

Multithreaded Monte Carlo Estimation of π with Pthreads

Mehdi Khameedeh 40131873

Full HW file and results are available at this [github link](#).

2025

Abstract

This report presents a comprehensive analysis of multithreaded implementations of the Monte Carlo method to estimate π using POSIX Threads (pthreads). The work implements four distinct approaches: (1) sequential baseline, (2) basic multithreaded with global counter synchronization, (3) optimized with thread-local counters, and (4) advanced byte-array partitioning strategy. The analysis demonstrates how synchronization strategy directly impacts performance scaling. Experimental validation across three sample sizes (100K, 1M, 10M) and five thread counts (1, 2, 4, 8, 16) reveals that thread-local counters achieve superior speedup ($6.90 \times$ at 16 threads for 10M samples) compared to basic mutex-protected counters ($7.60 \times$), while byte-array partitioning provides consistent performance but higher memory overhead. The study quantifies the cost of mutex contention in synchronization-heavy implementations and demonstrates that careful algorithm redesign to reduce synchronization overhead provides equivalent or superior scalability. Results confirm that on an 8-core system, efficiency peaks around 8 threads before diminishing returns, with speedup following theoretical predictions for work-aggregate patterns.

1 Setup and Methodology

The experimental framework executes multiple implementations of the Monte Carlo estimation algorithm across varying thread counts and sample sizes. All experiments were conducted on a single machine with consistent runtime environment, ensuring reproducible results and enabling direct performance comparison between implementations.

1.1 System Specification

All testing was performed on a single machine with sufficient computational resources. The experimental validation framework executed each implementation multiple times with deterministic output capture. The test suite was designed to systematically explore the relationship between thread count, sample size, and performance metrics including execution time, estimation accuracy, and speedup relative to the sequential baseline.

1.2 Algorithm Overview and Mathematical Foundation

The Monte Carlo method estimates π by exploiting the geometric relationship between a circle and a circumscribed square:

- **Setup:** Consider a unit square $[0, 1] \times [0, 1]$ with an inscribed quarter-circle of radius 1 centered at the origin.
- **Area Relationship:** Area of square = 1. Area of quarter-circle = $\pi/4$.

- **Probability:** The probability that a random point in the square falls within the circle is $P = \frac{\pi/4}{1} = \pi/4$.
- **Estimation:** If we generate n random points and count k points inside the circle, then $k/n \approx \pi/4$, yielding $\pi \approx 4k/n$.

The algorithm generates random points (x, y) with $x, y \in [0, 1]$ and checks the condition $x^2 + y^2 \leq 1$ to determine inclusion. Increasing sample size n improves accuracy; the error scales approximately as $O(1/\sqrt{n})$. Parallelization allows multiple threads to generate samples concurrently, with thread scalability determined by how efficiently results are aggregated.

1.3 Testing Configuration and Parameters

The test suite comprises 48 experiments spanning three sample sizes and five thread counts:

Table 1: Testing Configuration Parameters

Sample Size	Thread Counts	Total Configurations
100,000 samples	1, 2, 4, 8, 16	5
1,000,000 samples	1, 2, 4, 8, 16	5
10,000,000 samples	1, 2, 4, 8, 16	5
Per Implementation:	15 tests	
Total (4 implementations):	48 tests	

1.4 Validation Methodology

Each experiment measures:

- **Execution Time:** Elapsed time from program start to completion (seconds), measured using `gettimeofday()` for microsecond-level precision.
- **Accuracy:** The estimated value of π and absolute error relative to the true value ($\pi = 3.14159265\ldots$).
- **Speedup:** Execution time of sequential implementation divided by parallel implementation time, normalized to sequential baseline = 1.0.
- **Scalability:** How speedup increases with thread count, indicating whether implementations exploit parallelism effectively.

Results are collected in CSV format and analyzed to identify performance trends across implementations and thread counts.

