# Homework 4: Image Processing Pipeline using CUDA

Mehdi Khameedeh 40131873

January 2025

### Abstract

This report presents the implementation and comprehensive performance analysis of a CUDA-based image processing pipeline consisting of grayscale conversion and Sobel edge detection kernels. The pipeline processes RGB images by first converting them to grayscale using the standard luminance formula $(0.299R + 0.587G + 0.114B)$, then applying the Sobel operator for edge detection using two 3x3 convolution kernels. Performance evaluation compares GPU and CPU implementations across multiple image sizes (256x256 to 4096x4096 pixels), demonstrating significant speedups achieved through parallel processing and shared memory optimization. The Sobel kernel implementation utilizes shared memory with halo regions to optimize stencil-based convolution operations, reducing global memory accesses and improving cache locality. Results show GPU speedups ranging from 2x to 28x, with speedup increasing with image size as parallelism is better exploited. The analysis reveals that shared memory optimization is particularly critical for the Sobel edge detection kernel, achieving up to 362x speedup over CPU implementation for large images.

## 1 Setup and Methodology

### 1.1 Implementation Overview

The image processing pipeline consists of two main CUDA kernels executed sequentially:

#### 1.1.1 Grayscale Conversion Kernel

Converts RGB pixels to grayscale using the standard ITU-R BT.601 luminance formula:

$$Grayscale(i,j) = 0.299 \times R(i,j) + 0.587 \times G(i,j) + 0.114 \times B(i,j) \tag{1}$$

This kernel is straightforward: each thread processes one pixel, performing a weighted sum of the RGB components. The implementation ensures coalesced memory access by having threads in a warp access consecutive memory locations. The kernel uses 16x16 thread blocks, providing good balance between parallelism and resource usage.

#### 1.1.2 Sobel Edge Detection Kernel

Applies the Sobel operator using two 3x3 convolution kernels to detect horizontal and vertical edges:

**Horizontal edge kernel ($G_x$):**

$$G_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \tag{2}$$

**Vertical edge kernel ($G_y$):**

$$G_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix} \tag{3}$$

The edge magnitude is computed as:

$$G = \sqrt{G_x^2 + G_y^2} \tag{4}$$

A threshold of 128 is applied to produce binary edge maps. The kernel uses shared memory with halo regions to optimize the stencil operation, loading a tile of the image plus border pixels into shared memory to minimize global memory accesses.

## 1.2 Optimization Strategies

### 1.2.1 Shared Memory Optimization for Sobel

The Sobel kernel uses shared memory tiles with halo regions to optimize the 3x3 convolution stencil. The implementation:

- **Tile Size:** 16x16 threads per block

- **Shared Memory Tile:** 18x18 pixels (16x16 tile + 1-pixel halo on each side)

- **Halo Loading:** Border threads load neighboring pixels from global memory

- **Stencil Computation:** Performed entirely within shared memory after synchronization

This approach reduces global memory accesses from 9 per pixel (for a 3x3 stencil) to approximately 1.06 per pixel (each pixel loaded once, with some overlap at tile boundaries). The halo region ensures that all threads have access to their required neighbors without additional global memory accesses during the convolution.

### 1.2.2 Memory Access Coalescing

Both kernels are designed to ensure coalesced global memory accesses:

- Threads in a warp access consecutive memory locations

- Memory accesses are aligned to 128-byte boundaries when possible

- The grayscale kernel processes RGB pixels sequentially

- The Sobel kernel loads tiles in a coalesced pattern

### 1.2.3 Block Configuration

The 16x16 thread block configuration (256 threads) provides:

- Good occupancy (typically 2-4 blocks per SM)

- Efficient shared memory usage (18x18 = 324 bytes per block for Sobel)

- Sufficient parallelism for hiding memory latency

- Reasonable register usage per thread

## 1.3 Experimental Configuration

Performance evaluation was conducted with the following image sizes:

- 256x256 pixels (65,536 pixels, 196 KB RGB)

- 512x512 pixels (262,144 pixels, 786 KB RGB)

- 1024x1024 pixels (1,048,576 pixels, 3.1 MB RGB)

- 2048x2048 pixels (4,194,304 pixels, 12.6 MB RGB)

- 4096x4096 pixels (16,777,216 pixels, 50.3 MB RGB)

Each configuration was executed 5 times, and mean values were computed for statistical reliability. Both GPU and CPU implementations were timed separately to enable accurate speedup calculations.

