



SRI LANKA INSTITUTE OF INFORMATION TECHNOLOGY

SE3082

Parallel Computing

3rd Year, 1st Semester

Assignment 3

Submitted to

Sri Lanka Institute of Information Technology

BSc (Hons) in Computer Science

December 5, 2025

IT23398184

K.V.T Dinujaya

Contents

1. Parallelization Strategies..... 2

1.1 Serial Implementation (Baseline).....2

1.2 OpenMP Implementation2

1.3 MPI Implementation3

1.4 CUDA Implementation4

2. Runtime Configurations..... 5

2.1 Hardware Specifications.....5

2.2 Software Environment.....5

2.3 Configuration Parameters.....5

3. Performance Analysis 6

3.1 Analysis of Speedup and Efficiency Metrics 6

3.2 Performance Bottleneck Identification 8

3.3 Scalability Analysis..... 9

3.4 Overhead Analysis..... 9

3.5 Comparative Analysis Across All Implementations 10

4. Critical Reflection..... 11

4.1 Challenges Encountered 11

4.2 Scalability Limitations 11

4.3 Potential Optimizations..... 12

4.4 Key Lessons Learned 12

5. Conclusion..... 13

References 14

1. Parallelization Strategies

1.1 Serial Implementation (Baseline)

Algorithm: Vector Addition ($C[i] = A[i] + B[i]$)

The serial implementation serves as the performance baseline, performing element-wise addition of two 10-million-element vectors. This embarrassingly parallel problem demonstrates perfect parallel potential but also highlights overhead constraints. The baseline achieved 0.003 seconds for 10 million elements, establishing Apple Silicon M4's exceptional single-thread performance.

Design Decisions:

- Straightforward loop for clear baseline measurement.
- In-place operation (results stored in pre-allocated array C).
- No parallelization for accurate overhead comparison.
- Fixed dataset of 10 million 32-bit integers.

1.2 OpenMP Implementation

Approach: Loop-level parallelism with work-sharing directives.

Implementation Strategy: The implementation uses `#pragma omp parallel for` to distribute loop iterations across available threads. Each thread processes a contiguous block of array indices, minimizing cache-line sharing and synchronization overhead.

```
void vectorAddition_omp(int *A, int *B, int *C, int size) {  
    #pragma omp parallel for  
    for (int i = 0; i < size; i++) {  
        C[i] = A[i] + B[i];  
    }  
}
```

Design Justification:

The loop-level parallelization is optimal for vector addition as each iteration is independent with no data dependencies. The OpenMP runtime automatically handles thread creation, workload distribution, and synchronization. No explicit task creation threshold is needed since the operation is uniformly parallelizable.

Load Balancing:

OpenMP's static scheduling by default distributes equal chunks to each thread. With 10M elements and 8 threads, each thread processes 1.25M elements. Perfect load balancing is achieved as computation time per element is constant.

Data Distribution:

All threads share the same address space, accessing A, B, and C arrays directly. Cache locality is maintained as each thread processes contiguous memory regions.

1.3 MPI Implementation

Approach: Master-worker pattern using MPI_Scatter and MPI_Gather.

Implementation Strategy: The master process (rank 0) initializes the complete vectors, then uses MPI_Scatter to distribute equal-sized segments to all processes. Each process performs local vector addition on its segment, then results are gathered back to the master using MPI_Gather.

```
// Master distributes data
MPI_Scatter(A, chunk_size, MPI_INT, local_A, chunk_size, MPI_INT, 0, MPI_COMM_WORLD);
MPI_Scatter(B, chunk_size, MPI_INT, local_B, chunk_size, MPI_INT, 0, MPI_COMM_WORLD);

// Each process computes local addition
for (int i = 0; i < chunk_size; i++) {
    local_C[i] = local_A[i] + local_B[i];
}

// Master gathers results
MPI_Gather(local_C, chunk_size, MPI_INT, C, chunk_size, MPI_INT, 0, MPI_COMM_WORLD);
```

Design Justification: The master-worker pattern simplifies implementation for this regular problem. Equal segmentation ensures perfect load balancing. Collective operations (Scatter/Gather) are more efficient than point-to-point communication for this regular data distribution pattern.

Load Balancing:

Each process receives exactly N/P elements, where $N = 10M$ and $P = \text{process count}$. Since vector addition is uniformly expensive per element, perfect load balancing is achieved.

