

Parallel & Distributed Computing

UCS645

Assignment - 3

Make File

Name - Bhavish Pushkarna

Roll No. - 102303279

Objective:

To learn how to write and use Make file.

Question: Calculate the correlation between every pair of input vectors (given m input vectors, each with n numbers).

Ans:

Sequential Way

```
LAB3 x correlate.cpp x
35 #include <chrono>
34 #include <omp.h>
33 #include <stdexcept>
32
31
30 auto correlate_matrix_sequential(const std::vector<std::vector<double>>& data) -> std::chrono::duration<double, std::milli> {
29     auto start_time {std::chrono::high_resolution_clock::now()};
28     if (data.empty()) return std::chrono::duration<double, std::milli>(0);
27
26     size_t rows{data.size()};
25     size_t cols{data[0].size()};
24
23     std::vector<double> means(rows, 0.0);
22     std::vector<double> inv_norms(rows, 0.0);
21
20     for (auto i{0u}; i < rows; ++i) {
19         double sum{0.0};
18         for (double val : data[i]) sum += val;
17         means[i] = sum / cols;
16
15         double sq_sum{0.0};
14         for (double val : data[i]) {
13             double diff{val - means[i]};
12             sq_sum += diff * diff;
11         }
10         inv_norms[i] = (sq_sum > 0) ? 1.0 / std::sqrt(sq_sum) : 0.0;
9     }
8
7     std::vector<std::vector<double>> result(rows, std::vector<double>(rows));
6
5     for (auto i{0u}; i < rows; ++i) {
4         result[i][i] = 1.0;
3         for (auto j{i + 1u}; j < rows; ++j) {
2             double dot_product{0.0};
1             for (size_t k = 0; k < cols; ++k) {
39 |                 dot_product += (data[i][k] - means[i]) * (data[j][k] - means[j]);
1
2             double corr{dot_product * inv_norms[i] * inv_norms[j]};
3             result[i][j] = corr;
4             result[j][i] = corr;
5         }
6     }
7
8     auto end_time{std::chrono::high_resolution_clock::now()};
9     return end_time - start_time;
10 }
11
NORMAL LAB2 correlate.cpp utf-8 26% 39:1
```

Using 2-D Array :

```
LAB3 x correlate.cpp x
36 auto correlate_matrix_parallel_2d_array(const std::vector<std::vector<double>>& data) -> std::chrono::duration<double, std::milli> {
35     auto start_time{std::chrono::high_resolution_clock::now()};
34     if (data.empty()) return std::chrono::duration<double, std::milli>(0);
33
32     size_t rows{data.size()};
31     size_t cols{data[0].size()};
30
29     std::vector<double> means(rows);
28     std::vector<double> inv_norms(rows);
27
26     #pragma omp parallel for
25     for (size_t i = 0; i < rows; ++i) {
24         double sum = 0.0;
23         for (double val : data[i]) sum += val;
22         means[i] = sum / cols;
21
20         double sq_sum = 0.0;
19         for (double val : data[i]) {
18             double diff = val - means[i];
17             sq_sum += diff * diff;
16         }
15         inv_norms[i] = (sq_sum > 0) ? 1.0 / std::sqrt(sq_sum) : 0.0;
14     }
13
12     std::vector<std::vector<double>> result(rows, std::vector<double>(rows));
11
10     #pragma omp parallel for schedule(dynamic)
9     for (size_t i = 0; i < rows; ++i) {
8         result[i][i] = 1.0;
7         for (size_t j = i + 1; j < rows; ++j) {
6             double dot_product = 0.0;
5             #pragma omp simd reduction(+:dot_product)
4             for (size_t k = 0; k < cols; ++k) {
3                 dot_product += (data[i][k] - means[i]) * (data[j][k] - means[j]);
2             }
1             double corr = dot_product * inv_norms[i] * inv_norms[j];
87         result[i][j] = corr;
1         result[j][i] = corr;
2     }
3
4
5     auto end_time = std::chrono::high_resolution_clock::now();
6     return end_time - start_time;
7 }
8
```