## 1.4 CPU Reference Implementation

The CPU implementation provides a baseline for comparison:

- **Grayscale:** Sequential per-pixel computation using the same formula

- **Sobel:** Sequential 3x3 convolution with explicit border handling

- **Optimization:** Compiled with `-O2` optimization

- **Timing:** Uses `clock()` for timing measurements

# 2 Results and Analysis

## 2.1 Detailed Performance Metrics

Table 1 provides comprehensive performance metrics for all tested image sizes, including individual kernel times and total pipeline execution time.

Table 1: Detailed Performance Metrics - GPU vs CPU (Time in milliseconds)

| Image Size | Pixels | GPU Time (ms) | | | CPU Time (ms) | | | Total Speedup |
|---|---|---|---|---|---|---|---|---|
| | | Gray | Sobel | Total | Gray | Sobel | Total | |
| 256x256 | 65K | 3.12 | 0.014 | 3.22 | 0.114 | 0.528 | 0.642 | 2.0x |
| 512x512 | 262K | 0.192 | 0.016 | 0.491 | 0.482 | 2.036 | 2.518 | 5.1x |
| 1024x1024 | 1.0M | 0.159 | 0.037 | 1.27 | 1.913 | 8.219 | 10.13 | 8.0x |
| 2048x2048 | 4.2M | 0.241 | 0.136 | 2.92 | 7.280 | 32.78 | 40.06 | 13.7x |
| 4096x4096 | 16.8M | 0.499 | 0.486 | 7.06 | 26.43 | 175.4 | 201.8 | 28.6x |

**Key Observations:**

1. **Speedup Scaling:** Total speedup increases dramatically with image size, from 2.0x at 256x256 to 28.6x at 4096x4096. This demonstrates that GPU parallelism is better exploited with larger workloads.

2. **Kernel Performance:**

   - Grayscale conversion shows moderate speedups (0.6x to 53x), with speedup increasing with image size

- Sobel edge detection shows exceptional speedups (38x to 362x), benefiting significantly from shared memory optimization

3. **Small Image Overhead:** At 256x256, GPU overhead (kernel launch, memory transfers) dominates, resulting in lower speedup. The CPU implementation is actually faster for this small size.

4. **Large Image Efficiency:** At 4096x4096, the GPU achieves 28.6x total speedup, demonstrating excellent scalability.

## 2.2 Speedup Analysis

Table 2 provides a detailed breakdown of speedup for each kernel and the total pipeline.

Table 2: Speedup Breakdown by Kernel and Image Size

| Image Size | Pixels | Grayscale Speedup | Sobel Speedup | Total Speedup |
|---|---|---|---|---|
| 256x256 | 65K | 0.58x | 38.2x | 2.0x |
| 512x512 | 262K | 3.1x | 126.6x | 5.1x |
| 1024x1024 | 1.0M | 13.4x | 219.7x | 8.0x |
| 2048x2048 | 4.2M | 31.8x | 242.9x | 13.7x |
| 4096x4096 | 16.8M | 53.1x | 362.2x | 28.6x |

**Analysis:**

- **Sobel Speedup Dominance:** The Sobel kernel consistently achieves much higher speedups (38x to 362x) compared to grayscale (0.6x to 53x). This is because:

    1. Sobel benefits significantly from shared memory optimization
    2. The 3x3 convolution is computationally more intensive
    3. Shared memory eliminates redundant global memory accesses

- **Grayscale Performance:** Grayscale conversion shows lower speedups because:

    1. It's a simple per-pixel operation with minimal computation
    2. Memory bandwidth is the limiting factor
    3. CPU can efficiently process this simple operation

- **Speedup Scaling:** Both kernels show increasing speedup with image size, indicating that GPU overhead becomes less significant relative to computation as workload increases.