2 Introduction

Multithreading enables exploiting multicore processors to accelerate computationally intensive algorithms. The Monte Carlo estimation presents an ideal case study: the algorithm is naturally parallelizable (independent sample generation), yet requires careful synchronization when aggregating results. The challenge lies in minimizing synchronization overhead while maintaining correctness.

2.1 Synchronization Patterns in Parallel Algorithms

Parallel algorithms must protect access to shared data using synchronization mechanisms (mutexes, atomic operations). In the Monte Carlo case, the shared resource is the count of points inside the circle. Different synchronization strategies create different trade-offs:

- **Global Counter with Mutex:** Simplest approach but creates contention—every thread must acquire the mutex to update the counter, serializing the critical path.
- **Thread-Local Counters:** Each thread maintains its own counter (no synchronization), then results are aggregated after all threads complete. Eliminates contention during execution.
- **Array-Based Partitioning:** Store individual sample results (0 or 1) in array partitions, with threads writing to non-overlapping regions. Eliminates false sharing and synchronization.

The relative performance of these approaches depends on sample size, thread count, and contention pressure.

3 Implementation Variants

3.1 Implementation 1: Sequential Baseline

The sequential implementation establishes the performance baseline. It generates all samples serially and counts points inside the circle without any concurrency:

```
for (long long i = 0; i < samples; i++) {
    double x = drand48();
    double y = drand48();
    if (x * x + y * y <= 1.0)
        count_inside++;
}
```

Characteristics:

- No synchronization overhead
- No thread creation/destruction costs
- Reference point for speedup calculations

3.2 Implementation 2: Basic Multithreaded with Global Counter

The basic multithreaded version divides samples among threads, with each thread updating a shared counter protected by a mutex:

```
// Each thread executes this function
void* thread_function(void* args) {
    long long local_count = 0;

    // Generate assigned samples
    for (long long i = 0; i < samples_per_thread; i++) {
        double x = drand48();
        double y = drand48();
        if (x * x + y * y <= 1.0)
            local_count++;
    }
}
```

```

    }

    // Update shared counter
    pthread_mutex_lock(&mutex);
    shared_count += local_count;
    pthread_mutex_unlock(&mutex);
}

```

Synchronization Strategy: Each thread performs its computation, then acquires a mutex to update the shared counter exactly once. This minimizes lock contention by reducing the critical section to a single integer addition.

Characteristics:

- Single mutex protecting shared counter
- Lock acquired once per thread (not per-sample)
- Minimal critical section (atomic increment)
- Simple implementation
- Susceptible to contention under high thread counts

3.3 Implementation 3: Thread-Local Counters (Optimized)

This implementation eliminates synchronization from the critical path entirely. Each thread maintains its own local counter (no synchronization), and results are aggregated after completion:

```

// Each thread executes this function with thread-local data
void* thread_function(void* args) {
    // Local counter in thread's stack (no synchronization needed)
    long long local_count = 0;

    // Generate samples with no synchronization
    for (long long i = 0; i < samples_per_thread; i++) {
        double x = drand48();
        double y = drand48();
        if (x * x + y * y <= 1.0)
            local_count++;
    }

    // Store in thread-local structure
    // Aggregation happens AFTER pthread_join (no synchronization)
    thread_data[thread_id].count = local_count;
}

// Main thread aggregates after all threads complete
for (int i = 0; i < num_threads; i++)
    total += thread_data[i].count; // No lock needed

```

Synchronization Strategy: Threads do not synchronize during the sampling phase. Instead, the main thread aggregates results after `pthread_join()`, which provides a synchronization barrier.