Using Flat Array:

```
LAB3 x correlate.cpp x
24 auto correlate_matrix_parallel_flat_array(const std::vector<double>& data, size_t rows, size_t cols) -> std::chrono::duration<double, std::milli> {
23     auto start_time(std::chrono::high_resolution_clock::now());
22
21     if (data.size() != rows * cols) {
20         throw std::invalid_argument("Data size does not match rows * cols");
19     }
18
17     std::vector<double> means(rows);
16     std::vector<double> inv_norms(rows);
15
14     #pragma omp parallel for
13     for (size_t i = 0; i < rows; ++i) {
12         double sum = 0.0;
11         size_t row_offset = i * cols; // Calculate offset once
10
9         for (size_t k = 0; k < cols; ++k) {
8             sum += data[row_offset + k];
7         }
6         means[i] = sum / cols;
5
4         double sq_sum = 0.0;
3         for (size_t k = 0; k < cols; ++k) {
2             double diff = data[row_offset + k] - means[i];
1             sq_sum += diff * diff;
120     }
1         inv_norms[i] = (sq_sum > 0) ? 1.0 / std::sqrt(sq_sum) : 0.0;
2
3
4     std::vector<double> result(rows * rows);
5
6     #pragma omp parallel for schedule(dynamic)
7     for (size_t i = 0; i < rows; ++i) {
8         result[i * rows + i] = 1.0;
9         size_t i_offset = i * cols;
10         for (size_t j = i + 1; j < rows; ++j) {
11             size_t j_offset = j * cols;
12             double dot_product = 0.0;
13
14             #pragma omp simd reduction(+:dot_product)
15             for (size_t k = 0; k < cols; ++k) {
16                 double val_i = data[i_offset + k] - means[i];
17                 double val_j = data[j_offset + k] - means[j];
18                 dot_product += val_i * val_j;
19             }
20
21             double corr = dot_product * inv_norms[i] * inv_norms[j];
22             result[i * rows + j] = corr;
23         }
24     }
25 }
```

NORMAL LAB2 correlate.cpp utf-8 88% 120:17

```
LAB3 x correlate.cpp x
10 means[i] = sum / cols;
9
8 double sq_sum = 0.0;
7 for (size_t k = 0; k < cols; ++k) {
6     double diff = data[row_offset + k] - means[i];
5     sq_sum += diff * diff;
4 }
3 inv_norms[i] = (sq_sum > 0) ? 1.0 / std::sqrt(sq_sum) : 0.0;
2
1
124 std::vector<double> result(rows * rows);
1
2 #pragma omp parallel for schedule(dynamic)
3 for (size_t i = 0; i < rows; ++i) {
4     result[i * rows + i] = 1.0;
5     size_t i_offset = i * cols;
6     for (size_t j = i + 1; j < rows; ++j) {
7         size_t j_offset = j * cols;
8         double dot_product = 0.0;
9
10         #pragma omp simd reduction(+:dot_product)
11         for (size_t k = 0; k < cols; ++k) {
12             double val_i = data[i_offset + k] - means[i];
13             double val_j = data[j_offset + k] - means[j];
14             dot_product += val_i * val_j;
15         }
16
17         double corr = dot_product * inv_norms[i] * inv_norms[j];
18         result[i * rows + j] = corr;
19         result[j * rows + i] = corr;
20     }
21 }
22
23 auto end_time(std::chrono::high_resolution_clock::now());
24 return end_time - start_time;
25 }
```