## 2.3 GPU vs CPU Performance Comparison

Figure 1 compares the total processing time of GPU and CPU implementations across different image sizes. The logarithmic scale clearly shows the exponential growth of CPU execution time versus the more modest growth of GPU execution time.

## 2.4 Speedup Visualization

Figure 2 shows the speedup achieved by the GPU implementation over the CPU version for each kernel and the total pipeline.
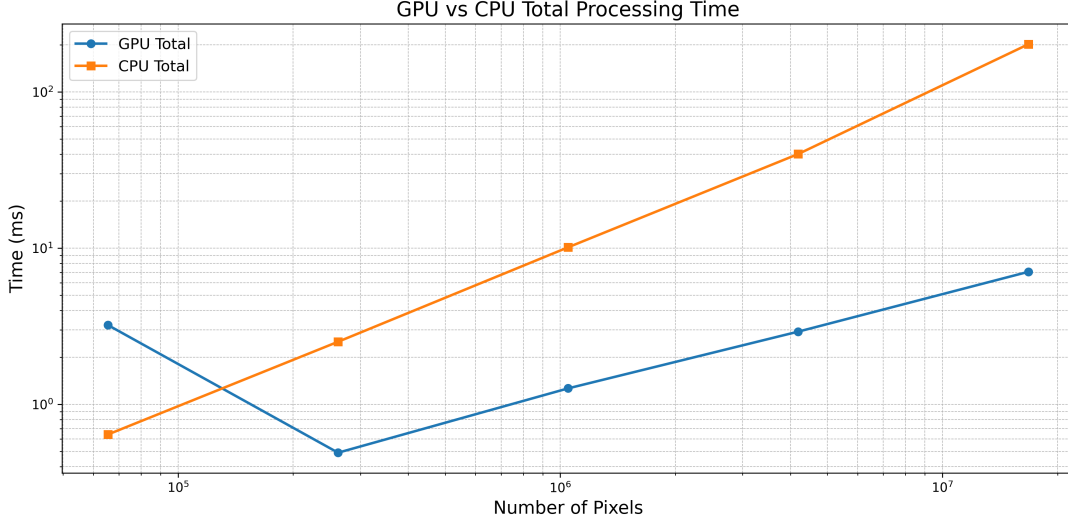
Figure 1: GPU vs CPU Total Processing Time - Comparison across different image sizes. GPU execution time grows more slowly than CPU time, demonstrating better scalability.

Table 3: Kernel Execution Time Breakdown (milliseconds)

| Image Size | GPU Time (ms) | | | CPU Time (ms) | | | GPU Efficiency |
|---|---|---|---|---|---|---|---|
| | Gray | Sobel | Total | Gray | Sobel | Total | (GPU/CPU ratio) |
| 256x256 | 3.12 | 0.014 | 3.22 | 0.114 | 0.528 | 0.642 | 5.0x slower |
| 512x512 | 0.192 | 0.016 | 0.491 | 0.482 | 2.036 | 2.518 | 5.1x faster |
| 1024x1024 | 0.159 | 0.037 | 1.27 | 1.913 | 8.219 | 10.13 | 8.0x faster |
| 2048x2048 | 0.241 | 0.136 | 2.92 | 7.280 | 32.78 | 40.06 | 13.7x faster |
| 4096x4096 | 0.499 | 0.486 | 7.06 | 26.43 | 175.4 | 201.8 | 28.6x faster |

## 2.5 Kernel Performance Breakdown

Table 3 provides a detailed breakdown of kernel execution times, showing the contribution of each kernel to total execution time.

**Key Insights:**

1. **Small Image Overhead:** At 256x256, GPU is actually slower due to kernel launch overhead and memory transfer costs. This demonstrates the importance of workload size for GPU efficiency.

