

Optimization Technique and Implementation Project Report

University of the Cumberlands

MSCS-532-M20: Algorithms and Data Structures

Professor Brandon Bass

Student: Suresh Ghimire

UC ID: 005038288

Date: August 9, 2025

Introduction

High-Performance Computing (HPC) systems depend on optimized algorithms and data structures to handle large-scale problems effectively. While parallelism is a key factor for performance, considering efficiency of the memory hierarchy is nearly as important (Hennessy & Patterson, 2019). The literature recently has shown through many different empirical studies that reformulations targeting data locality and data-structure level of efficiency can lead to significant performance improvements (Yoo et al., 2013).

In this project, we will evaluate one of these optimization approaches, which replaces pointer-based data structures with contiguous arrays and vectorized operations and preallocated buffers as a Python implementation. The goal is to replicate the performance improvement type and effect found in HPC applications and quantify these benefits with a controlled experiment.

Background

An empirical study of HPC performance bugs found that data structure optimization and re-ordering memory access are two of the ways to improve performance (Yoo et al., 2013). Transaction time (including memory management) can be improved specifically by replacing linked lists and fragmented memory with contiguous data structures (e.g., using an array rather than allocating memory for each entry or varying the memory structure for the HPC task). Better performing software has a lot of uses for core advances especially for cache utilization (ensuring contiguous data is retained nearby) by improving cache locality (temporal and data locality improves cache exploitation), thereby reducing TLB page misses (Jia et al., 2018); data LOC is optimized or vectorized and using SIMD (single-instruction, multiple-data) where data structures are memory layouts that can be SIMD friendly and high-level vectorized operations (Intel,

2022); avoid heap overflow inside a loop reduces dynamic memory allotment for fragmentation and having to repeat the overhead of allocation for subsequent iterations or memories (Lam et al., 2015).

Optimization Technique

The approach follows a system denoted as:

- **Data-locality:** Data-locality optimization provides data in contiguous memory, which increases spatial and temporal locality. If data is in consecutive addresses rendered by Address A1, Address A2, Address A3 etc. we can load and reuse contiguous cache lines in the processor, having fewer cache misses.

The empirical study shows that the reasons among the majority of performance bug fixes in HPC applications are due to poor data locality, typically Linked Lists or widely distributed objects or nested Python lists, which regularly cause cache misses whilst fetching their elements in non-contiguous memory locations.

Used the NumPy variable (numpy.ndarray), which are stored in contiguous memory, as opposed to python list-of-tuples.

- **Vectorization:** Vectorization substitutes scalar operations dependent on loops by exploiting SIMD (Single Instruction, Multiple Data) operations which process multiple data elements in a single CPU instruction. Vectorized operations reduce the overhead of investigating Python and take advantage of optimized C/BLAS back-end implementations. The data of the above empirical study indicated that loop unrolling and SIMD transformations are popular in HPC compute in order to maximize CPU vector registers.

- Application for this project:
 - Removed nested python loops within the distance computing
 - Used NumPy broadcasting to compute pairwise distances in bulk.
 - Pre-allocation - allocating memory buffers before acting.
- **Pre-allocation:** Pre-allocation takes care of all memory allocation in one go before any computation, rather than increasing these behaviors dynamically through processing. During computation malloc() behavior lead to multiple heap allocation with copying and as discussed in the methodological study, removing the on-the-fly allocation during computation could speed to other performance costs.

Implementation

Following code snippet shows the baseline implementation for the high-performance computation. The data will be generated on csv file

```
HighPerformanceComputingBenchmark.py > ...
1  import numpy as np
2  import math, random, timeit
3  import tracemalloc
4  import matplotlib.pyplot as plt
5  import pandas as pd
6
7  # Baseline: Python lists + loops
8  def baseline_pairwise(n=1000):
9      points = [(random.random(), random.random()) for _ in range(n)]
10     distances = []
11     for i in range(len(points)):
12         row = []
13         for j in range(len(points)):
14             dx = points[i][0] - points[j][0]
15             dy = points[i][1] - points[j][1]
16             row.append(math.sqrt(dx*dx + dy*dy))
17         distances.append(row)
18     return distances
19
20 # Optimized: NumPy vectorization + preallocation
21 def optimized_pairwise(n=1000):
22     points_np = np.random.rand(n, 2).astype(np.float32)
23     distances_np = np.empty((n, n), dtype=np.float32)
24     diffs = points_np[:, np.newaxis, :] - points_np[np.newaxis, :, :]
25     distances_np[:] = np.sqrt(np.sum(diffs**2, axis=-1))
26     return distances_np
27
```

```

28 # Helper: measure time and memory
29 def measure(func, *args, **kwargs):
30     tracemalloc.start()
31     start_time = timeit.default_timer()
32     _ = func(*args, **kwargs)
33     elapsed = timeit.default_timer() - start_time
34     current, peak = tracemalloc.get_traced_memory()
35     tracemalloc.stop()
36     return elapsed, peak / (1024*1024) # seconds, MB
37
38 # Run experiments
39 sizes = [200, 400, 600, 800, 1000]
40 baseline_times, optimized_times = [], []
41 baseline_mem, optimized_mem = [], []
42
43 for n in sizes:
44     print(f"Running baseline for n={n}...")
45     t, m = measure(baseline_pairwise, n)
46     baseline_times.append(t)
47     baseline_mem.append(m)
48
49     print(f"Running optimized for n={n}...")
50     t, m = measure(optimized_pairwise, n)
51     optimized_times.append(t)
52     optimized_mem.append(m)
53
54 # Save results table
55 df = pd.DataFrame({
56     "n": sizes,
57     "Baseline Time (s)": baseline_times,
58     "Optimized Time (s)": optimized_times,
59     "Baseline Mem (MB)": baseline_mem,
60     "Optimized Mem (MB)": optimized_mem,
61     "Speedup (x)": [b/o for b, o in zip(baseline_times, optimized_times)]
62 })
63
64 df.to_csv("performance_results.csv", index=False)