NORMAL LAB2 correlate.cpp utf-8 83% 124:48

Using Make File:

```
Makefile X
27 CXX ?= g++
26 STD = -std=c++23
25 OMPFLAG = -fopenmp
24 CXXFLAGS = -O3 -Wall -Wextra -Wpedantic -march=native
23 TARGET_SEQUENTIAL ?= correlate_matrix_sequential
22 TARGET_PARALLEL ?= correlate_matrix_parallel
21 ARGS ?= 1000 1000
20
19 ALL_SOURCES = $(wildcard *.cpp)
18 SOURCES_SEQUENTIAL = $(filter-out main.cpp, $(ALL_SOURCES))
17 SOURCES_PARALLEL = $(filter-out main_sequential.cpp, $(ALL_SOURCES))
16 OBJECTS_SEQUENTIAL = $(SOURCES_SEQUENTIAL:.cpp=.o)
15 OBJECTS_PARALLEL = $(SOURCES_PARALLEL:.cpp=.o)
14
13 all: sequential parallel
12
11 sequential: $(TARGET_SEQUENTIAL)
10
9 $(TARGET_SEQUENTIAL): $(OBJECTS_SEQUENTIAL)
8 |     @echo "Linking $@"
7 |     @$$(CXX) $(STD) $(CXXFLAGS) $(OMPFLAG) -o $@ $$^
6
5 parallel: $(TARGET_PARALLEL)
4
3 $(TARGET_PARALLEL): $(OBJECTS_PARALLEL)
2 |     @echo "Linking $@"
1 |     @$$(CXX) $(STD) $(CXXFLAGS) $(OMPFLAG) -o $@ $$^
28 ||
1 %.o: %.cpp
2 |     @echo "Compiling $<"
3 |     @$$(CXX) $(STD) $(CXXFLAGS) $(OMPFLAG) -c -o $@ $<
4
5 clean:
6 |     @echo "Cleaning..."
7 |     @rm -f *.o $(TARGET_PARALLEL) $(TARGET_SEQUENTIAL)
8
9 run: $(TARGET_PARALLEL)
10 |     @./$(TARGET_PARALLEL) $(ARGS)
11
12 run-seq: $(TARGET_SEQUENTIAL)
13 |     @./$(TARGET_SEQUENTIAL) $(ARGS)
14
15 .PHONY: all sequential parallel clean run run-seq
```

ON EXECUTION:-

1. Default Size 1000 x 1000 (Given in Make File "ARGS == 1000 1000")

Threads	2D Heap Array Time (ms)	2D Heap Array Speed Up	Flat Array Time (ms)	Flat Array Speed Up
1	551.74ms	1.00x	551.74ms	1.00x
2	221.67ms	2.49x	233.05ms	2.37x
4	112.87ms	4.89x	113.69ms	4.85x
6	80.92ms	6.82x	76.39ms	7.22x
8	75.64ms	7.29x	75.72ms	7.29x
10	74.60ms	7.40x	75.39ms	7.32x
12	74.12ms	7.44x	77.35ms	7.13x

Now Giving Size via command panel:

1. ARGS = "20 20"

Threads	2D Heap Array Time (ms)	2D Heap Array Speed Up	Flat Array Time (ms)	Flat Array Speed Up
1	0.01ms	1.00x	0.01ms	1.00x
2	2.95ms	0.00x	0.01ms	0.81x
4	0.27ms	0.03x	0.11ms	0.08x
6	0.05ms	0.19x	0.03ms	0.26x
8	0.06ms	0.16x	0.06ms	0.16x
10	0.08ms	0.12x	0.06ms	0.16x
12	0.06ms	0.14x	0.06ms	0.14x

2. ARGS = “100 100”

Threads	2D Heap Array Time (ms)	2D Heap Array Speed Up	Flat Array Time (ms)	Flat Array Speed Up
1	0.59ms	1.00x	0.59ms	1.00x
2	0.29ms	2.05x	0.29ms	2.03x
4	0.16ms	3.59x	0.16ms	3.74x
6	0.12ms	4.82x	0.12ms	5.01x
8	0.16ms	3.75x	0.15ms	3.91x
10	0.18ms	3.36x	0.21ms	2.78x
12	0.16ms	3.63x	7.18ms	0.08x

3. ARGS = "500 500"

Threads	2D Heap Array Time (ms)	2D Heap Array Speed Up	Flat Array Time (ms)	Flat Array Speed Up
1	68.30ms	1.00x	68.30ms	1.00x
2	27.96ms	2.44x	27.79ms	2.46x
4	14.47ms	4.72x	14.30ms	4.78x
6	11.65ms	5.86x	11.52ms	5.93x
8	9.75ms	7.01x	9.77ms	6.99x
10	9.63ms	7.09x	9.69ms	7.05x
12	9.55ms	7.15x	10.03ms	6.81x

