

EE4C07

## Advanced Computing Systems

### Lab 2: GPU-based Computing report

Yingli Ni                      4515021  
Javier James                4170253

Ye Wang                      4473582  
Luis Garcia Rosario        4062949

Friday 9<sup>th</sup> October, 2015

# Abstraction

This is the report for the lab of the course Advanced Computing Systems at the TU Delft. During this course...

# Contents

<b>Abstraction</b>	<b>i</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Background</b>	<b>2</b>
2.1 Hardware analysis . . . . .	2
2.2 Software analysis . . . . .	2
2.2.1 RGB2GRAY algorithm . . . . .	2
2.2.2 Histogram algorithm . . . . .	3
2.2.3 Contrast algorithm . . . . .	4
2.2.4 Smooth algorithm . . . . .	4
<b>3 Converting a color image to grayscale</b>	<b>5</b>
3.1 Implementation strategy . . . . .	5
3.2 Kernel speed-up achieved . . . . .	5
3.3 Validation of results . . . . .	6
3.4 Optimizations technique used . . . . .	6
3.5 Compiler options . . . . .	6
<b>4 Histogram computation</b>	<b>8</b>
4.1 Implementation strategy . . . . .	8
4.2 Kernel speed-up achieved . . . . .	8
4.3 Validation of result . . . . .	9
4.4 Problems . . . . .	9
4.5 Solutions . . . . .	9
<b>5 Contrast enhancement</b>	<b>10</b>
5.1 Implementation strategy . . . . .	10
5.2 Kernel speed-up achieved . . . . .	10
5.3 Validation of result . . . . .	10
5.4 Challenges encountered . . . . .	11
<b>6 Image smoothing</b>	<b>12</b>
6.1 Implementation Strategy . . . . .	12
6.2 Kernel speed-up achieved . . . . .	12
6.3 Validation of result . . . . .	13

<i>CONTENTS</i>	iii
<b>7 Final Results</b>	<b>14</b>
7.1 Software validation . . . . .	14
7.2 Validation of result . . . . .	14
<b>A Code for Converting a color image to grayscale</b>	<b>15</b>
<b>B Code for Histogram computation</b>	<b>17</b>
<b>C Code for Contrast enhancement</b>	<b>20</b>
<b>D Code for Image Smoothing</b>	<b>23</b>

# Chapter 1

## Introduction

Improving the speed-up of a software application may sometimes be essential, especially in cases where there are computational intensive algorithms. With all the millions of transistors available on a single chip, it is possible to have more cores running in parallel. Engineering such software is very difficult. During the Lab of Advanced Computing Systems we of Group 07 have analyzed and experimented in depth the speed-up achievable and the performance in a provided Software application running on the **NVIDIA Geforce GTX750Ti**. This Software application uses this NVIDIA GPU device to execute function kernels in parallel.

**Method and Experiment** During the experiments three images were tested individually, image04.bmp, image09.bmp and image15.jpg, having a size of 5.0 MB, 67.1MB, 5.0MB respectively. Each image were given as input to the cpu functions, and then to the gpu functions after. In this way we could analyze the effect of processing the entire image on the cpu vs gpu. The output result image smooth.bmp was then compared to see the difference. The % of error was denoted as the number of pixels above a set threshold, where both outputs differ from one another.

**Report Organization** The organization of this practical report is as following. Chapter 2 sets the mood by giving background information and analysis of the hardware and software. The potential speed-up and requirements will be discussed in short. Chapter ?? until 6 continues by going in depth analyzing the functions rgb2gray, histogram, contrast, and smoothing in chronological order. In each chapter we will give the speed-up attained for that kernel, the challenges encountered, and the validations of the results. Subsequently, chapter 7 provides the final results testing the entire application on the GPU and the achieved speed-up and image error. Finally, chapter ?? concludes this report by providing a summary of the most important points of the report and providing a discussion.

