

# COMS4040A & COMS7045A Assignment 3 – Report

Abdulkadir Dere - 752817 - Computer Science Hons

29 May 2020

## 1 Introduction

In this report, we will focus on the different design methodologies used to create 2D convolution on provided images. We will look at serial implementation, CUDA implementation using both global memory and constant memory, CUDA implementation using both shared memory and constant memory and CUDA implementation using texture memory. Each design methodology is defined and the process has been explained with the code.

## 2 Methodology

We will be using averaging, edge detection and sharpening masks for the convolution process.

-1	-1	-1
-1	9	-1
-1	-1	-1

-1	0	1
-2	0	2
-1	0	1

1	1	1
1	1	1
1	1	1

Table 1: Masks used for convolution: sharpening, edge detection and averaging respectively

### 2.1 Serial Computation

We will first implement the serial version of 2D image convolution algorithm using the CPU. Serial version of the algorithm will follow the following design methodology:

- Read the image using built-in functions. 2D image is returned as 1D image and assigned to a variable.
- Convert the read image from 1D to 2D for easier computation of convolution algorithm since the image is in 2D.

- Pad the image so mask we can compute the corners of the image. The image will be padded with zero padding.
- Apply the convolution mask/filter to the padded image. The convolution algorithm will be executed as a recursive algorithm. The process will start by selecting the first element of the array excluding the padded pixels. We will subtract padded size from the offset (half size of the mask) to find the starting row and column of the image. This pixel must be the same as the original pixel. We will apply the convolution mask to this pixel. (Algorithm 1)
- The convolution mask application process starts by identifying the neighbours of the given pixel location. We will compute neighbours according to the mask size since the image will be padded by the factor of half of the mask on all four sides. We will not run out-of-bounds when retrieving the pixel values for each neighbour since the image is padded. We will create a convolution results array with the dimensions of the mask size to compute the results. We will do element-wise multiplication between neighbours and the mask. All of these values will be added together to compute the convolution result of the given pixel location. (Algorithm 2)
- The convolution result for the pixel location will be saved in the output image and process will continue with the next pixel location until the algorithm iterates through all the pixels. (Algorithm 1)
- The padded area will be removed by un-padding the padded image to retrieve the results in original image dimensions.
- The resulting image will be converted from 2D array to 1D to save the resulting convolution image.

The code for the detailed methodology has been provided with comments. Pre-processing and post-processing functions are not included in the report. These functions can be found in the following file: *serialConvolution.cu* under *src* folder.

```
1 // 2D serial convolution method
2 double **serial_convolution(double **input, double **output){
3     int range = padded_size - offset;
4     // printf("range: %d \n", range);
5
6     for (int i = offset; i<range; i++){
7         for (int j = offset; j<range; j++){
8             output[i][j] = applyMask(input, i, j);
9         }
10    }
11    return output;
12 }
```

Listing 1: Add zero padding to the image

```
1 double applyMask(double **array, int row, int col){
2     int n_size = offset * 2 + 1;
```

```
4 // neighbours of given location
5 double **neighbours = allocateMatrix(n_size, n_size);

7 // dynamically get the neighbours range
8 int n1 = 0;
9 for (int r=row - 1; r <= row + offset; r++){
10     int n2 = 0;
11     for (int c=col - 1; c <= col + offset; c++){
12         neighbours[n1][n2] = array[r][c];
13         n2++;
14     }
15     n1++;
16 }

18 double **convolution = allocateMatrix(n_size, n_size);
19 double value = 0;
20 for (int r=0; r<3; r++){
21     for(int c=0; c<3; c++){
22         convolution[r][c] = mask[r][c] * neighbours[r][c];
23         value = value + convolution[r][c];
24     }
25 }
26 return value;
27 }
```

Listing 2: Apply convolution mask to the given pixel

## 2.2 CUDA implementation using both global memory and constant memory

Image convolution using CUDA C has been implemented using both global memory and constant memory. The convolution mask is constant throughout the convolution process. It is beneficial to cache the convolution mask in the constant memory as this informs the CUDA runtime that mask values will not change during kernel execution (Kirk and Hwu [2010]) (Algorithm 3). The constant memory will set the mask value as read-only and it will be broadcasted to all elements in the convolution kernel execution. The convolution process will be done in global memory using a CUDA kernel (Algorithm 4). The kernel parameters consist of original image, allocated space for resulting image, width and height of the original image. The row and column indexes will be computed to identify which index will be computed by which thread. The starting index for row and column will be calculated by subtracting the offset value because we want to ignore the padded area. This will let the convolution process to start and end at dimensions of the original image so

we don't run out-of-bounds. The kernel will calculate each convolution value by calculating all the elements within the mask filter size (dimension) and add all of them to get the convolution value for a specific pixel. The row and column indexes will be verified so that we are within the dimensions. The resulting value will be saved to the resulting image space.