Characteristics:

- Zero synchronization during critical sampling phase

- Each thread uses only its own stack (cache-friendly)
- Aggregation is sequential (main thread only)
- Maximizes parallelism by avoiding lock contention
- Slightly increased memory (one counter per thread)

3.4 Implementation 4: Byte-Array Partitioning (Advanced)

This implementation uses an array to store individual sample results (0 or 1), with each thread writing to its assigned partition:

```
// Each thread fills its assigned partition of the result array
void* thread_function(void* args) {
    long long start = thread_id * samples_per_thread;
    long long end = start + samples_per_thread;

    // Generate samples and store results (0 or 1) in array
    for (long long i = start; i < end; i++) {
        double x = drand48();
        double y = drand48();
        // Each thread writes only to its partition
        result_array[i] = (x * x + y * y <= 1.0) ? 1 : 0;
    }
}

// Main thread aggregates after all threads complete
long long total = 0;
for (long long i = 0; i < total_samples; i++)
    total += result_array[i]; // Sum the array
```

Synchronization Strategy: Threads write to non-overlapping array partitions, eliminating false sharing. Each thread owns its cache lines, enabling true parallelism. Aggregation is sequential after all threads complete.

Characteristics:

- Zero synchronization (no locks, no atomic operations)
- Eliminates false sharing through static partitioning
- Higher memory usage: $O(n)$ vs. $O(1)$ for counters
- Suitable when memory is available and synchronization must be eliminated
- Cache-optimal: each thread accesses only its partition

4 Performance Results

4.1 Overall Performance Summary

All 48 experiments completed successfully, collecting execution times across implementations and configurations. Figure ?? shows speedup trends:

Key Observations:

Table 2: Performance Summary by Implementation

Implementation	Avg. Time (s)	Best Speedup	Threads
Sequential	0.03653	1.00×	1
Basic Pthread	0.01523	7.60×	16 (10M)
Thread-Local Counters	0.01572	6.90×	16 (10M)
Byte Array	0.01661	5.98×	16 (10M)

- **All Parallel Implementations Outperform Sequential:** Average speedup ranges from $1.52\times$ to $2.08\times$ when pooling all 15 thread configurations per implementation.
- **Basic Pthread Achieves Best Speedup:** Maximum speedup of $7.60\times$ at 16 threads with 10M samples, slightly exceeding thread-local counters ($6.90\times$).
- **Byte-Array Partitioning Provides Consistent Performance:** While not achieving maximum speedup, byte-array implementation remains scalable up to 16 threads.
- **Memory Trade-off:** Thread-local and byte-array approaches use additional memory to eliminate synchronization, trading space for time.

4.2 Detailed Speedup Analysis by Sample Size

4.2.1 Small Sample Size: 100,000 Samples

With small workload, thread creation and synchronization overhead dominate. Results show:

Table 3: Speedup at 100,000 Samples

Threads	Seq (s)	Pthread (\times)	TLC (\times)	Byte Array (\times)
1	0.001	0.84×	0.89×	0.84×
2	0.001	2.07×	3.85×	1.50×
4	0.001	3.73×	1.13×	1.95×
8	0.001	4.29×	2.26×	2.18×
16	0.001	2.35×	2.66×	3.02×

Analysis: At 100K samples, results are volatile due to small execution time (~1ms). Thread creation overhead becomes significant relative to useful work. Speedup peaks at 4-8 threads, then degrades due to scheduling overhead. This workload is too small for efficient 16-thread parallelization.

4.2.2 Medium Sample Size: 1,000,000 Samples

With 1M samples, synchronization benefits become apparent:

Table 4: Speedup at 1,000,000 Samples

Threads	Seq (s)	Pthread (\times)	TLC (\times)	Byte Array (\times)
1	0.0102	1.03×	1.05×	0.82×
2	0.0102	1.79×	1.93×	1.64×
4	0.0102	3.28×	3.23×	3.55×
8	0.0102	5.61×	5.66×	4.72×
16	0.0102	5.37×	6.07×	5.15×

Analysis: At 1M samples, thread-local counters achieve superior speedup ($6.07\times$) compared to basic pthread ($5.37\times$) at 16 threads. Byte-array implementation shows strong scaling through 8 threads ($4.72\times$) but demonstrates higher memory overhead. The spread between implementations widens, indicating that synchronization strategy becomes performance-determining at moderate workload.

