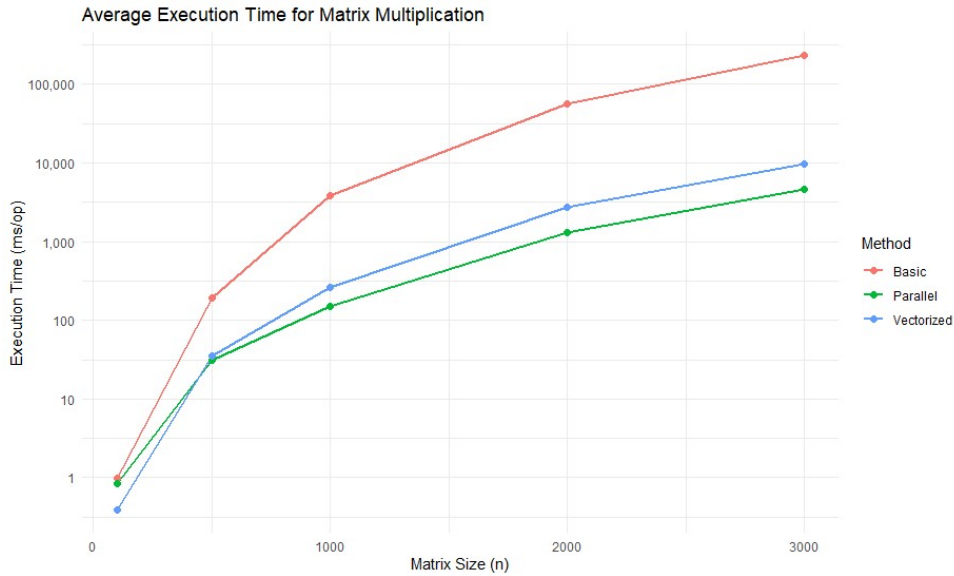


Figure 1: Average Execution Time (ms/op) for Matrix Multiplication

4.1.2 Memory Usage

Table 2 shows the memory usage during execution, and Figure 2 provides a graphical visualization of the results.

Table 2: Memory Usage (MB) for Matrix Multiplication

Method	n = 100	n = 500	n = 1000	n = 2000	n = 3000
Basic	0.06	1.91	7.61	30.62	69.35
Parallel	0.13	3.87	15.35	62.49	141.22
Vectorized	0.61	3.12	8.22	31.39	69.30

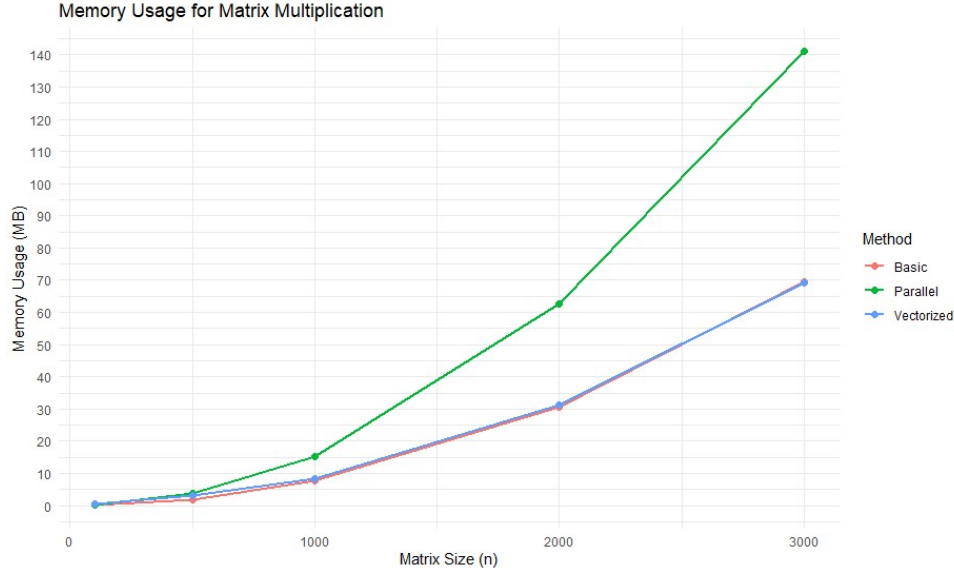


Figure 2: Memory Usage (MB) for Matrix Multiplication

4.2 Discussion

The results demonstrate significant differences in performance between the three methods:

- The basic algorithm exhibited the highest execution time and lowest memory usage due to its simplicity and lack of optimizations.
- The parallel implementation achieved a substantial reduction in execution time, particularly for larger matrices. However, this came at the cost of increased memory usage due to the overhead of task management and transposition.
- The vectorized implementation showed competitive execution times for smaller matrices, benefiting from SIMD optimizations. Memory usage was slightly higher than the basic algorithm but lower than the parallel implementation.

