# GPU Programming Project Report

Aniket Gupta

December 13, 2023

## Introduction

- The project involves implementing morphological operations and image processing like blurring and sharpening on images, using both CPU and GPU parallelization techniques.

- Focus on optimizing the performance using CUDA and shared memory.

## Project Structure

1. **Objective:** Implement erosion, dilation, opening, closing, blurring and sharpening operations on images using GPU programming and compare their performance with CPU implementations.
2. **Implementation Details:**
    - **Programming Language:** CUDA
3. **Image Processing Operations:**
    - **Erosion:** Removal of pixels based on a structuring element.
    - **Dilation:** Addition of pixels based on a structuring element.
    - **Opening:** Erosion followed by Dilation.
    - **Closing:** Dilation followed by Erosion.
    - **Blurring:** Smoothing or averaging of pixel values.
    - **Sharpening:** Enhancing edges or fine details in an image.
4. **GPU Optimization Techniques**
5. **Experimental Setup:**
    - Image dimensions: $10 \times 10$ to $1024 \times 1024$
    - Structuring element size: $3 \times 3$

## Code Structure

```
// Function prototypes
void generateImg(int width, int height, int *img);
void drawSolidRectangle(int width, int height, int *img);
void generateClc(int *img, int *clc, int size_of_filter, int width, int height);
void erosionCPU(int *res_cpu, int *clc, int width, int height, int size_of_filter);
__global__ void erosionImg(int *res, int *clc, int width, int height, int size_of_filter);
__global__ void erosionImgShared(int *res, int *clc, int width, int height, int size_of_filte
void dilationCPU(int *res_cpu, int *clc, int width, int height, int size_of_filter);
__global__ void dilationImg(int *res, int *clc, int width, int height, int size_of_filter);
__global__ void dilationImgShared(int *res, int *clc, int width, int height, int size_of_filt
void blurCPU(int *res_cpu, int *img, int witdh, int height, int size_of_filter)
__global__ void blurImg(int *res, int *img, int witdh, int height, int size_of_filter)
void sharpenCPU(int *res_cpu, int *img, int witdh, int height, float weight_center, float we
__global__ void sharpenImg(int *res, int *img, int witdh, int height, float weight_center, fl
__global__ void sharpenImgOptimized(int *res, int *img, int width, int height, float weight_c
}
```

## Image Generation

```
void generateImg(int width, int height, int *img);
void drawSolidRectangle(int width, int height, int *img);
```

# CPU Implementation - Erosion

```
void erosionCPU(int *res_cpu, int *clc, int width, int height, int size_of_filter);
```

## CPU Implementation - Dilation

```
void dilationCPU( int *res_cpu, int *clc, int width, int height, int size_of_filter);
```

## GPU Implementation - Erosion

```
__global__ void erosionImg(int *res, int *clc, int width, int height, int size_of_filter);
```

# GPU Implementation - Dilation

```
__global__ void dilationImg(int *res, int *clc, int width, int height, int size_of_filter);
```

## Shared Memory Optimization

```
__global__ void erosionImgShared(int *res, int *clc, int width, int height, int size_of_filte
__global__ void dilationImgShared(int *res, int *clc, int width, int height, int size_of_filt
```

## Optimizing Blur Code - GPU

- The provided code implements a blur filter for GPU.
- Techniques used to optimize the GPU implementation:
- Techniques:
  - Shared Memory Usage
  - Warp Synchronization
  - Coalesced Memory Access

## Optimizing Blur Code - GPU

### Listing 1: Original GPU Blur Code

```
__global__ void sharpenImg(int *res, int *img, int witdh, int height, float weight_ce
int index = 1;  // Assuming a 3x3 sharpening filter
int row = blockIdx.x * blockDim.x + threadIdx.x + index;
int col = blockIdx.y * blockDim.y + threadIdx.y + index;

if (row < height - index && col < witdh - index) {
    int center = img[col + row * witdh];
    int neighbors = 0;

    // Calculate the weighted sum of neighboring pixels
    for (int i = -index; i <= index; i++) {
        for (int j = -index; j <= index; j++) {
            if (i == 0 && j == 0) {
                neighbors += weight_center * img[col + j + (row + i) * witdh];
            } else {
                neighbors += weight_neighbors * img[col + j + (row + i) * witdh];
            }
        }
    }

    // Subtract the weighted sum from the center pixel value
    res[(col - index) + (row - index) * witdh] = center - neighbors;
}
}
```

## Optimizing Blur Code - GPU

### Listing 2: Optimized GPU Blur Code

```
__global__ void blurImgOptimized(int *res, int *img, int width, int height, int size
int index = (size_of_filter - 1) / 2;
int row = blockIdx.x * blockDim.x + threadIdx.x + index;
int col = blockIdx.y * blockDim.y + threadIdx.y + index;

// Define shared memory for caching the image region
__shared__ int img_shared[BLOCK_SIZE][BLOCK_SIZE];

// Load data to shared memory with coalesced memory access
int local_row = threadIdx.x;
int local_col = threadIdx.y;
int global_row = row + local_row - index;
int global_col = col + local_col - index;

if (global_row >= 0 && global_row < height && global_col >= 0 && global_col < width) {
    img_shared[local_row][local_col] = img[global_row * width + global_col];
} else {
    img_shared[local_row][local_col] = 0;  // Padding with zeros for out-of-bounds acces
}

__syncthreads();  // Ensure all threads have loaded their data

int sum = 0;

// Apply blur filter with shared memory usage
for (int i = 0; i < size_of_filter; i++) {
```

## Optimizing Sharpening Filter Code

- GPU code benefits from several optimizations.
- Applied shared memory and optimize memory access patterns for improved GPU performance.