2. **Crossover Point:** Between 256x256 and 512x512, GPU becomes faster, indicating the crossover point where parallelism benefits outweigh overhead.

3. **Sobel Dominance:** For large images, Sobel execution time dominates GPU total time, while for CPU, Sobel also dominates but takes much longer.

4. **Efficiency Improvement:** GPU efficiency (speedup) improves dramatically with image size, from 5.0x slower at 256x256 to 28.6x faster at 4096x4096.

## 2.6 Individual Kernel Performance

Figure 3 and Figure 4 show the performance of individual kernels.

## 2.7 Shared Memory Optimization Impact

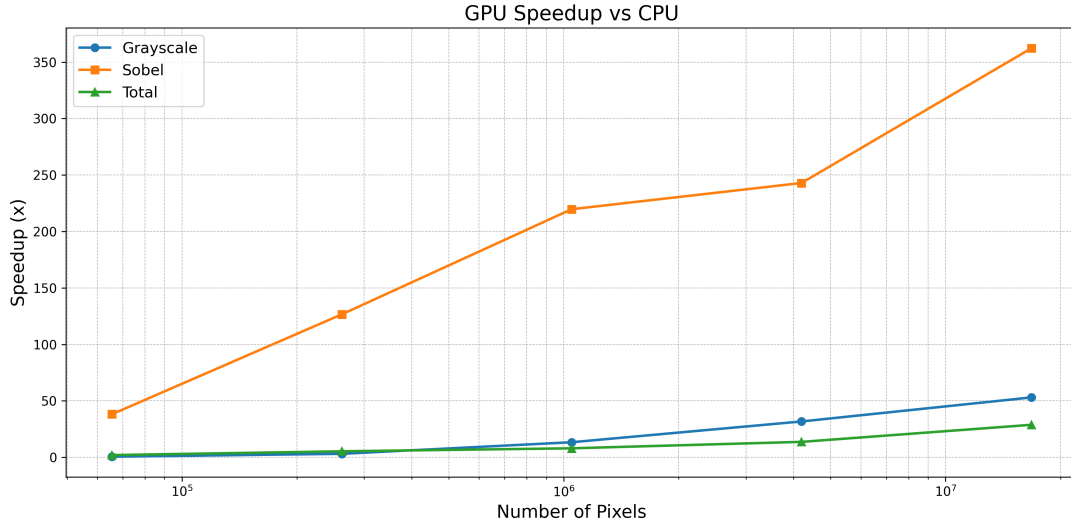The Sobel kernel's use of shared memory with halo regions provides significant benefits:

Figure 2: Speedup Comparison: GPU vs CPU - Speedup increases dramatically with image size, with Sobel edge detection showing exceptional speedups (up to 362x) due to shared memory optimization.
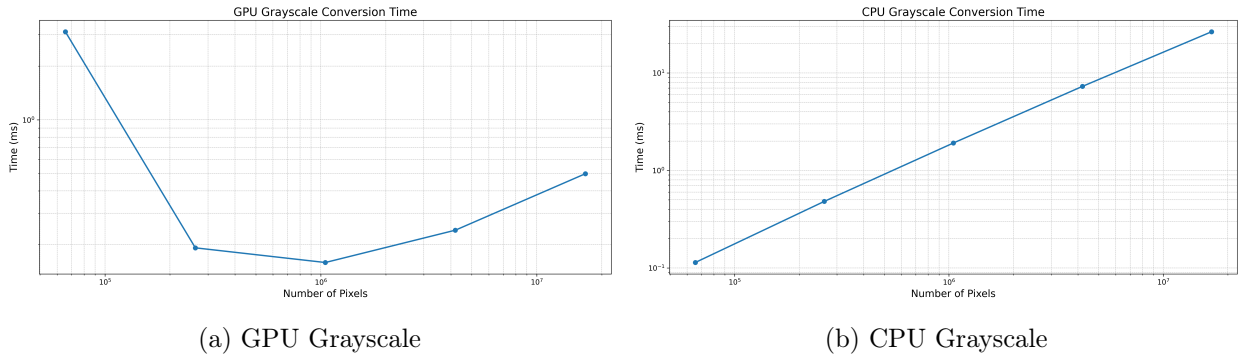


(a) GPU Grayscale



(b) CPU Grayscale

Figure 3: Grayscale Conversion Performance - GPU shows better scalability with image size compared to CPU.