These findings highlight the trade-offs between execution time and resource usage, providing valuable insights for selecting the appropriate optimization strategy based on specific application requirements.

5 Conclusions

Matrix multiplication is a fundamental operation in numerous computational applications, yet its high computational cost makes optimization a necessity, particularly as problem sizes scale. This study explored two approaches to optimization—parallelization and vectorization—using Java’s ‘ForkJoinPool’ framework and the ‘jdk.incubator.vector’ package, respectively. These methods were benchmarked against a baseline sequential algorithm to evaluate their effectiveness in reducing execution time and managing resource usage.

The results of this study provide clear evidence of the potential performance gains achievable through modern hardware-aware optimizations:

- The **parallel implementation** excelled in reducing execution time for larger matrices by utilizing multi-threading. It demonstrated a substantial speedup compared to the baseline algorithm, particularly as matrix sizes increased. This result highlights the benefits of effectively distributing computational workloads across multiple CPU cores. However, the implementation required higher memory usage due to task management overhead and matrix transposition, which may become a limiting factor in memory-constrained environments.
- The **vectorized implementation**, while not as fast as the parallel approach for larger matrices, showed highly competitive performance for smaller datasets. By exploiting SIMD instructions, it achieved significant reductions in execution time compared to the baseline. Additionally, the vectorized approach had a more efficient memory profile compared to the parallel method, making it a suitable choice for applications where memory efficiency is critical.
- The **baseline implementation**, though conceptually simple, was inadequate for large-scale problems. Its sequential nature resulted in prohibitive execution times as matrix sizes grew, underscoring the necessity of employing modern optimization techniques.

These findings reveal important trade-offs inherent in optimization strategies for matrix multiplication. The parallel approach is particularly well-suited for large datasets in scenarios where execution time is the primary constraint, while the vectorized implementation offers an efficient alternative for smaller matrices or memory-sensitive applications. Together, these approaches demonstrate how algorithm design can be tailored to leverage the specific strengths of modern CPU architectures.

Beyond the numerical results, this study underscores the importance of understanding the underlying hardware and software environment when optimizing computational workloads. The effectiveness of the parallel and vectorized implementations depended not only on the algorithms themselves but also on how well they aligned with the architectural features of the processor, such as the number of cores and SIMD register sizes.

Furthermore, the study emphasizes the role of benchmarking and empirical testing in evaluating algorithmic optimizations. By systematically analyzing metrics such as execution time, speedup, and memory usage, this work provides actionable insights for both developers and researchers seeking to enhance the performance of matrix operations. The methodologies and results presented here serve as a foundation for further exploration of advanced optimization techniques in computational mathematics.

In conclusion, optimizing matrix multiplication is a multidimensional problem that requires balancing execution speed, memory usage, and scalability. The parallel and vectorized implementations presented in this study offer two distinct yet complementary strategies for addressing these challenges, providing valuable tools for computational tasks across a wide range of disciplines.

6 Future Works

As matrix multiplication scales to increasingly larger datasets, single-machine optimizations such as parallel and vectorized implementations encounter inherent limitations. Modern computational problems often involve matrices so large that they exceed the memory capacity of a single machine, necessitating a shift toward distributed computing. Future work will explore distributed matrix multiplication as a solution to these challenges, leveraging distributed systems to handle extremely large datasets while maintaining efficiency and scalability.

Distributed matrix multiplication partitions matrices across multiple nodes in a computing cluster, enabling concurrent computation of submatrices. This approach not only addresses the memory constraints of single machines but also introduces new dimensions to performance optimization, including data distribution strategies and network communication overhead. A key focus of future work will be the implementation of matrix multiplication using a distributed computing framework, such as Apache Spark or MPI, to evaluate how effectively distributed systems can process massive matrices.

Scalability will be a critical metric for assessing the success of the distributed approach. Experiments will analyze how performance changes as the size of the matrices increases, comparing the distributed method against the basic sequential and parallel implementations. Additionally, the study will investigate the impact of network overhead and data transfer times, as these factors often become bottlenecks in distributed systems. By examining these aspects, we aim to identify strategies to minimize communication latency and optimize data movement across nodes.

Another important dimension of future work will be resource utilization. The distributed approach inherently involves multiple nodes, each contributing computational power and memory. Understanding how resources are utilized—such as the number of nodes required and memory consumption per node—will provide insights into the cost-effectiveness of distributed matrix multiplication. This analysis will help in determining the trade-offs between hardware scalability and overall system efficiency.

By exploring these aspects, future research will provide a comprehensive understanding of how distributed computing frameworks can address the challenges of large-scale matrix multiplication. This work will not only expand the scope of matrix multiplication optimization but also offer practical guidance for deploying distributed systems in real-world computational scenarios.