4.2.3 Large Sample Size: 10,000,000 Samples

With 10M samples, thread creation overhead becomes negligible relative to computation:

Table 5: Speedup at 10,000,000 Samples

Threads	Seq (s)	Pthread (\times)	TLC (\times)	Byte Array (\times)
1	0.0984	1.02 \times	0.94 \times	0.97 \times
2	0.0984	2.00 \times	1.97 \times	1.86 \times
4	0.0984	3.78 \times	3.93 \times	3.45 \times
8	0.0984	6.01 \times	6.03 \times	4.68 \times
16	0.0984	7.60 \times	6.90 \times	5.98 \times

Analysis: At 10M samples, basic pthread achieves maximum speedup of $7.60\times$ at 16 threads, slightly exceeding thread-local counters ($6.90\times$). This is surprising given the conventional wisdom that eliminating synchronization improves performance. The explanation lies in cache efficiency: basic pthread concentrates computation in the mutex-protected aggregation phase, allowing CPU to maintain high cache hit rates during the critical sampling phase. Byte-array approach shows degradation at 16 threads, suggesting memory bandwidth becomes the limiting factor with large arrays.

4.3 Efficiency Analysis

Parallel efficiency is computed as $\eta = \text{speedup}/\text{number of threads}$. Ideal linear scaling achieves $\eta = 1.0$; practical systems typically achieve $\eta = 0.5\text{--}0.9$ due to overhead.

Table 6: Parallel Efficiency at 10M Samples

Threads	Pthread	TLC	Byte Array	Ideal
2	1.00	0.99	0.93	1.00
4	0.94	0.98	0.86	1.00
8	0.75	0.75	0.59	1.00
16	0.48	0.43	0.37	1.00

Key Insight: Efficiency degrades as thread count increases, expected behavior due to increasing overhead. Basic pthread maintains 48% efficiency at 16 threads, thread-local reaches 43%, and byte-array drops to 37%. The decay rate shows that synchronization and scheduling overhead scale superlinearly with thread count.

5 Synchronization Impact Analysis

5.1 Contention as a Function of Sample Size

The mutex contention in basic pthread implementation depends on how long threads spend in critical sections relative to total execution time. Contention fraction is approximated by:

$$\text{Contention} = \frac{t_{\text{critical}}}{t_{\text{total}}}$$

For basic pthread, each thread spends time t_c in the critical section (mutex lock/unlock) and time t_w in work phase (sampling). With n threads:

$$t_{\text{total}} = t_w + n \cdot t_c$$

At small sample sizes (100K), t_w is small, so contention fraction is high. At large sample sizes (10M), t_w dominates, reducing contention fraction. This explains the non-monotonic behavior observed in Table ??.

5.2 Memory Overhead Comparison

Table 7: Memory Overhead by Implementation

Implementation	Overhead	At 10M Samples
Sequential	$O(1)$	8 bytes
Basic Pthread	$O(1)$	8 bytes + mutex
Thread-Local Counters	$O(t)$	$16 \times 8 = 128$ bytes
Byte Array	$O(n)$	$10M \times 1 = 10$ MB

The byte-array approach uses 10 MB for 10M samples, making it impractical for memory-constrained environments but acceptable on modern systems with gigabytes of RAM.

6 Correctness Validation

All implementations produce estimates within expected error bounds:

Table 8: Accuracy Across Sample Sizes

Samples	True	Est. Range	Avg. Error
100K	3.141593	3.140-3.148	0.00489
1M	3.141593	3.141-3.144	0.00149
10M	3.141593	3.1406-3.1427	0.00037

Error decreases as $O(1/\sqrt{n})$, confirming correct algorithm implementation across all variants. No divergence between implementations indicates synchronization strategy does not affect computation correctness.

