

GPU Sobel Edge Detection: Performance Analysis and Optimization

Thomas Giesbrecht
Avirup Bhattacharjee
Matthew Moran

CSE 4373/5373: General Purpose GPU Programming
Instructor: Dr. Alex J Dillhoff
April 23, 2024

Project Overview



Implementation of the Sobel edge detection algorithm

Four versions implemented:

- CPU Implementation
- GPU Naive Implementation
- GPU Optimized Implementation
- GPU Coarsened Implementation

Technologies used: **CUDA Toolkit, OpenCV**



CPU IMPLEMENTATION

- Standard convolution approach using Sobel kernels (Gx, Gy)
- Nested loops for pixel-by-pixel computation
- Gradient magnitude calculated and clamped
- Output image saved as `cpu_output_image.jpg`

```
7 void applySobel(const cv::Mat& input, cv::Mat& output) {
8     // Sobel kernels
9     int Gx[3][3] = {
10         {-1, 0, 1},
11         {-2, 0, 2},
12         {-1, 0, 1}
13     };
14
15     int Gy[3][3] = {
16         {-1, -2, -1},
17         { 0,  0,  0},
18         { 1,  2,  1}
19     };
20
21     // Initialize the output image
22     output = cv::Mat::zeros(input.size(), CV_8UC1);
23
24     // Apply the Sobel filter
25     for (int y = 1; y < input.rows - 1; ++y) {
26         for (int x = 1; x < input.cols - 1; ++x) {
27             int sumX = 0;
28             int sumY = 0;
29
30             // Convolution with Sobel kernels
31             for (int ky = -1; ky <= 1; ++ky) {
32                 for (int kx = -1; kx <= 1; ++kx) {
33                     int pixel = input.at<uchar>(y + ky, x + kx);
34                     sumX += pixel * Gx[ky + 1][kx + 1];
35                     sumY += pixel * Gy[ky + 1][kx + 1];
36                 }
37             }
38
39             // Compute gradient magnitude
40             int magnitude = std::sqrt(sumX * sumX + sumY * sumY);
41             magnitude = std::min(255, magnitude); // Clamp to 255
42
43             output.at<uchar>(y, x) = static_cast<uchar>(magnitude);
44         }
45     }
46 }
```

Naive GPU Implementation

- Basic CUDA kernel using direct global memory access
- Each thread computes gradient magnitude independently
- Coalesced global memory access
- Output image saved as `gpu_output_image_naive.jpg`

```
__global__ void sobelKernel(const unsigned char* input, unsigned char* output, int w, int h) {
    // Calculate the thread's position in the image
    int x = blockIdx.x * blockDim.x + threadIdx.x;
    int y = blockIdx.y * blockDim.y + threadIdx.y;

    // Sobel kernels
    int Gx[3][3] = {
        {-1, 0, 1},
        {-2, 0, 2},
        {-1, 0, 1}
    };

    int Gy[3][3] = {
        {-1, -2, -1},
        { 0,  0,  0},
        { 1,  2,  1}
    };

    // Ensure the thread is within bounds and not on the border
    if (x > 0 && x < width - 1 && y > 0 && y < height - 1) {
        int sumX = 0;
        int sumY = 0;

        // Apply the Sobel filter
        for (int ky = -1; ky <= 1; ++ky) {
            for (int kx = -1; kx <= 1; ++kx) {
                int pixel = input[(y + ky) * width + (x + kx)];
                sumX += pixel * Gx[ky + 1][kx + 1];
                sumY += pixel * Gy[ky + 1][kx + 1];
            }
        }

        // Compute gradient magnitude
        int magnitude = sqrtf(sumX * sumX + sumY * sumY);
        magnitude = min(255, magnitude); // Clamp to 255

        // Write the result to the output image
        output[y * width + x] = static_cast<unsigned char>(magnitude);
    }
}
```