These functions can be found in the following file: *globalConvolution.cu* under *src* folder.

```
1 // Convolution Mask Dimension
2 #define MASK_DIM 3
3 #define OFFSET (MASK_DIM/2)
4
5 // allocate mask in constant memory
6 __constant__ float d_mask_global[MASK_DIM * MASK_DIM];
```

Listing 3: Cache mask in to the constant memory

```
1 // 2D convolution using global and constant memory
2 __global__ void global_convolution(float *d_Data, float *d_result, int
   width, int height) {
3     // calculate the row and column index to compute for each thread
4     int row = blockIdx.y * blockDim.y + threadIdx.y;
5     int col = blockIdx.x * blockDim.x + threadIdx.x;
6
7     // Starting index for convolution so we can ignore the padded area
8     int i_row = row - OFFSET;
9     int i_col = col - OFFSET;
10
11    // convolution value to be calculated for each pixel's row and
       column
12    double value = 0;
13    // iterate over all rows and column using the mask dimension.
14    // this will calculate all the neighbours and origin pixel and sum
       these values to give
15    // us the value of the origin pixel
16    for (int i = 0; i < MASK_DIM; i++) {
17        for (int j = 0; j < MASK_DIM; j++) {
18            if ((i_row + i) >= 0 && (i_row + i) < height && (i_col + j) >= 0
               && (i_col + j) < width) {
19                // sum all the values within the range of the mask to get
               origin pixel's value
20                value += d_Data[(i_row + i) * width + (i_col + j)] *

```

```

        d_mask_global[i * MASK_DIM + j];
21     }
22 }
23 }
24 // write back convolution result
25 d_result[row * width + col] = value;
26 }

```

Listing 4: 2D Convolution using the global memory

### 2.3 CUDA implementation using both shared memory and constant memory

In this section, we will look at image convolution using both shared memory and constant memory. The constant memory is used for caching the masks for read-only purposes as it is required by every element during the convolution process (5). The kernel parameters consist of the original image, the resulting image, width and height of the original image. Block sized tiles are created for shared memory tiling process. Each tile's row and column indexes are computed for the utilisation by threads. We need the boundaries of the tile so we don't run out-of-bounds and each thread works only on its own tile and does not modify other tiles values. Tiles are loaded from global memory to shared memory for faster access to tile data. Thread barrier has been set so all the threads need to finish loading data before the kernel continues with the process. We will iterate through the rows and columns within the mask dimensions and also check if we are still within the tile bounds. The convolution masks will be applied to every element within the tile within the dimensions of the mask. Hence results can be summed to retrieve the output result for a given pixel.

These functions can be found in the following file: *sharedConvolution.cu* under *src* folder.

```

1 // Convolution Mask Dimension
2 #define MASK_DIM 3
3 #define OFFSET (MASK_DIM/2)
4
5 #define TILE_WIDTH 16
6 #define RADIUS 2
7 #define BLOCK_WIDTH (TILE_WIDTH+(2*RADIUS))
8
9 #define DIAMETER (RADIUS*2+1) // mask diameter
10 #define SIZE (RADIUS*DIAMETER) // mask size
11
12 // allocate mask in constant memory
13 __constant__ float d_mask_shared[MASK_DIM * MASK_DIM];

```

Listing 5: Cache mask in to the constant memory

```
1  __global__ void shared_convolution(float* dData, float* dResult,
   unsigned int width, unsigned int height){

3  // create tile in shared memory for the convolution
4  __shared__ float shared[BLOCK_WIDTH * BLOCK_WIDTH];

6  // for simplicity to use threadIdx
7  int tx = threadIdx.x;
8  int ty = threadIdx.y;
9  int bx = blockIdx.x;
10 int by = blockIdx.y;

12 // get row and column index of pixels in the tile
13 int col = bx * TILE_WIDTH + tx - RADIUS;
14 int row = by * TILE_WIDTH + ty - RADIUS;

16 // Find the last and first pixel locations within the image
17 col = max(0, col);
18 col = min(col, width-1);
19 row = max(row, 0);
20 row = min(row, height-1);

22 // load the tile pixels from the global memory into shared memory
23 // this will help us to reduce global memory access by the factor
   of 1/TILE_WIDTH
24 // ignore any pixels which are out-of-bounds (i.e. padded area)
25 unsigned int index = row * width + col;
26 unsigned int block_index = ty * blockDim.y + tx;
27 shared[block_index] = dData[index];

29 // thread barrier to wait for all the threads to finish loading
   from
30 // global memory to shared memory
31 __syncthreads();

33 // Elementwise multiplication of pixel and mask values and add all
   of the values within the mask
34 // range to get output value of one pixel. Verify that we are not
   working out-of-bounds of the image
35 // We will iterate over rows and columns within the mask
   dimensions (i.e. all the neighbours)
```

```

36     float value = 0;
37     if (((tx >= RADIUS) && (tx < BLOCK_WIDTH-RADIUS)) && ((ty >= RADIUS)
        && (ty <= BLOCK_WIDTH-RADIUS))) {
38         for(int i = 0; i < MASK_DIM; i++) {
39             for(int j = 0; j < MASK_DIM; j++) {
40                 value += shared[block_index + (i * blockDim.x) + j] *
                    d_mask_shared[i * 3 + j];
41             }
42         }
43         dResult[index] = value;
44     }
45 }