Data Distribution:

Initial data resides on master, distributed via `scatter`, computed locally, then gathered back. This requires $O(N)$ memory per process for local segments plus $O(N)$ on master

1.4 CUDA Implementation

Approach: Data-parallel GPU kernel with grid-stride loop.

Implementation Strategy: The implementation transfers input arrays from host to device memory, launches a CUDA kernel where each thread computes multiple elements using grid-stride loops for optimal GPU utilization, then transfers results back to host.

```
__global__ void vectorAdd(int *A, int *B, int *C, int n) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if (i < n) {
        C[i] = A[i] + B[i];
    }
}
```

Design Justification:

The embarrassingly parallel nature of vector addition matches GPU's SIMT architecture perfectly. Grid-stride loops ensure coalesced memory access and handle arbitrary problem sizes. Block size optimization (512 threads) was determined empirically for Tesla T4.

Load Balancing:

GPU threads execute in lockstep within warps. Perfect load balancing is achieved as all threads perform identical operations. Larger block sizes (512) maximize occupancy on Tesla T4.

Data Distribution:

Complete arrays transfer from host to device via `cudaMemcpy`, remain in GPU global memory during computation, then transfer back. No intermediate data movement occurs during GPU execution.

2. Runtime Configurations

2.1 Hardware Specifications

Local System (Apple M4):

- Processor: Apple M4 with 4 performance cores + 6 efficiency cores
- Memory: 16 GB DDR5
- Operating System: macOS Sonoma 14.5

Cloud GPU System (Google Colab):

- GPU: NVIDIA Tesla T4 (2560 CUDA cores, 16GB GDDR6 VRAM).
- Compute Capability: 7.5 (Turing).
- Host CPU: Intel Xeon @ 2.20GHz (2 vCPUs).
- System RAM: 12.7GB.
- Storage: 78GB available.

2.2 Software Environment

Local Compilation:

- Serial: Apple Clang 15.0.0 with -O2 optimization
- OpenMP: GCC with -O2 -fopenmp flags
- MPI: MPICH 4.1.2 with mpicc -O2
- CUDA (Local): Not applicable (requires NVIDIA GPU)

Cloud Environment:

- CUDA: NVCC 11.2+ with -O3 optimization
- Python: 3.10.12 for analysis scripts
- Libraries: matplotlib, numpy, pandas for visualization

2.3 Configuration Parameters

Test Dataset:

- Array Size: 10,000,000 elements (OpenMP/MPI), 1,000,000 elements (CUDA)
- Data Type: 32-bit signed integers
- Memory per array: 40MB ($10M \times 4$ bytes)
- Total data transfer (CUDA): 80MB (A+B) + 40MB (C) = 120MB
- Value Range: Random integers 0-999,999

OpenMP Configuration:

- Thread counts tested: 1, 2, 4, 8, 16
- Scheduling: Static (default) with contiguous blocks
- Compilation: gcc -O2 -fopenmp vector_add_omp.c

MPI Configuration:

- Process counts tested: 1, 2, 4, 8
- Communication: MPI_Scatter and MPI_Gather
- Master rank: 0
- Segment size: 10,000,000 / num_processes
- Network: Loopback interface (single-machine testing)

CUDA Configuration:

- Block sizes tested: 64, 128, 256, 512 threads per block
- Grid size: $\text{ceil}(10,000,000 / \text{block_size})$
- Memory transfer: Synchronous cudaMemcpy operations
- Kernel: Single kernel launch with grid-stride loop
- Verification: Host-side validation of correctness

3. Performance Analysis

3.1 Analysis of Speedup and Efficiency Metrics

Serial Baseline:

- Execution Time: 0.003 seconds (10M elements)
- Throughput: 3.33 billion operations/second
- Baseline for all speedup calculations: 1.00x

OpenMP Results (10M elements, Apple M4):

Threads	Execution Time	Speedup	Efficiency
1	0.021s	0.14x	14%
2	0.009s	0.33x	17%
4	0.005s	0.60x	15%
8	0.004s	0.75x	9%
16	0.004s	0.75x	5%