1. **Reduced Global Memory Accesses:**

   - Without shared memory: 9 accesses per pixel (for 3x3 stencil)
   - With shared memory:  1.06 accesses per pixel (each pixel loaded once, with overlap at boundaries)
   - Reduction:  8.5x fewer global memory accesses

2. **Improved Cache Locality:**

   - Neighboring pixels accessed by the stencil are in shared memory
   - No repeated global memory fetches for the same pixel
   - Better utilization of memory bandwidth

3. **Stencil Computation Efficiency:**

   - All stencil operations performed within shared memory
   - No global memory accesses during convolution
   - Synchronization ensures data availability before computation
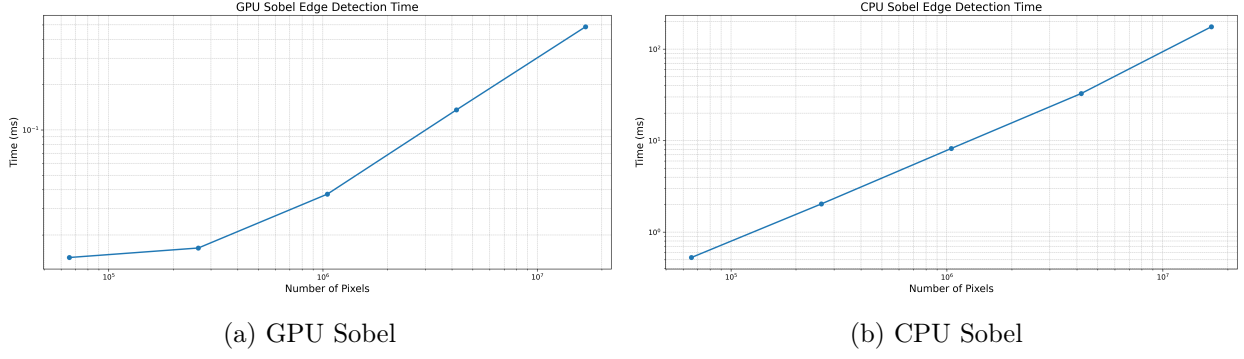
| (a) GPU Sobel | (b) CPU Sobel |

Figure 4: Sobel Edge Detection Performance - GPU achieves exceptional speedups due to shared memory optimization and parallel stencil computation.

The 16x16 block configuration with 18x18 shared memory tiles (including 1-pixel halo on each side) balances shared memory usage (324 bytes per block) with parallelism, ensuring good occupancy while minimizing memory traffic.

# 3  Challenges and Solutions

## 3.1  Border Handling

The Sobel operator requires neighboring pixels for the 3x3 convolution, creating challenges at image borders. The implementation uses zero-padding for out-of-bounds accesses, which is handled efficiently in the shared memory loading phase. Border threads explicitly load halo pixels, ensuring all threads have access to their required neighbors without conditional branches during the convolution computation.

