

Introduction to CUDA Parallel Programming

Homework Assignment 6

- April, 2025
- NTNU
- 41173058H

Problem Statement

Implement a pseudo-random number generator to generate random numbers in $(0, \infty)$ with exponential distribution e^{-x} , creating a dataset of 81,920,000 entries. Compare histogram computations using CPU, GPU with global memory, and GPU with shared memory implementations. Analyze their performance characteristics and determine optimal block sizes. Plot the histogram alongside the theoretical probability distribution curve.

Mathematical Foundation

Exponential Distribution

The exponential distribution with rate parameter λ has the probability density function:

$$f(x) = \lambda e^{-\lambda x}, \quad x \geq 0$$

For $\lambda = 1$ (standard exponential distribution):

$$f(x) = e^{-x}$$

Random Number Generation

The inverse transform method generates exponential random variables using:

$$X = -\frac{1}{\lambda} \ln(U)$$

where U is a uniform random variable on $(0, 1)$.

Source Code Analysis

Random Number Generation

```
void generate_exponential_data(std::vector<float>& data, float lambda = 1.0f) {  
    std::mt19937 rng(12345); // Mersenne Twister with fixed seed  
    std::exponential_distribution<float> exp_dist(lambda);  
    for (size_t i = 0; i < data.size(); ++i) {
```

```

        data[i] = exp_dist(rng);
    }
}

```

Key Features:

1. **Mersenne Twister:** High-quality pseudo-random number generator
2. **Fixed Seed:** Ensures reproducible results across runs
3. **Standard Library:** Uses C++11 `std::exponential_distribution`

CPU Histogram Implementation

```

void cpu_histogram(const std::vector<float>& data, std::vector<int>& hist, float bin_width)
{
    std::fill(hist.begin(), hist.end(), 0);
    for (size_t i = 0; i < data.size(); ++i) {
        int bin = std::min(int(data[i] / bin_width), NUM_BINS - 1);
        hist[bin]++;
    }
}

```

Algorithm Characteristics:

- **Sequential Processing:** $O(N)$ time complexity
- **Bin Calculation:** Linear mapping with overflow protection
- **Memory Access:** Sequential, cache-friendly pattern

GPU Global Memory Implementation

```

__global__ void histogram_global_kernel(const float* data, int* hist, int n, float bin_width)
{
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if (idx < n) {
        int bin = min(int(data[idx] / bin_width), NUM_BINS - 1);
        atomicAdd(&hist[bin], 1);
    }
}

```

Implementation Features:

1. **Global Memory Atomics:** Direct updates to global histogram
2. **Thread Mapping:** One thread per data element
3. **Atomic Contention:** High contention for popular bins
4. **Memory Bandwidth:** Limited by global memory latency

GPU Shared Memory Implementation

```

__global__ void histogram_shared_kernel(const float* data, int* hist, int n, float bin_width)
{
    __shared__ int local_hist[NUM_BINS];
    int tid = threadIdx.x;

```

```

// Initialize shared memory histogram
if (tid < NUM_BINS) local_hist[tid] = 0;
__syncthreads();

// Compute local histogram
int idx = blockIdx.x * blockDim.x + threadIdx.x;
if (idx < n) {
    int bin = min(int(data[idx] / bin_width), NUM_BINS - 1);
    atomicAdd(&local_hist[bin], 1);
}
__syncthreads();

// Merge to global histogram
if (tid < NUM_BINS) {
    atomicAdd(&hist[tid], local_hist[tid]);
}
}

```

Optimization Strategy:

1. **Two-Level Reduction:** Block-local then global aggregation
2. **Shared Memory Atomics:** Faster than global memory atomics
3. **Reduced Global Contention:** Only NUM_BINS global atomic operations per block
4. **Memory Hierarchy Utilization:** Exploits shared memory bandwidth

Results and Analysis

Performance Comparison

Block Size	GPU Global (ms)	GPU Shared (ms)	Global Speedup	Shared Speedup
8	44.509	35.799	3.66×	4.55×
16	33.818	19.047	4.82×	8.55×
32	33.713	19.126	4.83×	8.51×
64	33.881	19.170	4.81×	8.50×
128	33.845	9.549	4.81×	17.05×
256	33.782	5.022	4.82×	32.43×
512	33.605	3.883	4.85×	41.94×
1024	33.561	3.986	4.85×	40.85×

