

Analysis and Benchmarking Report of Matrix Operations

1. Introduction

This report compares five baseline implementations:

- **Matrix-Vector Multiplication:**
 - `multiply_mv_row_major`
 - `multiply_mv_col_major`
- **Matrix-Matrix Multiplication:**
 - `multiply_mm_naive`
 - `multiply_mm_transposed_b`
 - `multiply_mm_tiled`

Each function has been benchmarked using a custom framework based on `std::chrono` with various problem sizes (small, medium, and large square matrices). In addition to presenting raw benchmark data, we analyze memory access patterns (cache locality), discuss the potential impact of memory alignment and inlining, and summarize profiling insights. Finally, we outline potential optimization strategies.

2. Benchmarking Framework & Results

2.1 Benchmarking Framework

A benchmarking framework was developed in C++ using `std::chrono` to measure execution time. Key characteristics include:

- **Multiple Runs:** Each function is executed multiple times (e.g., 5 runs) to compute the average execution time and standard deviation.
- **Variable Problem Sizes:** Benchmarks were run on small (100×100), medium (500×500), and large (1000×1000) matrices.

- **Preventing Compiler Optimization:** A global sum (`g_sum`) is updated with results to avoid code elimination.
- **Clear Output:** Results are printed to the console and can be summarized in tables or graphs.

2.2 Benchmark Results

Below is a summary table based on the provided sample output:

Function	Small (N=100)	Medium (N=500)	Large (N=1000)
multiply_mv_row_major	48 μ s (avg)	336 μ s (avg)	1052 μ s (avg)
multiply_mv_col_major	41.6 μ s (avg)	337.8 μ s (avg)	1041.8 μ s (avg)
multiply_mm_naive	3121.4 μ s (avg)	142941 μ s (avg)	1.24962 $\times 10^6$ μ s (avg)
multiply_mm_transposed_b	2125.8 μ s (avg)	142714 μ s (avg)	1.16019 $\times 10^6$ μ s (avg)
multiply_mm_tiled	1889.6 μ s (avg)	144811 μ s (avg)	1.23635 $\times 10^6$ μ s (avg)

Note: Times are averages over 5 runs.

These results provide a baseline for comparing further optimizations.

3. Cache Locality Analysis

3.1 Matrix-Vector Multiplication

- **Row-Major Implementation:**
Accesses elements sequentially from each row (i.e., `matrix[i * cols + j]`), which typically leverages spatial locality. This often results in fewer cache misses when the matrix is stored in row-major order.
- **Column-Major Implementation:**
The inner loop iterates over rows for each column (accessing `matrix[j * rows + i]`), which may lead to non-sequential memory access. Despite this, the provided benchmarks show similar performance. This similarity might be due to small working sets or caching of the vector values.

3.2 Matrix-Matrix Multiplication

- **Naive Implementation:**
The inner-most loop accesses matrix A in a row-major fashion but accesses matrix B with a stride (i.e., $B[k * \text{colsB} + j]$), potentially leading to inefficient cache usage.
- **Transposed-B Implementation:**
By transposing matrix B prior to multiplication, the inner loop then accesses contiguous memory locations ($\text{matrixB_transposed}[j * \text{rowsB} + k]$), which improves cache utilization.
- **Tiled Implementation:**
The algorithm divides the matrices into blocks (tiles) that ideally fit into the cache. This blocking reduces the frequency of cache misses and improves overall performance.

Benchmarks show that for smaller matrices the transposed and tiled versions outperform the naive approach. However, as the matrix size increases, the performance differences become less pronounced, highlighting that cache effects can depend on both matrix size and system cache characteristics.

4. Memory Alignment

4.1 Impact of Alignment

Memory alignment ensures that data structures (matrices and vectors) begin at addresses that are multiples of a specific boundary (e.g., 64 bytes). This can:

- **Improve Cache Efficiency:** Aligned memory can reduce cache line splits.
- **Enable SIMD:** Aligned data is often required for efficient vectorized operations.

4.2 Experimental Considerations

To study the impact:

- **Aligned Allocation:** Modify memory allocation using functions like `aligned_alloc` (C11/C++17) or custom allocators to align data on a 64-byte boundary.
- **Benchmark Comparison:** Run the same benchmarks using aligned and unaligned memory to quantify improvements. Performance gains are most significant in memory-bound operations, especially with larger matrices.

5. Inlining and Compiler Optimizations

5.1 Inlining

- **Use Case:**
Inlining small, frequently called helper functions can reduce function call overhead.
- **Observation:**
For our implementations (which involve nested loops), inlining may offer marginal benefits unless the helper functions are very small and called in inner loops.

5.2 Compiler Optimizations

- **Optimization Levels:**
Compiling with aggressive optimizations (e.g., `-O3`) may:
 - Unroll loops,
 - Vectorize arithmetic operations,
 - Inline functions automatically.
- **Trade-Offs:**
While higher optimization levels generally improve performance, they can also make debugging and profiling more challenging.
- **Analysis:**
Comparing builds compiled with and without such optimizations (e.g., `-O0` vs. `-O3`) helps identify which parts of the code are most affected by these changes.

6. Optimization Strategies and Implementation

6.1 Brainstorming and Optimization

Based on the analysis, potential optimization strategies include:

- **Loop Reordering:**
Reorder inner loops to access contiguous memory.
- **Blocking/Tiling:**
The tiled implementation is an example that divides the computation into blocks to improve cache reuse.
- **Memory Alignment:**
Ensure that all matrices and vectors are aligned on cache line boundaries.

- **Vectorization:**
Explore SIMD instructions to perform operations on multiple data points simultaneously.

6.2 Implementation Example: Tiled Multiplication

The tiled multiplication (`multiply_mm_tiled`) is already a step toward improved cache reuse. Its design:

- Divides the matrices into sub-blocks based on a configurable `BLOCK_SIZE`.
- Works on smaller blocks that ideally reside in the cache, thus reducing cache misses.
- Benchmarks for the small matrix ($N=100$) indicate that it outperforms the naive approach by nearly 40% compared to the $3121.4\ \mu s$ (naive) versus $1889.6\ \mu s$ (tiled).

6.3 Benchmarking the Optimized Version

After implementing an optimization (such as memory alignment or further loop unrolling/vectorization), re-run the benchmarks to quantify improvements. Document:

- **Performance Gain:** Percentage improvement over the baseline.
- **Conditions:** For which matrix sizes the improvements are most significant.
- **Trade-offs:** Any increased complexity or maintenance costs.

7. Conclusion

The baseline implementations provide a strong starting point for understanding the impact of cache locality and memory access patterns on performance. Key takeaways include:

- **Matrix-Vector Multiplication:**
Although both row-major and column-major implementations achieve the same results, row-major tends to benefit from contiguous memory accesses. The differences may be less pronounced for small matrices but become more critical as sizes grow.
- **Matrix-Matrix Multiplication:**
The transposed and tiled versions improve cache utilization by reducing stride-related cache misses. For small matrices, these optimizations can significantly reduce execution time.
- **Memory Alignment & Inlining:**
Optimizing memory alignment can further improve cache performance, and judicious use of inlining may reduce overhead in frequently called functions.

- **Profiling & Optimization:**

Profiling helps identify hotspots—primarily in the inner loops—and guides further optimizations such as blocking, vectorization, or even parallelization.

Future work should include rigorous benchmarking of aligned memory versions, experimenting with compiler flags, and exploring advanced optimizations such as auto-vectorization or multi-threading. These steps will help the team further enhance performance while balancing maintainability and complexity.