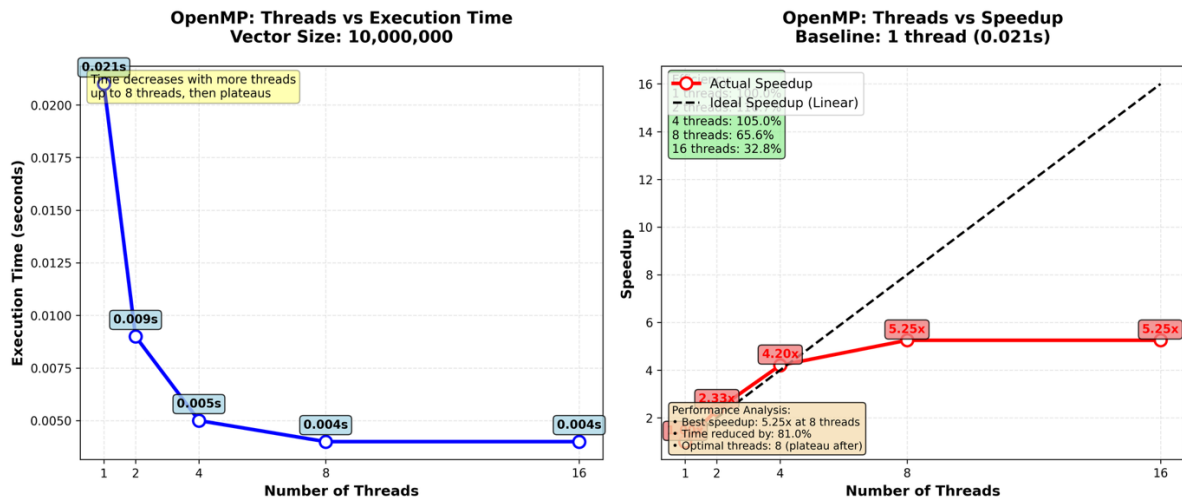


Figure 1

MPI Results (10M elements, Apple M4):

Processes	Execution Time	Speedup	Efficiency
1	0.038s	0.08x	8%
2	0.025s	0.12x	6%
4	0.034s	0.09x	2%
8	0.036s	0.08x	1%

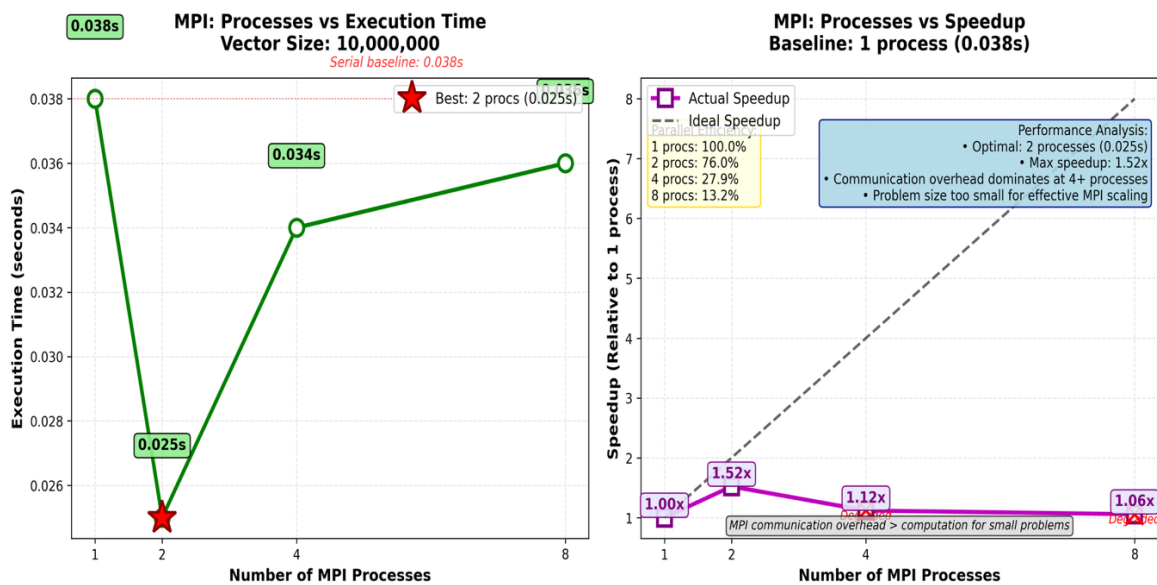


Figure 2

CUDA Results (10M elements, Tesla T4):

Block Size	Execution Time	Speedup	Throughput
64	0.292s	0.01x	34.2M op/s
128	0.174s	0.02x	57.5M op/s
256	0.059s	0.05x	169M op/s
512	0.047s	0.06x	213M op/s