```

```

63
64 df.to_csv("performance_results.csv", index=False)
65 print("\nSaved results table as performance_results.csv")
66
67 # Plot results
68 plt.figure(figsize=(8,5))
69 plt.plot(sizes, baseline_times, 'r-o', label="Baseline (lists)")
70 plt.plot(sizes, optimized_times, 'g-o', label="Optimized (NumPy)")
71 plt.xlabel("Number of Points (n)")
72 plt.ylabel("Execution Time (seconds)")
73 plt.title("Baseline vs Optimized Performance")
74 plt.legend()
75 plt.grid(True)
76 plt.savefig("performance_comparison.png")
77 print("Saved performance plot as performance_comparison.png")
78
79 # Display summary
80 print("\n\tBenchmark Summary")
81 print(df)
82

```

Results and analysis

Following screenshot shows the results of the above script

```

PS C:\Users\sures\Desktop\University of the Cumberland\MSCS-532-DSA\MSCS532_Final_Project> python -u "c:\Users\sures\Desktop\University of the Cumberland\MSCS-532-DSA\MS
CS532_Final_Project\HighPerformanceComputingBenchmark.py"
Running baseline for n=200...
Running optimized for n=200...
Running baseline for n=400...
Running optimized for n=400...
Running baseline for n=600...
Running optimized for n=600...
Running baseline for n=800...
Running optimized for n=800...
Running baseline for n=1000...
Running optimized for n=1000...

Saved results table as performance_results.csv
Saved performance plot as performance_comparison.png

Benchmark Summary
n Baseline Time (s) Optimized Time (s) Baseline Mem (MB) Optimized Mem (MB) Speedup (x)
0 200 0.150320 0.004106 1.253510 0.950352 36.612441
1 400 0.641152 0.008210 4.933083 3.697742 78.091166
2 600 1.581484 0.021948 11.390175 8.276905 72.056301
3 800 2.236443 0.026916 19.994423 14.687122 83.088195
4 1000 3.062934 0.031531 31.400978 22.928394 97.139157

```

The benchmark was run on datasets of 200, 400, 600, 800, and 1,000 data points. For each dataset size, the baseline implementation (Python lists with nested loops) and optimized implementation (NumPy arrays with vectorization and pre-allocation) were also run. Execution time and peak memory usage were captured as well as the calculated speedup factor.

Following table shows the output of the above computational script

N(points)	Baseline Time(s)	Optimized Time(s)	Baseline Memory (MB)	Optimized Memory (MB)	Speedup
200	0.150320	0.004106	1.253510	0.950352	36.61
400	0.641152	0.008210	4.933083	3.697742	78.09
600	1.581484	0.021948	11.390175	8.276095	72.06
800	2.236443	0.026916	19.994423	14.687122	83.09

1000	3.062934	0.031531	31.400978	22.928394	97.14
------	----------	----------	-----------	-----------	-------

Observation

Reduction in Time to Complete Task:

- The optimized version exhibits speedup ranging from $36\times$ (for the smallest dataset) to over $97\times$ (for the largest dataset).
- The largest gains are associated with larger datasets because of reductions in Python interpreter overhead as well as better cache efficiency for vectorized operations.

Memory Savings:

- Memory usage is approximately 20–30 % lower for all dataset sizes.
- This is due to contiguous storage in NumPy arrays, which eliminates the overhead of Python objects for each element.

Scalability:

- The performance improvement gap increases with size of input dataset, thereby demonstrating that the proportional benefits of vectorization and data-locality are bigger for larger computations.
- The optimized approach is scalable with respect to computation time but has a much smaller constant factor.

These results are consistent with Yoo et al. (2013), who cite optimization of data structures, loop transformations, and modifications to memory access patterns, as major contributors to the performance gains attained in HPC applications.

Conclusion

This project illustrated that, when deploying data-locality optimization, vectorization, and pre-allocation in tandem, it had significant performance improvements of up to $97\times$ faster execution time and 20–30% less memory than a baseline Python list-based implementation. These results are in line with the research as reported in the Empirical Study of High Performance Computing Performance Bugs (Yoo et al., 2013), which indicated that data structure optimizations and memory access optimizations tend to be two of the best they found for HPC. The scalability observed demonstrates that while small datasets may not benefit as much, more efficient memory layout and bulk computation will not only lead to improved performance at a large scale, but actually yield performance that could be equal to a hardware upgrade (e.g., reduced the memory requirements that normally requires a few hundred GB of RAM), thus representing a really beneficial strategy to employ for research prototypes and production scale scientific computing modules.

References:

Cormen, T. H., Leiserson, C. E., Rivest, R. L., & Stein, C. (2022). *Introduction to Algorithms* (4th ed.). Random House Publishing Services.

<https://reader2.yuzu.com/books/9780262367509>

Hennessy, J. L., & Patterson, D. A. (2019). *Computer architecture: A quantitative approach* (6th ed.). Morgan Kaufmann.

Intel. (2022). *Intel® 64 and IA-32 architectures optimization reference manual*. Intel Corporation.

Jia, Z., Zaharia, M., & Aiken, A. (2018). Beyond data and model parallelism for deep neural networks. *Proceedings of the 2nd Conference on Systems and Machine Learning*.

Lam, S. K., Pitrou, A., & Seibert, S. (2015). Numba: A LLVM-based Python JIT compiler. *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC*, 1–6.

Yoo, R. M., Romano, A., & Amarasinghe, S. (2013). An empirical study of high performance computing performance bugs. *Proceedings of the ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, 15–26.