## Chapter 2

# Background

This chapter focuses on the hardware and software analysis. Section 2.1 begins analysing the target GPU hardware providing its specification. Continually, section 2.2 analysis each software function to see its potential speed-up.

### 2.1 Hardware analysis

The target hardware is the NVIDIA Geforce GTX750Ti/Tesla C2075 GPU, which is present on the similde machine. This GPU uses first-generation NVIDIA Maxwell architecture. and the performance up to twice of previous generation graphics cards, however, it is only half the power consumption of previous generation products.the GTX 750TiGPU graphics card consits of 640 CUDA processor core and 1020MHz base frequency as well as 1085MHz enhance the frequency.

### 2.2 Software analysis

In this section we briefly describe the functionality of each function, display its potential speed-up as given by Amdahl's law in equation 2.2.1. As the number of cores increase to infinity the speed-up is upper bounded by the sequential part, which is given by equation 2.2.2.

$$Speedup = \frac{1}{1 - f + \frac{f}{n}} \quad (2.2.1)$$

$$\leq \frac{1}{1 - f} \quad (2.2.2)$$

#### 2.2.1 RGB2GRAY algorithm

**Function Description** This is the original code that can converst the image to gray. According to the code , we can get cuda code and kernel function.

We see the pixel as the matrix which define y and x as the height and width of the image.we assume different color has different pixels function and float type in the same place.So, we can just use function and add these float index together. The most important thing is that we need to static transform the type of the variable because the input image type is char when the r,g,b is float. In the end , we need to change float to char.

```

1 for ( int y = 0; y < height; y++ )
2 { for ( int x = 0; x < width; x++ )
3     { float grayPix = 0.0f;
4         float r = static_cast< float >(inputImage[(y * width)
5             + x]);
6         float g = static_cast< float >(inputImage[(width *
7             height) + (y * width) + x]);
8         float b = static_cast< float >(inputImage[(2 * width *
9             height) + (y * width) + x]);
10        grayPix = (0.3f * r) + (0.59f * g) + (0.11f * b);
11        grayImage[(y * width) + x] = static_cast< unsigned
12            char >(grayPix);
13    }
14 }

```

**Listing 2.1:** ... Yingli put a caption namehere

This code has three parts. First is the for loop. Pointer moves through the image height. Second is the for loop that moves through the image width, the third part is calculating function and adding the variable. The whole function is based on the loop function. This is the basic theory of the RGB to gray image. We need to use the main function to program kernel function and recall the kernel function at the cuda function.

**Max speed-up** The for loop which a kernel was created from was 98% of the entire function. Therefore, according to equation 2.2.2 there is a maximum speed is bounded by x50. This is shown in 2.2.3.

$$Speedup_{max-rgb2gray} = \frac{1}{1-0.98} = \frac{1}{0.02} = 50 \quad (2.2.3)$$

## 2.2.2 Histogram algorithm

**Function Description** In histogram part, we have to static how many times each grayscale appears in the image. That means, we should have 256 counters for each grayscale and visit every point in the image in order. These counters use a series of continuous addresses, so we could change the value of it easily. But this would take a long time if the size of the image is large. I will show what the function it should have from software aspect by an example in CPU version.

```

1 for ( int y = 0; y < height; y++ )
2 {
3     for ( int x = 0; x < width; x++ )
4     {
5         histogram[static_cast< unsigned int >(grayImage[(y * width) +
6             x])] += 1;
7     }
8 }

```

**Listing 2.2:** Histogram algorithm analysis

This code could be seen as three parts. First is the for loop. Pointer moves through the image. Second is visiting memory which has stored grayscale of each point in the image. Third is change the value of the counters. Array histogram is used as the counter to record how many times each grayscale appears in the image.

Other parts of histogram is to visualize statistical results. There is no need to use parallel operation to optimize it.