Optimized GPU Implementation

- Utilizes shared memory to minimize global memory accesses
- 18x18 shared memory tile for 16x16 pixel blocks, loading additional border pixels
- Thread synchronization (`__syncthreads()`) for data consistency
- Output image saved as `gpu_output_image_optimized.jpg`

```
1 // Global variables
2 __global__ void sobelKernelOptimized(const unsigned char* input, unsigned char* out)
3 {
4     // Shared memory for the tile
5     __shared__ unsigned char sharedMem[18][18]; // Block size (16x16) + 2 for border
6
7     // Calculate global thread coordinates
8     int x = blockIdx.x * blockDim.x + threadIdx.x;
9     int y = blockIdx.y * blockDim.y + threadIdx.y;
10
11     // Calculate shared memory coordinates
12     int sharedX = threadIdx.x + 1;
13     int sharedY = threadIdx.y + 1;
14
15     // Load data into shared memory
16     if (x < width && y < height) {
17         sharedMem[sharedY][sharedX] = input[y * width + x];
18     }
19
20     // Load border pixels into shared memory
21     if (threadIdx.x == 0 && x > 0) {
22         sharedMem[sharedY][0] = input[y * width + (x - 1)];
23     }
24     if (threadIdx.x == blockDim.x - 1 && x < width - 1) {
25         sharedMem[sharedY][sharedX + 1] = input[y * width + (x + 1)];
26     }
27     if (threadIdx.y == 0 && y > 0) {
28         sharedMem[0][sharedX] = input[(y - 1) * width + x];
29     }
30     if (threadIdx.y == blockDim.y - 1 && y < height - 1) {
31         sharedMem[sharedY + 1][sharedX] = input[(y + 1) * width + x];
32     }
33
34     // Load corner pixels
35     if (threadIdx.x == 0 && threadIdx.y == 0 && x > 0 && y > 0) {
36         sharedMem[0][0] = input[(y - 1) * width + (x - 1)];
37     }
38     if (threadIdx.x == blockDim.x - 1 && threadIdx.y == 0 && x < width - 1 && y > 0) {
39         sharedMem[0][sharedX + 1] = input[(y - 1) * width + (x + 1)];
40     }
41     if (threadIdx.x == 0 && threadIdx.y == blockDim.y - 1 && x > 0 && y < height - 1) {
42         sharedMem[sharedY + 1][0] = input[(y + 1) * width + (x - 1)];
43     }
44     if (threadIdx.x == blockDim.x - 1 && threadIdx.y == blockDim.y - 1 && x < width - 1 && y < height - 1) {
45         sharedMem[sharedY + 1][sharedX + 1] = input[(y + 1) * width + (x + 1)];
46     }
47
48     __syncthreads();
49 }
```

Performance Findings

- Optimized GPU implementation unexpectedly slower than naive GPU
- Main reasons:
 - Compute-bound workload
 - Shared memory overhead
 - Synchronization overhead (`__syncthreads()` latency)
 - Modern GPU efficient global memory handling

Recommendations

- **Reduce shared memory overhead:**
 - Minimize unnecessary data in shared memory
 - Adjust shared memory tile sizes
- **Decrease synchronization cost:**
 - Implement warp-level synchronization (`__shfl_sync()`)
- **Hybrid approach:**
 - Naive kernel for compute-bound tasks
 - Optimized kernel for memory-bound tasks

Visual Outputs Comparison



Fig: input_image.jpg



Fig: cpu_output_image.jpg

Visual Outputs Comparison(contd.)



Fig: gpu_output_image_naive.jpg



Fig: gpu_output_image_optimized.jpg

Thread Coarsening

After updating for Thread coarsening (4x coarsening factor) the Optimized almost caught up to the naive (rather than twice as slow, is about 5% slower). However, the output is slightly different



Final Summary and Takeaways

- **Key Insight:** Optimization must align with workload nature — compute-bound cases favor simpler designs
- **Surprising Result:** Naive GPU kernel outperformed optimized version due to lower overhead
- **Improvement Noted:** Thread coarsening helped optimized kernel almost match naive performance
- **Conclusion:** Always profile and analyze before optimizing — simple designs can outperform complex ones on modern GPUs
- **Future Work:** Test across GPU architectures, refine warp-level coordination, and explore dynamic kernel strategies

Thank You!