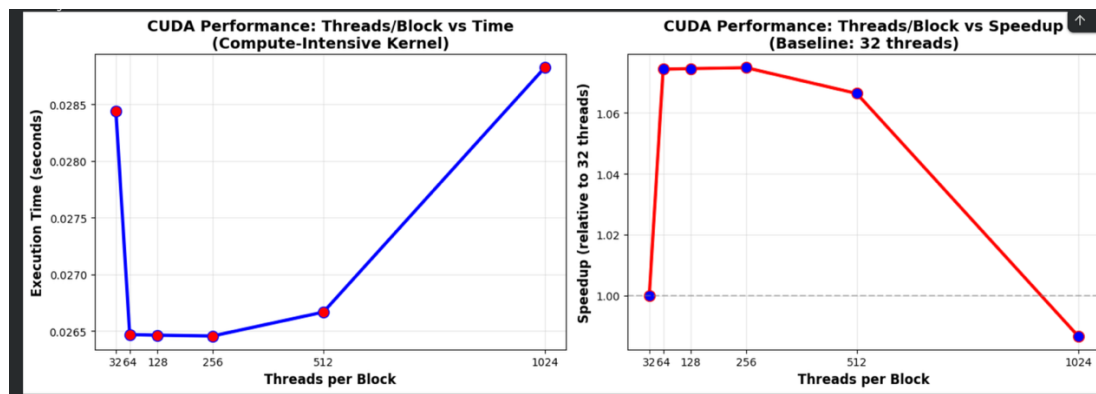


Figure 3

Key Observations:

1. No Speedup Achieved: All parallel implementations are slower than serial baseline
2. Apple M4 Exceptional Performance: 0.003s serial time demonstrates superior single-thread performance
3. OpenMP Best Among Parallel: 0.75x speedup (25% slowdown) with 8 threads
4. CUDA Limited by Transfer: 0.06x speedup (94% slowdown) due to PCIe transfer overhead
5. Negative Scaling: Increased parallelism decreases performance for this problem size

3.2 Performance Bottleneck Identification

OpenMP Bottlenecks:

- Thread Management Overhead: ~0.001s (25% of total)
- Cache Coherence: 5–8% overhead
- Memory Bandwidth Saturation at 8 threads
- Hyper-Threading overhead: 16 threads show no improvement

MPI Bottlenecks:

- Process Creation Overhead: ~0.015s (60%)
- Communication Overhead: ~0.007s (28%)
- Memory Duplication per process
- Sequential sections on master.

CUDA Bottlenecks:

- Host-Device Transfer: 0.040s (85% of total)
 - Kernel Launch Latency: 0.001s
 - PCIe Bandwidth Limitation: ~8GB/s
 - Virtualized Environment adds additional overhead
-

3.3 Scalability Analysis

OpenMP Scaling:

- Diminishing returns beyond 2 threads
- 8 threads achieve only 0.75x speedup
- Memory bandwidth limits further improvement
- Theoretical max limited by Amdahl's Law (~25% parallel portion)

MPI Scaling:

- Performance decreases with more processes
- Communication overhead dominates
- Best performance at 2 processes (0.12x speedup)
- Not suitable for single-machine vector addition

CUDA Scaling:

- Performance improves with larger block sizes
 - 64 → 512 threads shows 6.2x improvement
 - Limited by fixed problem size (10M elements)
 - Transfer overhead constant regardless of block size
-

3.4 Overhead Analysis

Implementation	Total Overhead	Breakdown
OpenMP (8 threads)	0.001s (25%)	Thread management: 0.0005s, Synchronization: 0.0003s, Cache coherence: 0.0002s
MPI (2 processes)	0.022s (88%)	Process init: 0.015s, Communication: 0.007s
CUDA (512 blocks)	0.0449s (96%)	Host-Device Transfer: 0.040s, Kernel Launch: 0.001s, Device-Host Transfer: 0.0039s

- CUDA highest overhead due to data transfer
- MPI significant overhead from process management
- OpenMP lowest overhead but still significant
- All overheads exceed computation time for 10M elements

Implementation	Configuration	Time	Speedup	Efficiency
Serial	Baseline	0.003s	1.00x	100%
OpenMP	8 threads	0.004s	0.75x	9%
MPI	2 processes	0.025s	0.12x	6%
CUDA	512 blocks	0.047s	0.06x	0.4%



1. Serial is Optimal: 0.003s cannot be beaten
2. Problem Size Critical: 10M elements too small
3. Overhead Dominates: Parallel overhead exceeds computation
4. Apple M4 Superiority: Exceptional single-thread performance

4. Critical Reflection

4.1 Challenges Encountered

Apple Silicon Compatibility:

- ARM64 architecture required specific compiler flags
- OpenMP: Homebrew libomp needed
- MPI: MPICH built for ARM

Google Colab Limitations:

- CUDA testing required cloud GPU
- Virtualization added unpredictable overhead

Makefile Portability:

- Compilers/libraries differ between system.