**Max speed-up** The for loop which a kernel was created from was 23% of the entire function. Therefore, according to equation 2.2.2 there is a maximum speed is bounded by x50. This is shown in 2.2.4.

$$Speedup_{max-histogram} = \frac{1}{1-0.23} = \frac{1}{0.77} = 1.30 \quad (2.2.4)$$

### 2.2.3 Contrast algorithm

**Function Description** The Contrast routines takes as input a grayImage and histogram of length 256 and a contrast threshold value. A contrast is set on the grayImage base on the histogram and the contrast threshold.

The function begins by finding the index for which the index value is no longer less than the threshold value. This is done two times, from firsts index to the last, then from the last to the first. Respectively, these index will be min and max. Once these index are found, each pixel will be updated.

**Max speed-up** The for loop which a kernel was created from was 46% of the entire function. Therefore, according to equation 2.2.2 there is a maximum speed is bounded by x1.85. This is shown in 2.2.5.

$$Speedup_{max-contrast} = \frac{1}{1-0.46} = \frac{1}{0.54} = 1.85 \quad (2.2.5)$$

### 2.2.4 Smooth algorithm

**Function Description** In phase 4 of the image processing pipeline we can see how a triangular smoothing algorithm is used to remove noise from a gray-scale image. Each pixel of an input image is replaced by a weighted average of its neighboring pixels to remove structures from the image.

**Max speed-up** The for loop which a kernel was created from was 81% of the entire function. Therefore, according to equation 2.2.2 there is a maximum speed is bounded by x5.26. This is shown in 2.2.6.

$$Speedup_{max-smooth} = \frac{1}{1-0.81} = \frac{1}{0.19} = 5.26 \quad (2.2.6)$$



## Chapter 3

# Converting a color image to grayscale

### 3.1 Implementation strategy

In RGB2GRAY part, we have to use the function to convert the RGB image to gray image. Firstly, we can use thread to program code. We define the index and the size of it. on the other hand, we can see the input image as the matrix and just use thread to define and calculation. The thread is typical processing by one core of the GPU. And the block is Composed by multiple threads (which can be expressed as a one-dimensional, two-dimensional, three-dimensional, and then elaborate on the specific below). Each block is executed in parallel, inter-block can not communicate, there is no execution order. we need also notice that the number of thread blocks is limited to no more than 65,535 because of the hardware limitations. And the grid is composed by multiply blocks. In this part of code we use thread to define the index. Here is the code.

```
1 unsigned int index = blockIdx.x * blockDim.x + threadIdx.x;  
2 if (index < ImageSize)
```

This part code is defined and programmed in the kernel function. it will be easy to calculate because we don't defined two or more thread and the calculate function is also linear calculation. the speed of the RGB is depended on the different size of the image.

### 3.2 Kernel speed-up achieved

Input image	CPU Time (sec)	GPU time (sec)	Speed-up
image 04	0.013690	0.001081	x12.67
image 09	0.483947	0.007024	x68.90
image 15	0.053946	0.001512	x35.67

**Table 3.2.1:** Kernel speed up using RGB2GRAY on GPU

Slight description of how function works

What was changed to create parallelism

challenges encountered

### 3.3 Validation of results

Input image	Pixels above threshold	error %
image 04	864	0
image 09	0	0
image 15	624	0

**Table 3.3.1:** RGB2GRAY error comparison CPU vs GPU

### 3.4 Optimizations technique used

When we use parallel computing to solve this problem, we want to use threads to finish the computing in O(1). Like following code.

```
1 unsigned int index = blockIdx.x * blockDim.x + threadIdx.x;
```

we can assume that When the use of GPU accelerated image processing, the RGB image format generally 3-channel image, each through one byte, or 24-bit pixel image. we can expand the input image size to 3\* width\*height. While cuda access data elements, if each thread access 8bit, 16bit, 32bit, when 64bit, corresponding to the length of the data segment can be 32Byte, 64Byte, 128Byte, 128Byte. So we can meet the requirements to access global memory consolidation, improving access to global memory access performance. However, If each thread to access 24bit will not meet the requirements to access global memory consolidation, affecting the performance of the global memory access, so we can find solution that we expand the data size when we deal with the 24bit data as 32bit data access. So in the RGB part we assume the block size is 512 in order to suit for different sizes image.