Baseline Performance:

- **CPU Histogram Time:** 162.851 ms
- **Optimal GPU Global:** Block size 1024 (33.561 ms, 4.85× speedup)
- **Optimal GPU Shared:** Block size 512 (3.883 ms, 41.94× speedup)

Performance Analysis

Why Shared Memory Dramatically Outperforms Global Memory Atomic Operation Efficiency:

- **Shared Memory Atomics:** ~20-100× faster than global memory atomics
- **Memory Latency:** Shared memory ~1-2 cycles vs. global memory ~400-800 cycles
- **Bandwidth:** Shared memory ~1.5 TB/s vs. global memory ~900 GB/s

Contention Reduction:

- **Global Method:** All threads across all blocks compete for same histogram bins
- **Shared Method:** Only threads within same block compete, then single atomic per bin per block

Block Size Impact Analysis Small Block Sizes (8-64):

- **Poor Occupancy:** Insufficient parallelism to hide memory latency
- **Underutilization:** GPU resources not fully utilized
- **Higher Overhead:** More blocks increase kernel launch overhead

Medium Block Sizes (128-256):

- **Improved Occupancy:** Better resource utilization
- **Balanced Contention:** Optimal balance between parallelism and atomic conflicts
- **Shared Memory Efficiency:** Better utilization of shared memory bandwidth

Large Block Sizes (512-1024):

- **Optimal for Shared Memory:** Block size 512 provides best performance
- **Diminishing Returns:** Block size 1024 shows slight performance degradation
- **Resource Constraints:** Register pressure and shared memory limits

Exponential Distribution Characteristics Data Distribution Impact:

- **Heavy Concentration:** Most values fall in first few bins (exponential decay)
- **Atomic Contention:** High contention for bins near zero
- **Cache Behavior:** Poor spatial locality due to exponential distribution

Numerical Accuracy Validation

Correctness Verification:

- **CPU vs GPU (global) mismatched bins:** 0
- **CPU vs GPU (shared) mismatched bins:** 0

- **Perfect Agreement:** All implementations produce identical results

Error Sources Eliminated:

1. **Floating-Point Precision:** Consistent across all implementations
2. **Atomic Race Conditions:** Proper synchronization prevents data races
3. **Bin Calculation:** Identical binning logic across all versions

Discussion

Algorithm Efficiency Analysis

Computational Complexity

- **CPU:** $O(N)$ sequential processing
- **GPU Global:** $O(N/P)$ parallel processing with atomic contention
- **GPU Shared:** $O(N/P + B)$ where B is bins per block reduction

Memory Access Patterns CPU Implementation:

```
// Sequential access, cache-friendly
for (size_t i = 0; i < data.size(); ++i) {
    hist[bin]++; // Potential cache misses for histogram array
}
```

GPU Global Implementation:

```
// Coalesced input reads, random histogram writes
atomicAdd(&hist[bin], 1); // High contention, poor cache behavior
```

GPU Shared Implementation:

```
// Coalesced input reads, local shared memory updates
atomicAdd(&local_hist[bin], 1); // Low latency, reduced contention
atomicAdd(&hist[tid], local_hist[tid]); // Minimal global updates
```

Performance Optimization Strategies

Current Optimizations

1. **Shared Memory Utilization:** Reduces global memory traffic by factor of `block_size`
2. **Atomic Operation Reduction:** Minimizes expensive global atomics
3. **Memory Coalescing:** Ensures efficient data loading patterns
4. **Block-Level Aggregation:** Exploits memory hierarchy effectively

Further Optimization Opportunities Warp-Level Primitives:

```
// Use warp shuffle for intra-warp reduction
__device__ int warpReduceSum(int val) {
    for (int offset = warpSize/2; offset > 0; offset /= 2) {
```

```

        val += __shfl_down_sync(0xFFFFFFFF, val, offset);
    }
    return val;
}

```

Multi-Pass Strategy:

```

// For very large bin counts, use multiple passes
if (NUM_BINS > SHARED_MEMORY_LIMIT) {
    // Process histogram in chunks that fit in shared memory
    for (int chunk = 0; chunk < NUM_BINS; chunk += CHUNK_SIZE) {
        // Process bins [chunk, chunk+CHUNK_SIZE)
    }
}

