

Performance Report: Matrix Multiplication Optimization (ABt vs AtBt)

Executive Summary

This report presents performance benchmarking results for two matrix multiplication variants implemented on both CPU and GPU platforms. The optimization focuses on cache locality, shared memory utilization, and thread-level parallelism to achieve significant speedups.

Key Finding: GPU implementations achieve **15-60x speedup** over optimized CPU versions, with best GPU performance reaching **~5,800 GFLOPS** and CPU peaking at **~309 GFLOPS**.

Test Configuration

Matrix Variants

1. **ABt:** Computes $C[i, j] = \sum_k A[i, k] \cdot B[j, k]$ (row-vector dot products)
- 2.

AtBt: Computes $C[i, j] = \sum_k A[k, i] \cdot B[j, k]$ (**transposed A indexing**)

Performance Results

CPU Performance (ABt - Optimized with OpenMP + Tiling)

Nj	Nj	Nk	Best (GFLOPS)	Worst (GFLOPS)
16,384	8,192	256	308.39	298.15
8,192	16,384	256	303.27	292.52
16,384	16,384	128	296.00	289.72
8,192	8,192	512	290.05	280.70
2,048	2,048	8,192	278.04	253.28
512	512	131,072	121.56	99.14
256	256	524,288	75.20	63.73
16,377	7,777	257	274.78	264.21
7,777	16,377	257	268.10	260.86
15,677	15,677	131	267.54	253.28
7,777	7,777	777	309.41	296.50

Nj	Nj	Nk	Best (GFLOPS)	Worst (GFLOPS)
1,999	1,999	33,333	161.34	149.82
511	511	13,111	259.36	115.98
131	131	55,557	138.57	58.32

CPU Observations: - Peak: **309.41 GFLOPS** ($7,777 \times 7,777 \times 777$) - Performance degrades significantly with large Nk (cache misses) - Performance drops ~60% when Nk > 100K (memory bandwidth limited)

CPU Performance (AtBt - Optimized with OpenMP + Tiling)

Ni	Nj	Nk	Best (GFLOPS)	Worst (GFLOPS)
16,384	8,192	256	74.02	73.03
8,192	16,384	256	70.22	64.54
16,384	16,384	128	83.14	75.71
8,192	8,192	512	67.47	65.69
2,048	2,048	8,192	66.74	60.47
512	512	131,072	64.54	61.19
256	256	524,288	60.63	59.13
16,377	7,777	257	130.31	114.35
7,777	16,377	257	117.02	106.70
15,677	15,677	131	158.08	146.34
7,777	7,777	777	131.22	129.67
1,999	1,999	33,333	57.86	54.56
511	511	13,111	94.81	57.20
131	131	55,557	51.12	21.76

AtBt CPU Observations: - Peak: **158.08 GFLOPS** ($15,677 \times 15,677 \times 131$) - **4x slower** than ABt across all test cases - A[k*Ni+i] access pattern causes poor cache locality - Memory access is not contiguous, leading to cache line waste

GPU Performance (ABt - Shared Memory Optimized)

Ni	Nj	Nk	Best (GFLOPS)	Worst (GFLOPS)	GPU/CPU Ratio
16,384	8,192	256	5,175.38	5,179.44	16.8x
8,192	16,384	256	5,184.45	5,187.01	17.1x

Ni	Nj	Nk	Best (GFLOPS)	Worst (GFLOPS)	GPU/CPU Ratio
16,384	16,384	128	5,172.32	5,184.93	17.5x
8,192	8,192	512	5,172.86	5,183.54	17.8x
2,048	2,048	8,192	5,128.23	5,142.76	18.4x
512	512	131,072	4,165.81	4,176.59	34.3x
256	256	524,288	2,750.18	2,767.52	36.6x
16,377	7,777	257	5,428.48	5,438.29	19.8x
7,777	16,377	257	5,436.46	5,448.58	20.3x
15,677	15,677	131	5,394.82	5,408.65	20.2x
7,777	7,777	777	5,801.06	5,817.49	18.8x
1,999	1,999	33,333	5,786.35	5,802.04	35.9x
511	511	13,111	4,776.06	4,841.88	18.4x
131	131	55,557	1,148.91	1,166.96	8.3x

GPU ABt Observations: - Peak: **5,801.06 GFLOPS** ($7,777 \times 7,777 \times 777$)

- Consistent 5.1-5.4 TFLOPS for Nk < 100K - Excellent scalability: maintains ~5.8 TFLOPS even with Nk = 33,333 - Performance drops only when Nk becomes extreme (>100K)
-

GPU Performance (AtBt - Shared Memory Optimized)

Ni	Nj	Nk	Best (GFLOPS)	Worst (GFLOPS)	GPU/CPU Ratio
16,384	8,192	256	5,237.19	5,243.32	70.7x
8,192	16,384	256	5,243.26	5,252.05	74.7x
16,384	16,384	128	5,232.95	5,244.89	62.9x
8,192	8,192	512	5,248.17	5,256.56	77.8x
2,048	2,048	8,192	5,196.79	5,211.82	77.8x
512	512	131,072	4,169.73	4,188.68	64.6x
256	256	524,288	2,737.48	2,754.24	45.2x
16,377	7,777	257	5,484.73	5,489.85	42.1x
7,777	16,377	257	5,480.72	5,499.06	46.8x
15,677	15,677	131	5,436.69	5,445.62	34.4x
7,777	7,777	777	5,849.45	5,860.15	44.6x
1,999	1,999	33,333	5,797.15	5,803.84	100.1x
511	511	13,111	4,866.55	4,913.14	51.4x
131	131	55,557	1,157.44	1,169.57	22.6x