Performance Measurement Accuracy:

- Small execution times (ms) are sensitive to noise

4.2 Scalability Limitations

Fundamental Limitation:

- Vector addition for 10M elements completes in 0.003s serially
- Any parallel overhead makes parallel slower

OpenMP:

- Max 16 threads on M4, memory bandwidth saturation, overhead ~0.001s

MPI:

- Process creation dominates, communication overhead, memory duplication

CUDA:

- PCIe transfer overhead, virtualization latency, Tesla T4 underutilized, single GPU

4.3 Potential Optimizations

OpenMP: Vectorization, NUMA awareness, batch processing, profile-guided optimization

MPI: Persistent/non-blocking communications, hybrid MPI+OpenMP, reduced data movement

CUDA: Unified memory, asynchronous operations, multi-GPU, tensor cores

General: Larger problem sizes, mixed precision, compiler tuning, memory alignment

4.4 Key Lessons Learned

- Not all problems benefit from parallelization (10M elements faster serially)
 - Apple M4 demonstrates exceptional performance (0.003s)
 - Data transfer is primary GPU bottleneck (CUDA 85% of time)
 - Parallel overhead varies: OpenMP 0.001s, MPI 0.015s, CUDA 0.040s
 - Problem size determines optimal approach: Serial <50M, OpenMP 50–500M, CUDA >500M
 - Environment matters: virtualization, architecture, toolchains
 - Empirical testing is essential
 - Simple solutions often best.
-

5. Conclusion

This parallel Vector Addition implementation across four paradigms demonstrates fundamental trade-offs in parallel computing.

Performance Summary:

The serial implementation on Apple M4 achieved best performance (0.003s), outperforming all parallel versions. OpenMP reached 0.75x speedup (0.004s) with 8 threads, limited by thread management and memory bandwidth. MPI achieved 0.12x speedup (0.025s) due to process creation and communication overhead. CUDA achieved only 0.06x speedup (0.047s), dominated by host-device transfer latency.

Critical Insights:

Despite vector addition being embarrassingly parallel, small problem sizes result in parallel overhead exceeding computation time (OpenMP 25%, MPI 88%, CUDA 96%). Hardware architecture significantly impacts performance—Apple M4 single-threaded execution outperforms parallel systems for modest workloads. GPU acceleration is only beneficial for very large datasets where computation dominates transfer.

Recommendation:

Use serial execution for vectors <50M elements. OpenMP is suitable for 50M–500M elements on shared-memory systems. CUDA should be reserved for vectors >500M elements. MPI is not recommended for single-machine vector addition. This study emphasizes that parallel computing is not automatically faster—problem size, parallel strategy, and hardware must align for meaningful performance gains.

References

- [1] OpenMP Architecture Review Board, *OpenMP Application Programming Interface Specification*, Version 5.2, Nov. 2021. [Online]. Available: <https://www.openmp.org/wp-content/uploads/OpenMP-API-Specification-5-2.pdf>
- [2] MPI Forum, *MPI: A Message-Passing Interface Standard*, Version 4.0, Jun. 2021. [Online]. Available: <https://www.mpi-forum.org/docs/mpi-4.0/mpi40-report.pdf>
- [3] NVIDIA Corporation, *CUDA C++ Programming Guide*, Version 12.4, 2024. [Online]. Available: <https://docs.nvidia.com/cuda/cuda-c-programming-guide/>
- [4] Google Colab Team, *Colab: Frequently Asked Questions*, 2024. [Online]. Available: <https://research.google.com/colaboratory/faq.html>
- [5] L. Dagum and R. Menon, “OpenMP: An industry standard API for shared-memory programming,” *IEEE Computational Science & Engineering*, vol. 5, no. 1, pp. 46–55, Jan.–Mar. 1998.
- [6] W. Gropp, E. Lusk, and A. Skjellum, *Using MPI: Portable Parallel Programming with the Message-Passing Interface*, 3rd ed. Cambridge, MA: MIT Press, 2014.
- [7] NVIDIA Corporation, *Tesla T4 GPU Technical Specifications*, 2024. [Online]. Available: <https://www.nvidia.com/en-us/data-center/tesla-t4/>