### 3.5 Compiler options

Without any optimization option, the compiler's goal is to reduce the cost of compilation and to make debugging produce the expected results. So, for our program, we need to choice the best options. a slightly different set of optimizations may be enabled at each -O level than those listed here. But for our program we can find the best option in the make file document, we can find the fellow code:

```
1 debug      := 0
2   Ifneq    ($(debug), 1)
3     CFLAGS      := -O3 -g0
4     NVCCFLAGS   += -O3
5 else
6     CFLAGS      := -O0 -g3 -Wall
7     NVCCFLAGS   += -O0
8 End
```

So -O3 is the option that we used during the debug. we can find the definition about -O3 option is Optimizing yet more. And -O3 turns on all optimizations specified by -O2 and also turns on the -finline-

functions, funswitch-loops, -fpredictive-commoning, -fgcse-after-reload, -ftree-loop-vectorize, -ftree-loop-distribute-patterns, -ftree-slp-vectorize, -fvect-cost-model, -ftree-partial-pre and -fipa-cp-clone options and such functions. We can see -O3 use lots of functions and for us this is the best option during the debug. Here is the description of some functions. -ftree-loop-distribute-patterns Perform loop distribution of patterns that can be code generated with calls to a library. This flag is enabled by default at -O3. This pass distributes the initialization loops and generates a call to memset zero. For example, the loop

```

1  DO I = 1, N
2      A(I) = 0
3      B(I) = A(I) + I
4  ENDDO    is transformed to
5      DO I = 1, N
6          A(I) = 0          ENDDO
7      DO I = 1, N
8          B(I) = A(I) + I
9  ENDDO

```

the initialization loop is transformed into a call to memset zero. -ftree-loop-vectorize Perform loop vectorization on trees. This flag is enabled by default at -O3 and when -ftree-vectorize is enabled. -ftree-slp-vectorize Perform basic block vectorization on trees. This flag is enabled by default at -O3 and when -ftree-vectorize is enabled.

Compiling multiple files at once to a single output file mode allows the compiler to use information gained from all of the files when compiling each of them. Conclusion is that we invoke gcc with -O3 as our best option because it use and transfer lots of functions that we can promote our program.

## Chapter 4

# Histogram computation

### 4.1 Implementation strategy

What we use to implement GPU version is atomic operation. In concurrent programming, an is atomic, linearizable, indivisible or interruptible if it appears to the rest of the system to occur instantaneously. Atomicity is a guarantee of isolation from concurrent processes. When a thread uses atomic operation to read or write a memory unit, this memory unit would be locked during this time. Other threads would have no ability to read or write this memory unit. They have to wait until the first thread have finished its operation and unlock the memory unit. There are lots of atomic operations in cuda. We use atomicAdd() in this part of the program. The code shows below.

```
1 unsigned int index = blockIdx.x * blockDim.x + threadIdx.x;
2     if (index < ImageSize)
3     {
4         unsigned char Item = d_grayImage[index];
5         atomicAdd (&(device_histogram[Item]), 1);
6     }
```

**Listing 4.1:** ..

We will shorten the running time because of using atomicAdd() function. The time complexity of serial computing is  $O(n^2/N)$ , n is the size of each image and N is the number of threads. This also waste some computing resource of GPU. No matter how many threads can be used, the speed of histogram is limited by the size of the histogram.

### 4.2 Kernel speed-up achieved

Input image	CPU Time (sec)	GPU time (sec)	Speed-up
image 04	0.003524	0.003260	x1.08
image 09	0.165254	0.078573	x2.10
image 15	0.008705	0.005942	x1.47

**Table 4.2.1:** Kernel speed up using Histogram on GPU

### 4.3 Validation of result

Input image	Pixels above threshold	error %
image 04	0	0
image 09	0	0
image 15	0	0

**Table 4.3.1:** Histogram error comparison CPU vs GPU

### 4.4 Problems

In histogram part, we are asked to count points number of each grayscale. When we use CPU do serial arithmetic, it is very obvious that we can access memory of the grayscale of each point in the image in order and then add the counter. The time complexity of serial computing is  $O(n^2)$ ,  $n$  is the size of each image.