```

Cooperative Groups:

```

#include <cooperative_groups.h>
// Use cooperative groups for more flexible synchronization
namespace cg = cooperative_groups;
auto block = cg::this_thread_block();

```

Scalability Analysis

Data Size Scaling

- **Memory Bandwidth Bound:** Performance scales with GPU memory bandwidth
- **Atomic Contention:** Increases with data size for popular bins
- **Cache Effects:** Larger datasets may exceed cache capacity

Bin Count Scaling

- **Shared Memory Limit:** 128 bins \times 4 bytes = 512 bytes (well within 48KB limit)
- **Atomic Conflicts:** Reduce with more bins (better distribution)
- **Memory Footprint:** Linear scaling with bin count

Real-World Applications

Scientific Computing:

- **Monte Carlo Simulations:** Statistical analysis of random processes
- **Image Processing:** Intensity histograms for enhancement algorithms
- **Signal Processing:** Amplitude distribution analysis

Machine Learning:

- **Feature Engineering:** Data distribution analysis
- **Anomaly Detection:** Outlier identification through histogram analysis
- **Preprocessing:** Data normalization and scaling

Theoretical vs. Experimental Validation

Expected Distribution Properties

For exponential distribution with $\lambda = 1$:

- **Mean:** $\mu = 1/\lambda = 1$
- **Variance:** $\sigma^2 = 1/\lambda^2 = 1$
- **Mode:** 0 (maximum at $x = 0$)
- **Median:** $\ln(2) \approx 0.693$

Histogram Visualization Strategy

Plotting Implementation:

```
import matplotlib.pyplot as plt
import numpy as np

# Theoretical exponential PDF
x = np.linspace(0, max_value, 1000)
theoretical_pdf = np.exp(-x)

# Normalize histogram to probability density
bin_centers = np.arange(NUM_BINS) * bin_width + bin_width/2
histogram_density = histogram_counts / (DATA_SIZE * bin_width)

# Plot comparison
plt.figure(figsize=(12, 8))
plt.bar(bin_centers, histogram_density, width=bin_width*0.8,
        alpha=0.7, label='Computed Histogram')
plt.plot(x, theoretical_pdf, 'r-', linewidth=2,
        label='Theoretical exp(-x)')
plt.xlabel('x')
plt.ylabel('Probability Density')
plt.title('Exponential Distribution: Computed vs. Theoretical')
plt.legend()
plt.grid(True, alpha=0.3)
plt.show()
```

Conclusion

The GPU shared memory implementation achieves exceptional performance with a **41.94× speedup** over CPU computation using the optimal block size of 512 threads. This dramatic improvement results from:

Key Performance Factors:

1. **Shared Memory Efficiency:** 20-100× faster atomic operations than global memory

2. **Contention Reduction:** Block-local aggregation minimizes global atomic conflicts
3. **Memory Hierarchy Utilization:** Exploits GPU's memory subsystem effectively
4. **Optimal Block Size:** 512 threads balance parallelism with resource constraints

Technical Achievements:

- **Perfect Numerical Accuracy:** All implementations produce identical results
- **Scalable Architecture:** Efficient handling of 81.92M data points
- **Comprehensive Analysis:** Systematic block size optimization
- **Robust Implementation:** Proper synchronization and error handling

Recommendations for Production Use:

1. **Use shared memory approach** for histogram computations with moderate bin counts
2. **Optimize block size** based on specific GPU architecture and problem characteristics
3. **Consider multi-pass strategies** for very large bin counts exceeding shared memory
4. **Implement proper validation** to ensure numerical correctness

The implementation demonstrates the effectiveness of GPU computing for embarrassingly parallel problems and highlights the critical importance of memory hierarchy optimization in achieving peak performance.

Submission Guidelines

Submit your homework report including source codes, results, and discussions (without *.exe files). Prepare the discussion file using a typesetting system (LaTeX, Word, etc.) and convert to PDF. Compress all files into a gzipped tar file named with your student number and problem set number (e.g., `r05202043_HW6.tar.gz`). Send your homework from your NTU/NTNU/NTUST email account to `twchiu@phys.ntu.edu.tw` before 17:00, June 11, 2025. If email attachment fails, place the file in your `twcp1` home directory and send notification email.