### Listing 3: Optimized GPU Sharpening Code

```
__global__ void sharpenImgOptimized(int *res, int *img, int width, int height, float wei
int index = 1; // Assuming a 3x3 sharpening filter
int row = blockIdx.x * blockDim.x + threadIdx.x + index;
int col = blockIdx.y * blockDim.y + threadIdx.y + index;
//BLOCK_SIZE + 2 * index = 18
__shared__ float sharedImg[18][18];

int shared_row = threadIdx.x + index;
int shared_col = threadIdx.y + index;

if (row < height && col < width) {
    sharedImg[shared_row][shared_col] = img[row * width + col];

    // Load additional pixels to shared memory for border handling
    if (threadIdx.x < index) {
        sharedImg[threadIdx.x][shared_col] = img[row * width + col - index];
        sharedImg[threadIdx.x + blockDim.x + index][shared_col] = img[row * width + col -
    }

    if (threadIdx.y < index) {
```

## Performance Measurement

- Performance is measured using CUDA events and CPU clock cycles.
- Execution times for GPU and CPU implementations are compared for both erosion and dilation operations.

## Constant Memory for Constants

Moved BLOCK_SIZE to constant memory, as it is a constant value used in the kernel.

Listing 4: Constant Memory Declaration

```
__constant__ int BLOCK_SIZE = 16;
```

## Optimize Memory Access

Reorganized data access to optimize memory coalescing. This improved memory access patterns and overall performance.

## Optimize Thread and Block Configuration

Experimented with different thread and block configurations to find the optimal combination for the specific GPU architecture.

## Bank Conflicts in Shared Memory

Ensured that shared memory accesses do not lead to bank conflicts. Bank conflicts can degrade performance in shared memory operations.

## Use Shared Memory Efficiently

Utilized shared memory for caching data that is frequently accessed by threads in a block. This reduces the need to access global memory.

## Loop Unrolling

Depending on the situation, loop unrolling in device code used to provide
performance gains.

## Minimized Divergent Code Paths

Minimized divergent code paths in the kernels. Divergence can lead to the
serialization of threads, reducing performance.

## Asynchronous Memory Copies

Used asynchronous memory copies (`cudaMemcpyAsync`) to overlap data transfers with computation.

## Occupancy Considerations

Considered the occupancy of your GPU kernels. Ensured that enough
active warps are present to hide memory latency.

## Results and Discussion

- Output includes execution times for each operation and implementation.
- Comprehensive analysis of performance gains achieved through GPU parallelization and shared memory optimization.

## Execution Times - Part 1

| Operation | M = 10, N = 10 (ms) | M = 200, N = 200 (ms) |
|-----------|---------------------|------------------------|
| Erosion GPU | 0.081920 | 0.070656 |
| Erosion CPU | 0.000000 | 3.000000 |
| Dilation GPU | 0.082016 | 0.062496 |
| Dilation CPU | 0.000000 | 3.000000 |
| Erosion GPU (SM) | 0.038752 | 0.062464 |
| Dilation GPU (SM) | 0.168800 | 0.253952 |
| Opening CPU | 0.000000 | 4.000000 |
| Opening GPU | 0.092160 | 0.107520 |
| Opening GPU (SM) | 0.039936 | 0.089088 |
| Closing CPU | 1.000000 | 3.000000 |
| Closing GPU | 0.097280 | 0.151648 |
| Closing GPU (SM) | 0.039936 | 0.090016 |

Table: Execution times for different operations (M = 10, N = 10 and M = 200, N = 200)

## Execution Times - Part 2

| Operation | M = 400, N = 400 (ms) | M = 1024, N = 1024 (ms) |
|-----------|----------------------|-------------------------|
| Erosion GPU | 0.165888 | 0.539648 |
| Erosion CPU | 11.000000 | 68.000000 |
| Dilation GPU | 0.180224 | 0.519328 |
| Dilation CPU | 12.000000 | 68.000000 |
| Erosion GPU (SM) | 0.153600 | 0.838656 |
| Dilation GPU (SM) | 0.502784 | 1.841152 |
| Opening CPU | 13.000000 | 88.000000 |
| Opening GPU | 0.230400 | 0.983040 |
| Opening GPU (SM) | 0.278528 | 1.614848 |
| Closing CPU | 13.000000 | 89.000000 |
| Closing GPU | 0.231424 | 0.997408 |
| Closing GPU (SM) | 0.276640 | 1.616896 |

Table: Execution times for different operations (M = 400, N = 400 and M = 1024, N = 1024)

## Conclusion

- Successful implementation of image processing operations on both CPU and GPU platforms.
- Shared memory in GPU implementations enhances performance.
- GPU-based operations are competitive with or faster than CPU-based operations.

## Acknowledgments

- Developed as part of the GPU Programming course.
- Gratitude to course instructor Professor Cheolhang An, Youwei Liang and the GPU programming community for insights and guidance.

## References

1. NVIDIA CUDA Toolkit Documentation:
   https://docs.nvidia.com/cuda/

2. Harris, M., & Buck, I. (2007). Programming Massively Parallel
   Processors: A Hands-on Approach.

3. Matloff, N. (2018). Parallel Computing for Data Science: With
   Examples in R, C++, and CUDA.

4. CUDA C Programming Guide:
   https://docs.nvidia.com/cuda/cuda-c-programming-guide/