```

Listing 6: 2D Convolution using the shared memory

The dimensions of the block are set to block width (7). Block dimension is tile width summed with two times the offset. The offset is the padded area of the image. The offset is multiplied by two to cater for the padding of opposite sides of the image (e.g. left and right). Hence, the block dimension is set to block width so all the blocks are padded and pixels at the block corners can be computed. The padded area of each block is ignored in the kernel when are calculating the starting row and columns. The grid dimensions are set to cater for the minimum number of grids required to compute the image convolution for the whole image.

```

1 #define TILE_WIDTH 16
2 #define OFFSET 2
3 #define BLOCK_WIDTH (TILE_WIDTH + (2 * OFFSET))
4 dim3 dimGrid(BLOCKS, BLOCKS);
5 dim3 dimBlock(BLOCK_WIDTH, BLOCK_WIDTH);

```

Listing 7: Shared memory kernel dimensions

## 3 Experiment

### 3.1 Experiment Setup

Experiments are conducted on a cluster. The details for the CUDA device are listed.

CUDA Device 0

Major revision number: 6

Minor revision number: 1

Name: GeForce GTX 1060 6GB

Total global memory: 6371475456  
Total shared memory per block: 49152  
Total registers per block: 65536  
Warp size: 32  
Maximum memory pitch: 2147483647  
Maximum threads per block: 1024  
Maximum dimension 0 of block: 1024  
Maximum dimension 1 of block: 1024  
Maximum dimension 2 of block: 64  
Maximum dimension 0 of grid: 2147483647  
Maximum dimension 1 of grid: 65535  
Maximum dimension 2 of grid: 65535  
Clock rate: 1784500  
Total constant memory: 65536  
Texture alignment: 512  
Concurrent copy and execution: Yes  
Number of multiprocessors: 10  
Kernel execution timeout: Yes

### 3.2 Experiment Results

Experimental results are shown for the serial, global memory and shared memory implementations for the image convolution.



Image	image21.pgm	lena_bw.pgm
Matrix Size	512x512	512x512
Tile Size	16x16	16x16
Serial Convolution Time (ms)	154.862244	152.380386
Global Memory Time (ms)	0.093184	0.094880
Shared Memory Time (ms)	0.030656	0.032704
Speedup of global memory kernel (ms)	1661.9	1606.03
Speedup of shared memory kernel (ms)	5051.61	4659.38
Throughput of serial implementation (GFLOPS)	0.0152348	0.0154829
Throughput of global memory implementation (GFLOPS)	25.3187	24.8661
Throughput of shared memory implementation (GFLOPS)	76.9603	72.1409
Performance improvement: global over serial	1661.9x	1606.03x
Performance improvement: shared over serial	5051.61x	4659.38x
Performance improvement: shared over global	3.03967x	2.90117x

Table 2: Results of the convolutions applied to the given images.

Experimental results are shown for the serial, global memory and shared memory implementations for the image convolution. The results for the different implementations are the same. Hence, the different image results are produced by the different filters used during the convolution process (figure 1). The results for the *lena\_bw* image is more clear (figure 2) and the effects can be visualised better. The different design implementations do not affect the resulting image, as expected.

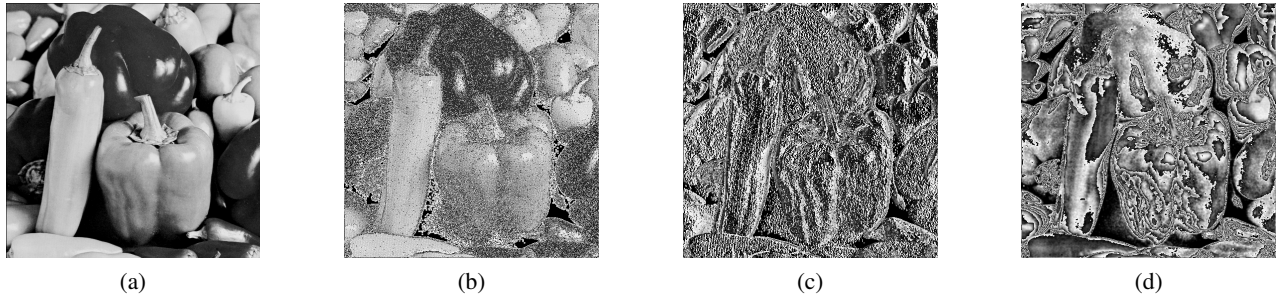


Figure 1: Convolution results for the *image21* image (a) Original Image (b) Sharpening mask (c) Edge detection mask (d) Averaging mask