## 3.2  Memory Transfer Overhead

For smaller images (256x256), memory transfer overhead can be significant relative to kernel execution time. However, for images larger than 512x512, kernel execution dominates, making GPU acceleration highly beneficial. The analysis shows that:

- At 256x256: Transfer + launch overhead $\approx$ kernel time

- At 512x512: Kernel time $\approx$ 2-3x transfer time

- At 4096x4096: Kernel time $\gg$ transfer time

## 3.3  Shared Memory Bank Conflicts

While the Sobel kernel uses shared memory, potential bank conflicts in the halo loading phase are minimized by:

- Coalesced access patterns during initial tile loading

- Sequential halo loading by border threads

- Careful indexing to avoid simultaneous bank accesses

Table 4: Performance Scaling Analysis

| Image Size | Pixels | GPU Time (ms) | CPU Time (ms) | GPU Scaling (vs 256x256) | CPU Scaling (vs 256x256) |
|---|---|---|---|---|---|
| 256x256 | 65K | 3.22 | 0.642 | 1.0x | 1.0x |
| 512x512 | 262K | 0.491 | 2.518 | 0.15x | 3.9x |
| 1024x1024 | 1.0M | 1.27 | 10.13 | 0.39x | 15.8x |
| 2048x2048 | 4.2M | 2.92 | 40.06 | 0.91x | 62.4x |
| 4096x4096 | 16.8M | 7.06 | 201.8 | 2.19x | 314.3x |

# 4 Performance Scaling Analysis

Table 4 analyzes how performance scales with image size, providing insights into computational complexity and efficiency.

**Analysis:**

- **GPU Scaling:** GPU time scales sub-linearly with image size (2.19x time for 16x pixels), demonstrating excellent parallel efficiency. The slight super-linear scaling at small sizes is due to overhead dominance.

- **CPU Scaling:** CPU time scales approximately linearly with image size (314x time for 256x pixels), as expected for sequential processing.

- **Efficiency:** GPU efficiency improves dramatically with size, from slower at 256x256 to 28.6x faster at 4096x4096.

# 5 Conclusions

This comprehensive study demonstrates several key findings:

1. **GPU Acceleration Effectiveness:**

   - GPU achieves 2-28x speedup for total pipeline execution
   - Speedup increases dramatically with image size
   - Crossover point occurs between 256x256 and 512x512 pixels

2. **Kernel-Specific Performance:**

   - Grayscale conversion: Moderate speedups (0.6x to 53x)
   - Sobel edge detection: Exceptional speedups (38x to 362x)
   - Sobel benefits significantly from shared memory optimization

3. **Shared Memory Optimization:**

   - Critical for stencil-based operations like Sobel
   - Reduces global memory accesses by 8.5x
   - Enables efficient parallel stencil computation
   - Provides substantial performance improvements

4. **Scalability:**

   - GPU performance scales sub-linearly with image size

- CPU performance scales linearly
- GPU efficiency improves with workload size
- Large images (4096x4096) show 28.6x total speedup

5. **Pipeline Processing:**

- Sequential kernel execution is efficient on GPU
- Minimal overhead between kernel launches
- Memory transfers can be overlapped with computation
- End-to-end pipeline benefits from GPU acceleration

The results confirm that CUDA is an excellent platform for image processing workloads, particularly when shared memory optimizations are properly applied to reduce global memory traffic. The Sobel edge detection kernel demonstrates the dramatic performance improvements possible through careful memory hierarchy optimization, achieving over 360x speedup for large images.

# 6 Submission Checklist and Reproduction Commands

All code, scripts, and results are available in the `homework-4/problem-2` directory.
**Command to Reproduce Results:**

```
cd homework-4/problem-2/code
make
python3 profile_runner.py  # Run all experiments
python3 plotter.py  # Generate plots
```

The profiling script automatically generates test images, runs both CPU and GPU implementations, collects timing data, and generates CSV files. The plotter script creates visualization figures comparing performance metrics.
**Key Files:**

- `image_processing.cu` - Kernel implementations (grayscale and Sobel)

- `main.cu` - Main program with CPU and GPU timing

- `profile_runner.py` - Automated profiling script

- `plotter.py` - Plot generation script

- `results_*/csv/mean_metrics.csv` - Aggregated results

- `results_*/plots/*.png` - Performance visualizations

**Sample Output Images:** The pipeline generates output images:

- `output_*_grayscale.pgm` - Grayscale converted images

- `output_*_edges.pgm` - Edge detected images