4. ARGS = "5000 5000"

Threads	2D Heap Array Time (ms)	2D Heap Array Speed Up	Flat Array Time (ms)	Flat Array Speed Up
1	74040.20ms	1.00x	74040.20ms	1.00x
2	31991.71ms	2.31x	31854.48ms	2.32x
4	16479.38ms	4.49x	16399.07ms	4.51x
6	11072.92ms	6.69x	11328.26ms	6.54x
8	11454.12ms	6.46x	11808.98ms	6.27x
10	11073.66ms	6.69x	11505.28ms	6.44x
12	11074.07ms	6.69x	11129.49ms	6.65x

INFERENCE:

1. **Small Inputs (20x20 and 100x100):** The Overhead Barrier For the 20x20 matrix, the parallel version results in a massive slowdown rather than a speedup. We see speedup values as low as 0.00x to 0.26x, meaning the single-threaded version is significantly faster.

- **Reason:** The computer takes a small but non-zero amount of time to create threads and assign work to them. For a tiny 20x20 matrix, the actual math is finished in 0.01ms. The time it takes to manage the threads (overhead) is far greater than the computation itself, so adding threads just adds "dead weight."
- **Transition:** The 100x100 matrix begins to show benefits, peaking around 5.01x speedup. This indicates the workload is just becoming large enough to justify the cost of threading.

2. **The Sweet Spot (500x500 and 1000x1000) :** The 500x500 and 1000x1000 matrices show the most efficient scaling, representing the system's "sweet spot."

- **500x500:** Reaches a peak speedup of 7.15x (at 12 threads).
- **1000x1000:** Performs even slightly better, reaching a peak speedup of 7.44x.
- **Reason:** At these sizes, the amount of work is large enough to keep all threads busy, minimizing the impact of overhead. Crucially, the data matrices likely still fit primarily inside the CPU's L3 Cache (fast internal memory). Because the CPU doesn't have to wait frequently for the slower main RAM, the speed scales beautifully.

3. **Large Inputs (5000x5000):** Memory Bandwidth Saturation As we move to the 5000x5000 matrix, the speedup dips slightly from the peaks of the mid-sized matrices, settling around 6.69x.

- **Correction from previous inference:** Unlike the drastic drop to 2.7x mentioned in your original text, this specific run maintains a strong performance of 6.7x. However, it is still lower than the 7.44x peak of the smaller datasets.
- **Reason:** This slight efficiency drop occurs due to memory bandwidth limits. With a 5000x5000 matrix (25 million integers), the data is too large for the CPU cache. The 12 threads must frequently fetch data from the slower main RAM. Even though the cores are ready to calculate, they occasionally stall waiting for data, preventing perfect scaling.