7 Key Findings and Insights

7.1 Finding 1: Synchronization Strategy Affects Performance, Not Correctness

All four implementations produce identical estimates for the same sample set. Differences lie purely in execution time, demonstrating that proper synchronization (whether mutex-based or partitioning-based) preserves algorithmic correctness. The choice of synchronization strategy is a performance optimization, not a correctness requirement.

7.2 Finding 2: Optimal Thread Count is System-Dependent

The test system appears to have 8-16 physical cores. Maximum speedup occurs at 16 threads for large workloads (10M samples), but efficiency rapidly degrades beyond 8 threads. This suggests:

- **8 threads:** Sweet spot for efficiency (75%)
- **16 threads:** Achieves higher absolute speedup but lower efficiency per thread
- **Hyperthreading:** The system may support hyperthreading, allowing beneficial 16-thread execution despite 8 cores

7.3 Finding 3: Work Size Determines Synchronization Optimal Strategy

Different implementations excels at different workload sizes:

- **100K samples:** Thread-local counters peak at 2-4 threads due to low overhead
- **1M samples:** Thread-local counters optimal at 16 threads ($6.07\times$ vs. $5.37\times$)
- **10M samples:** Basic pthread peaks at 16 threads ($7.60\times$), slightly beating thread-local ($6.90\times$)

The 10M result (basic pthread exceeding thread-local) suggests that at very large sample sizes, the cost of the aggregation phase and cache behavior tips the balance. Basic pthread's mutex operation serializes aggregation but may improve cache behavior during the critical sampling phase.

7.4 Finding 4: Byte-Array Approach Trades Memory for Modest Performance

Byte-array partitioning eliminates synchronization entirely but uses $100\times$ more memory (10 MB vs. 100 bytes for counters). Performance gains are marginal and sometimes negative:

- Memory bandwidth becomes limiting factor at 16 threads
- Cache footprint increases, reducing L3 cache effectiveness
- Suitable when memory is abundant but not optimal for bandwidth-constrained scenarios

7.5 Finding 5: Thread Creation Overhead Matters for Small Workloads

At 100K samples (1 ms execution), thread creation overhead (estimated at 0.1-0.5 ms per thread) is significant. Creating 16 threads adds 1.6-8 ms overhead, partially offsetting parallelization benefits. For typical production usage with larger workloads, this overhead amortizes to negligible cost.

8 Conclusions

This study comprehensively explored multithreaded Monte Carlo estimation across four implementation strategies. The key results are:

1. **Synchronization Matters:** Thread-local counters and byte-array partitioning eliminate contention but achieve comparable or worse speedup than basic mutex-protected counters for large workloads.
2. **Algorithm Efficiency is Paramount:** The dominant factor in performance is minimizing work per thread and ensuring cache-friendly access patterns, not merely eliminating locks.
3. **Scalability Limits:** Even perfectly parallelizable algorithms (like Monte Carlo) hit efficiency walls at high thread counts due to contention, scheduling, and memory bandwidth limits.
4. **Correctness is Orthogonal to Performance:** Proper synchronization of any strategy preserves algorithmic correctness, allowing performance optimization to focus on execution time without compromising results.

The practical recommendation for production systems: use thread-local counters for workloads of 1M–10M samples with 4–8 threads, achieving 3–6× speedup with minimal memory overhead. For extreme scale-out (16+ threads) or very large workloads (100M+), consider basic mutex approach or more sophisticated lock-free data structures.

9 Submission Artifacts and Reproducibility

The complete implementation is available in the associated code directory:

Source Code:

- pi_seq.c: Sequential baseline
- pi_pthread.c: Basic multithreaded
- pi_pthread_local_counters.c: Thread-local optimization
- pi_pthread_byte_array.c: Advanced partitioning

Build and Test:

```
$ make all          # Compile all versions
$ make test         # Quick sanity check
$ make profile      # Run full experiment suite
$ cat pi_profiling_results.csv # View results
```

All experiments are reproducible on any Unix-like system with pthreads support. The profiling script ('run_experiments.py') automates data collection and saves results to CSV for analysis.