When we use parallel computing to solve this problem, we want to use threads to finish the computing in  $O(1)$ . Like following code.

```
unsigned int index = blockIdx.x * blockDim.x + threadIdx.x;
histogram[index]++;
```

But this is only an illusion. It could not be so easy. The first challenge we met is the memory conflict. To explain this challenge, we have to show the principle of parallel computing at first. In one unit time, every thread would have their own work. They would access the memory, process the data or modify the memory at the same time. In this case, we have  $n^2$  memory units to read and 256 memory units to write. Because these are only 256 kinds of grayscale in enactment, some points may have the same grayscale. If we operate what we write above, it must have a lot of memory conflict. This means more than one threads would write in a memory at same time. We won't get accurate result after these chaotic operations.

### 4.5 Solutions

So what we use to solve this problem is atomic operation. In concurrent programming, an is atomic, linearizable, indivisible or uninterruptible if it appears to the rest of the system to occur instantaneously. Atomicity is a guarantee of isolation from concurrent processes. When a thread uses atomic operation to read or write a memory unit, this memory unit would be locked during this time. Other threads would have no ability to read or write this memory unit. They have to wait until the first thread have finished its operation and unlock the memory unit. There are lots of atomic operations in cuda. We use `atomicAdd()` in this part of the program. The code shows below. `unsigned int index = blockIdx.x * blockDim.x + threadIdx.x;` —`atomicMax()` would be added later— We will shorten the running time because of using `atomicAdd()` function. The time complexity of serial computing is  $O()$ ,  $n$  is the size of each image and  $N$  is the number of threads.

## Chapter 5

# Contrast enhancement

### 5.1 Implementation strategy

The pixel update which is done for each pixel using a for loop of  $O(n^2)$  have been used as the function kernel. Each thread on the GPU will be responsible for updating one pixel. This function was executed with no problem because they were no data dependencies.

### 5.2 Kernel speed-up achieved

The baseline speed for the cpu function was 0.061069 seconds. After we parallelized this for loop we got a speed up of

Input image	CPU Time (sec)	GPU time (sec)	Speed-up
image 04	0.013598	0.001125	x12.09
image 09	0.514815	0.011060	x46.547
image 15	0.063011	0.001814	x34.73

**Table 5.2.1:** Kernel speed up using Contrast on GPU

### 5.3 Validation of result

Input image	Pixels above threshold	error %
image 04	0	0
image 09	0	0
image 15	0	0

**Table 5.3.1:** Contrast error comparison CPU vs GPU

## 5.4 Challenges encountered

This routine was trivial. Calculations could have been done using a 1D array on the kernel. Furthermore, there was no data dependencies, between the indexes.

## Chapter 6

# Image smoothing

### 6.1 Implementation Strategy

In the image smoothing phase the application takes a noisy gray-scale image as an input, remove small-scale unwanted structure, and return the resulting filtered/smoothed image as output. The sequential time complexity of the smoothing algorithm only is  $O(n^4)$  due to the 4 layered for-loop in the function. This is a computational intensive function for the CPU, but then due to its parallel nature (operations performed per pixels) it's ideal for the GPU.

The easiest approach was to map the image's width and height to CUDA's two-dimensional index for threads where and operation on each pixel is approached by at least one thread. We started with a 256 threads (a 16x16 block) thread block. The number of blocks per grid was dynamically determined by the width and the height of the input gray-scaled noisy image. After copying the image from host to device memory with cuda memcpy function, the kernel was invoked with a number of thread blocks equal to  $width/threadIdx.x$  in one dimension and a number of  $height/threadIdx.y$  thread blocks in the other dimension.