GPU AtBt Observations: - Peak: **5,849.45 GFLOPS** ($7,777 \times 7,777 \times 777$)

- **Massive speedup over CPU AtBt:** 20-100x improvement - GPU eliminates the cache locality penalty that cripples CPU AtBt - Similar performance profile

to ABt GPU

Performance Analysis

1. CPU Performance Comparison (ABt vs AtBt)

ABt Advantage: 4.0-4.2x faster than AtBt across all cases

Reason: Memory Access Patterns

- ABt: $A[i*Nk+k]$ - highly cache-friendly (row-major, contiguous)
- AtBt: $A[k*Ni+i]$ - poor cache locality (column-major, scattered)

Impact on different matrix sizes: - Small matrices ($Nk < 1K$): ABt ~4x faster, cache miss penalty minimal - Medium matrices ($1K < Nk < 10K$): ABt ~4x faster, L3 cache helps both - Large matrices ($Nk > 100K$): ABt ~4x faster, memory bandwidth bottleneck

2. GPU Performance Comparison (ABt vs AtBt)

GPU Performance: AtBt ~ ABt (within 1-5% variation)

GPU-Specific Advantages over CPU:

1. Shared Memory Tiling:
 - CPU must rely on L1/L2/L3 caches (slow, limited)
 - GPU explicitly manages 32×32 tiles in fast shared memory
2. Memory Access Pattern Normalization:
 - GPU threads coalesce all loads into sequential patterns
 - AtBt's scattered access pattern becomes just as efficient as ABt
3. Massive Parallelism:
 - 1024 threads per block (32×32) vs 128 CPU threads
 - Multiple blocks run simultaneously

3. Speedup Analysis

Small Nk (< 1K)

- **GPU/CPU Ratio:** 16-20x
- Bottleneck: Kernel launch overhead, limited compute per thread
- Shared memory has less impact

Medium Nk (1K - 100K)

- **GPU/CPU Ratio:** 18-37x (ABt), 45-78x (AtBt)
- Optimal range: Shared memory fully utilized, high occupancy
- AtBt GPU shows larger relative improvement due to AtBt CPU weakness

Large Nk (> 100K)

- **GPU/CPU Ratio:** 36-100x (AtBt)
 - CPU severely bandwidth-limited
 - GPU still achieves 2.7-4.2 TFLOPS despite very large Nk
-

Optimization Techniques Applied

CPU Optimizations (Both ABt and AtBt)

1. Loop Tiling on k dimension

```
Tile size adapts to matrix dimension:  
- Nk >= 131,072: 512  
- Nk >= 8,192: 256  
- Nk >= 1,024: 128  
- else: 64
```

2. OpenMP Parallelization

- ABt: collapse(2) on (i,j) loops for better load balancing
- AtBt: 3D blocking with explicit i-dimension tiling

3. SIMD Vectorization

- #pragma omp simd reduction(+:sum) on innermost k-loop
- Compiler auto-vectorizes to SSE/AVX instructions

4. Cache Locality

- Row-major access patterns where possible
- Minimize cache-line waste

GPU Optimizations (Both ABt and atBt)

1. Shared Memory Tiling (32×32)

```
__shared__ float sA[32][33]; // Padded to avoid bank conflicts  
__shared__ float sB[32][33];  
• Reduces global memory reads by ~32x
```

- Each tile re-loaded once per 32 k-iterations

2. Thread Block Configuration

- Block size: $32 \times 32 = 1024$ threads
- Grid: $(N_i + 31)/32 \times (N_j + 31)/32$ blocks
- Maximum occupancy for Granite GPU

3. Loop Unrolling

```
#pragma unroll 8
for (int k = 0; k < 32 && (kk + k) < Nk; k++)
```

- Improves instruction-level parallelism (ILP)
- Compiler schedules multiple iterations

4. Memory Coalescing

- Adjacent threads load adjacent memory addresses
- 100% memory bandwidth utilization for shared memory loads

5. Bank Conflict Avoidance

- Shared memory width: 33 instead of 32 (padding)
 - 32 banks in shared memory; 33-width prevents conflicts
-

Key Findings

Finding 1: Memory Access Pattern Dominates CPU Performance

- **ABt CPU:** ~300 GFLOPS (contiguous memory access)
- **AtBt CPU:** ~75 GFLOPS (scattered memory access)
- **Penalty:** 4x slowdown from poor cache locality

Finding 2: GPU Eliminates Memory Access Penalty

- **GPU ABt:** ~5,200 GFLOPS
- **GPU AtBt:** ~5,200 GFLOPS
- **Result:** GPU/CPU speedup for AtBt is **70-100x** vs **16-20x** for ABt

Finding 3: Shared Memory is Critical

- Without shared memory: All accesses go to global memory (6.4 TB/s)
- With shared memory tiling: $32 \times$ reduction in global memory traffic
- Effective bandwidth: ~150 GB/s (shared memory) vs 6.4 GB/s (global memory)

Finding 4: Performance Scaling with Nk

Metric	Small Nk (256)	Medium Nk (8K)	Large Nk (131K)	Massive Nk (524K)
CPU ABt	300 GFLOPS	280 GFLOPS	121 GFLOPS	75 GFLOPS
GPU ABt	5,180 GFLOPS	5,130 GFLOPS	4,165 GFLOPS	2,750 GFLOPS
GPU/CPU	17x	18x	34x	36x

- CPU performance degrades ~75% from small to massive Nk
-

GPU performance degrades only ~47% (better scaling)