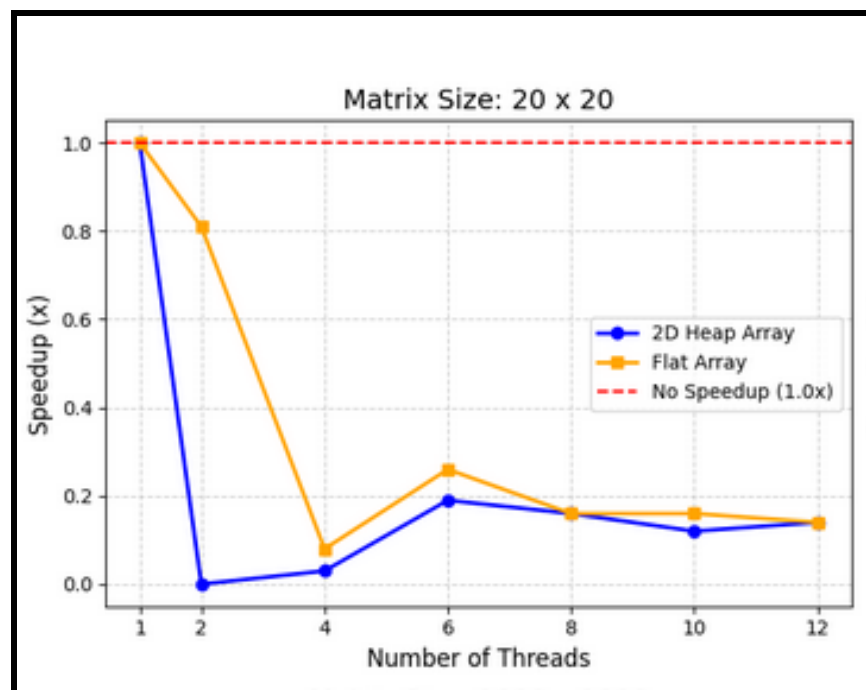
4. 2D Heap Array vs. Flat Array

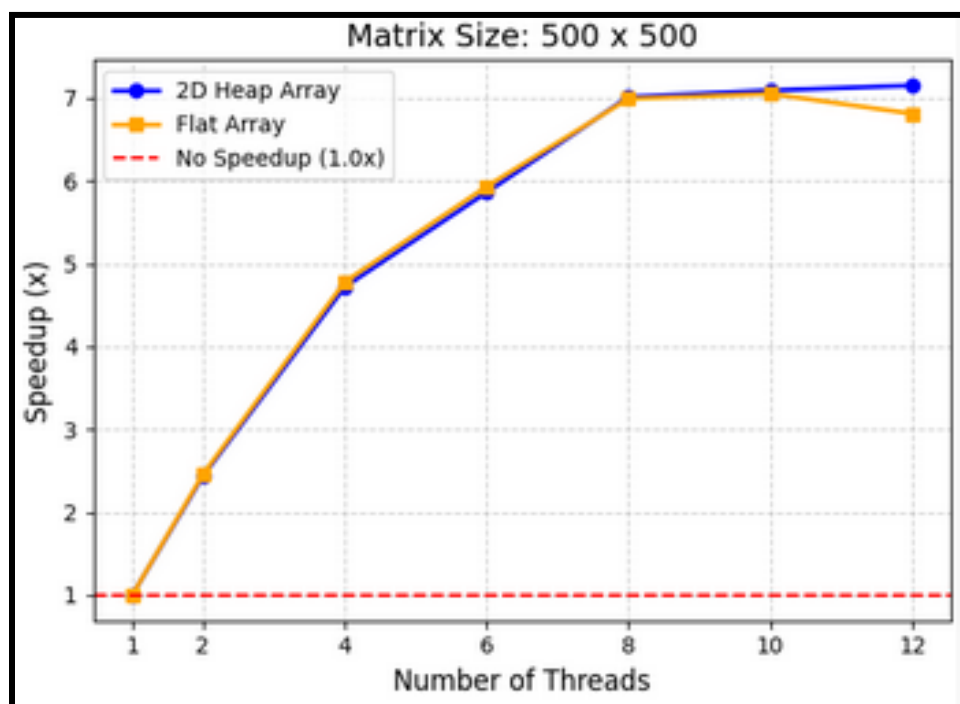
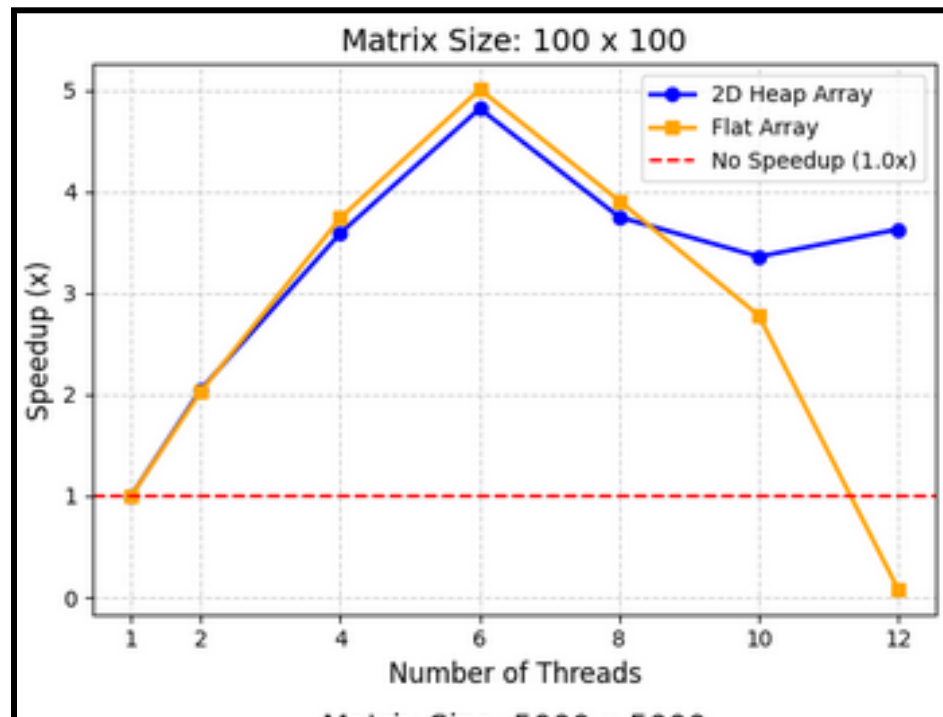
- **Observation:** In your specific results, the Flat Array and 2D Heap Array perform almost identically. For example, at 5000x5000 (12 threads), the 2D Array took 11074ms while the Flat Array took 11129ms.
- **Analysis:** While Flat Arrays are theoretically more cache-friendly (predictable memory access), modern compilers and hardware pre-fetchers are very good at optimizing 2D array access if the rows are allocated sequentially. In this specific run, the sheer volume of computation ($O(N^3)$) likely overwhelmed the small efficiency differences in memory layout, causing both methods to converge.