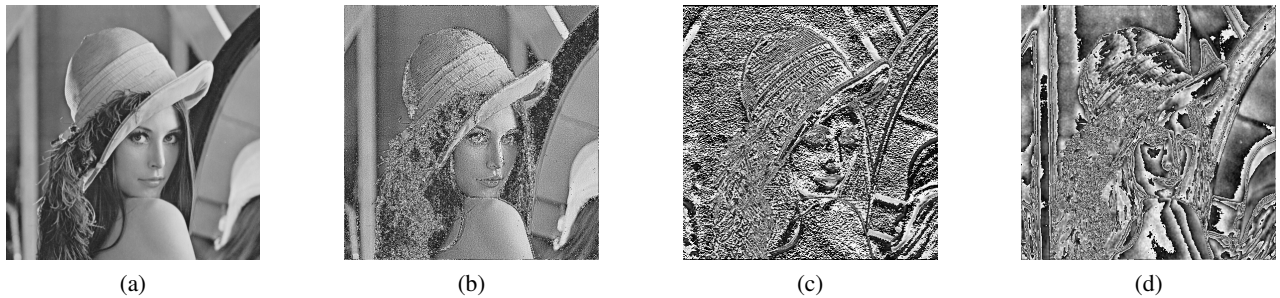


Figure 2: Convolution results for the *lena\_bw* image (a) Original Image (b) Sharpening mask (c) Edge detection mask (d) Averaging mask

The resulting images for all the kernels can be found in the *results* folder.

The averaging filter with 16 elements (4x4 size) results are as follows:

Serial Image Convolution Time: 156.733444 ms

Global Memory Time elapsed: 0.153600 ms

Shared Memory Time elapsed: 0.034880 ms

We can see that the global memory version takes almost double the kernel with 3x3 size as shown in table 2. However, the shared memory version still performs similar results compared to the smaller sized averaging filter.

### 3.3 Summary of Results

As we can see from table 2 that global memory version has greatly improved on the serial implementation. The shared memory version has improved on global memory and serial implementations. However, global memory implementation is still not implemented efficiently since every thread calculates its own output. This results in global memory access for each thread. The global memory version can be improved by implementing threads to compute a tile of the image similar to shared memory implementation. This will decrease access to global memory.

## 4 Questions and Answers

(a) How many floating point operations are performed in your convolution kernel using global memory? The floating point operations are calculated by multiplying image dimensions with mask dimensions. Hence,  $FLOP = width * height * mask\_width * mask\_height$ . All the provided images are the size of 512 x 512 and mask size is kept at 3x3 in most cases. Therefore,  $FLOP = 512 * 512 * 3 * 3 = 2359296$ .

(b) How many global memory reads are performed by your kernel using global memory and kernel using shared memory, respectively? The global memory version performs a global memory read for each thread since every thread calculates its own output. This is inefficient as discussed before. The shared memory version has less memory reads since we load the tiles from global memory to shared memory. The shared memory reads are  $\frac{1}{TILE\_WIDTH}$  of the global memory kernel because we are accessing the global memory only to load the tile from global memory to shared memory. Hence, the memory access decreases by the factor of  $\frac{1}{TILE\_WIDTH}$  of the global memory.

(c) How many global memory writes are performed by your convolution kernel using shared memory? The global memory writes are conducted by each thread. Hence, the global memory write will be  $512 \times 512 = 262144$  writes to the global memory for an image with the dimensions of 512x512. The threads are restricted in a way that only the threads within the limits of the image can write to the global memory. Hence, we are calculating the number of threads which has write access to the output image.

(d) What would happen to the performance of your kernel using shared and constant memory when the size of the 'averaging' mask increases (say, to a substantial large size)? The issue is that mask is cached in the constant memory. So any mask size which exceeds the limits of cache memory capacity will not allow the program to compile and run. Hence, it will not work. Another is the kernel dimensions are dependent on the tile width which is also affected by the mask size. Hence, any changes to mask size will break the limitations of the kernel. Hence, CUDA will not be able to compile and run the application.

## 5 Conclusion

In this report, we have discussed and showcased different image convolution methods using various device memories. The usage of shared memory shows significant performance increases compared to global memory and serial implementations. The usage of constant memory allowed us to access convolution masks as read-only data and sped-up the process. These can be visualised in the given tables and results.

## References

Kirk, D. B. and Hwu, W.-m. W. (2010). *Programming Massively Parallel Processors: A Hands-on Approach*. Morgan Kaufmann Publishers Inc., 2nd edition.