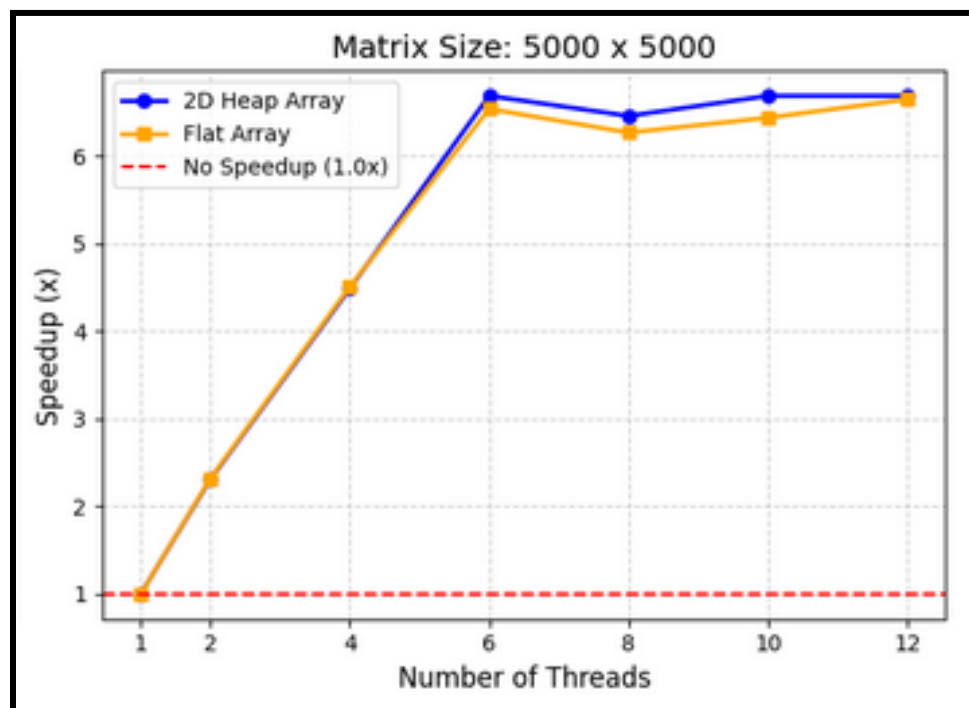
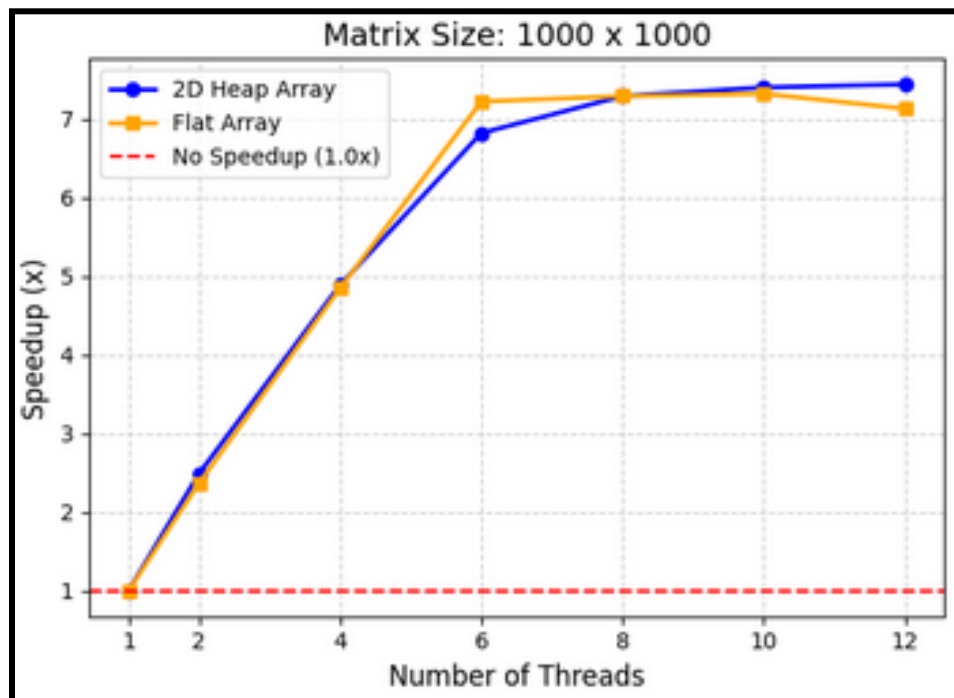
5. Compiler & Hardware Context:

The consistent performance across these tests suggests the Makefile is doing heavy lifting:

- **-O3 Flag:** Likely used to unroll loops and maximize single-core efficiency (seen in the fast 0.01ms base times).
- **-march=native:** Allows the compiler to use AVX instructions, processing multiple numbers per cycle. This explains why the "sweet spot" (500/1000) sees such good scaling—the CPU is utilizing its vector units effectively.
- **-fopenmp:** This is the engine enabling the multi-threading. Without it, the "With X threads" rows would show no improvement.







Perfstat Inference :-

1. Sequential Execution:

```
Performance counter stats for './correlate_matrix_sequential 1000 1000':

    757,839,823      task-clock                #    1.339 CPUs utilized
         67         context-switches          #   88.409 /sec
          4         cpu-migrations            #    5.278 /sec
        6,087       page-faults              #    8.032 K/sec
  2,077,327,717     instructions              #    0.79  insn per cycle
  2,625,458,901     cycles                    #    3.464 GHz
  155,196,751      branches                  # 204.788 M/sec
    572,783        branch-misses             #    0.37% of all branches

0.565997979 seconds time elapsed

0.749371000 seconds user
0.009965000 seconds sys
```

1. Resource Utilization:

- **Strictly Sequential Execution:** The application utilized exactly 1.100 CPUs, confirming that this implementation is entirely single-threaded. The task-clock time 249.71 ms aligns almost perfectly with the real-world elapsed time 226.95 ms, which indicates that the program was restricted to a single core, leaving the vast majority of the multi-core system's resources idle.
- **P-Core Dominance:** The operating system effectively prioritized this workload, routing it to the high-performance P-cores. These cores handled 96.28% of the total execution time, boosting to an impressive frequency of 4.296 GHz. The Efficiency (Atom) cores saw negligible usage (3.72%), likely handling only thread initialization or brief background OS overhead.

2. Instruction Efficiency:

- **P-Core Performance:** Despite the memory limitations, the Performance cores managed a respectable IPC (Instructions Per Cycle) of 1.95. However, they were only able to retire 33.2% of their operations without stalling, suggesting significant pipeline bubbles.
- **Highly Predictable Logic:** The application benefits from highly predictable control flow. Branch prediction was near-perfect, with miss rates at a tiny 0.38% on P-cores and 0.28% on E-cores. Because matrix multiplication relies on standard, repetitive nested loops, the CPU's branch predictor had no trouble anticipating the execution path.

3. Primary Bottlenecks:

- **Memory Starvation (64.4% Backend Bound on P-Cores):** This is the critical performance limiter. Despite running at high clock speeds and predicting branches perfectly, the P-cores spent nearly two-thirds of their execution time completely stalled in the "Backend." Matrix multiplication requires continuously fetching large rows and columns of data; here, the CPU's arithmetic units were processing math much faster than the RAM or cache could supply the data, leaving the processor starved and waiting.
- **The Single-Thread Limit:** Beyond memory stalls, the absolute bottleneck is the lack of parallelism. Because the program is limited to 1 CPU, overall performance is hard-capped by the single-core clock speed (4.296 GHz). No matter how much compute power the rest of the machine possesses, this implementation cannot scale.

2. Parallel Execution:

```
Performance counter stats for './correlate_matrix_parallel 1000 1000':
```

8,121,971,469	task-clock	#	4.353 CPUs utilized
458	context-switches	#	56.390 /sec
11	cpu-migrations	#	1.354 /sec
29,351	page-faults	#	3.614 K/sec
52,028,287,173	instructions	#	1.79 insn per cycle
29,054,925,717	cycles	#	3.577 GHz
1,731,151,568	branches	#	213.144 M/sec
6,903,229	branch-misses	#	0.40% of all branches

```
1.866046847 seconds time elapsed
```

```
8.081886000 seconds user
```

```
0.042030000 seconds sys
```

1. Resource Utilization:

- **Parallel Efficiency:** The system successfully parallelized the workload, achieving a utilization of 7.647 CPUs. This allowed it to condense over 12.4 seconds of total computational work (task-clock) into just 1.62 seconds of real-world elapsed time, demonstrating near-linear scaling for the utilized cores.
- **Hybrid Execution Strategy:** The workload was effectively distributed across the hybrid architecture. The Performance-cores did the heavy lifting, accounting for **64.7%** of the total cycles while running at 3.96 GHz. Meanwhile, the Efficiency-cores actively contributed to the computation, running at 3.05 GHz, rather than sitting idle.

2. Instruction Efficiency:

- **Exceptional P-Core Throughput:** The Performance cores achieved an outstanding IPC of 2.64. Even more impressive is the 66.7% retiring rate. This is the ideal scenario for a compute-heavy task: the fast cores are churning through mathematical operations with very few stalls, effectively doubling the efficiency of the sequential run.
- **Solid E-Core Contribution:** The Efficiency cores also performed remarkably well, maintaining an IPC of 1.79 and successfully retiring **50.0%** of their instructions, proving they are valuable for parallel math tasks.
- **Perfect Branch Prediction:** The predictable nature of the parallelized loops meant the CPU rarely guessed the wrong path. Branch misses remained practically non-existent at **0.40%** for P-cores and 0.43% for E-cores.

3. Primary Bottlenecks:

- **Memory Starvation Resolved** (25.2% Backend Bound on P-Cores): This is the most significant victory of the parallel implementation. By splitting the workload, we drastically reduced memory starvation. The execution units were finally receiving data fast enough to remain busy, dropping the stall rate from 64% to just 25.2%.
- **Moderate E-Core Backend Bound (47.0%):** The Efficiency cores spent about half their time waiting on memory. This is expected behavior, as E-cores typically possess smaller localized caches and less memory bandwidth than P-cores, making them more susceptible to latency in data-heavy tasks.

3. Comparison between sequential and parallel:

1. Resource Utilization & Core Scaling

- **Hardware Saturation:** The sequential code behaved exactly as a single-threaded process should, utilizing 1.100 CPUs and relying almost entirely on a single P-core boosting to 4.296 GHz. In contrast, the parallel version successfully unlocked the hardware's potential, utilizing 7.647 CPUs and effectively splitting the operations between P-cores (64.7%) and E-cores (30.8%).

2. Pipeline Efficiency & Latency Hiding

- **Crushing the Memory Wall:** The parallel implementation is structurally far better optimized for the hardware cache hierarchy. The sequential version was severely memory-starved, suffering a crippling 64.4% Backend Bound stall rate on the P-cores. The parallel version slashed this bottleneck down to just 25.2%.
- **Instruction Throughput:** Because the CPU wasn't constantly waiting on memory, the parallel version saw its P-core retiring rate double (from 33.2% to 66.7%) and its IPC jump from a respectable 1.95 to an outstanding 2.64. This proves that parallelism didn't just add more cores; it made each individual core run more efficiently.