### 6.2 Kernel speed-up achieved

The execution times of the tests for the cases with CUDA parallelizations and the cases without any CUDA parallelizations are given in Table 6.2.1. The speed-up achieved was calculated by:  $Speedup = T_{old}/T_{new} = 0.547109/0.009494 \approx 57.6$ .

Input image	CPU Time (sec)	GPU time (sec)	Speed-up
image 04	0.129558	0.002254	x57.48
image 09	6.103695	0.078858	x77.40
image 15	0.547109	0.009494	x57.63

**Table 6.2.1:** Kernel speed up using Smooth on GPU



### 6.3 Validation of result

The % of error was denoted in the table below as the number of pixels above a set threshold. We validated the resulting output image using the 'compare' application that tests the sequential output vs the parallel output.

Input image	Pixels above threshold	error %
image 04	0	0
image 09	0	0
image 15	0	0

**Table 6.3.1:** Smoothing error comparison CPU vs GPU

## Chapter 7

# Final Results

In this chapter we will show the final results of our experiments.

### 7.1 Software validation

To test the software, each cpu (gpu) function output will go to the input of the next cpu (gpu) function in the pipeline. However, the same input image is given to both cpu and gpu processing pipeline. The output image, smooth.bmp, from the cpu and gpu was compared with one another.

Input image	CPU Time (sec)	GPU time (sec)	Speed-up
image 04	0.511147	1.539530	x0.33
image 09	15.718547	29.692304	x0.52
image 15	1.523320	3.626021	x0.42

**Table 7.1.1:** Overall speed up

### 7.2 Validation of result

Input image	Pixels above threshold	error %
image 04	197562	3
image 09	28792137	14
image 15	238785	1

**Table 7.2.1:** all functions error comparison CPU vs GPU

## Appendix A

# Code for Converting a color image to grayscale

```
1 __global__ void rgb2grayCudaKernel(unsigned char *d_inputImage,
2   unsigned char *d_grayImage, int ImageSize)
3 {
4     unsigned index = blockIdx.x * blockDim.x + threadIdx.x;
5     if(index < ImageSize)
6     {
7         float grayPix = 0.0f;
8         float r = static_cast< float >(d_inputImage[
9             index]);
10        float g = static_cast< float >(d_inputImage[
11            ImageSize + index]);
12        float b = static_cast< float >(d_inputImage
13            [(2 * ImageSize) + index]);
14        grayPix = (0.3f * r) + (0.59f * g) + (0.11f *
15            b);
16        d_grayImage[index] = static_cast< unsigned
17            char >(grayPix);
18    }
19 }
20
21 void rgb2grayCuda(unsigned char *inputImage, unsigned char *
22   grayImage, const int width, const int height)
23 {
24     NSTimer kernelTime = NSTimer("kernelTime", false, false);
25     NSTimer timer_setUp_Comm = NSTimer("
26         setUp_Communication_Timnee", false, false);
27     memset(reinterpret_cast< void * >(grayImage), 0, width *
28         height * sizeof(unsigned char));
```

```

23
24     unsigned char* d_grayImage= NULL;
25     unsigned char* d_inputImage= NULL;
26     int ImageSize = width * height;
27
28     timer_setUp_Comm.start();
29     checkCudaCall(cudaMalloc( (void **) &d_inputImage, (3*
30         ImageSize) ));
31     checkCudaCall(cudaMalloc((void **) &d_grayImage, ImageSize)
32         );
33     checkCudaCall(cudaMemcpy(d_grayImage, grayImage, ImageSize,
34         cudaMemcpyHostToDevice));
35     checkCudaCall(cudaMemcpy(d_inputImage, inputImage, 3*
36         ImageSize, cudaMemcpyHostToDevice));
37     timer_setUp_Comm.stop();
38
39     time_setup_comm[0]= timer_setUp_Comm.getElapsed();
40
41     kernelTime.start();
42     dim3 dimBlock(512);
43     dim3 dimGrid(ImageSize/(int)dimBlock.x);
44     rgb2grayCudaKernel<<<dimGrid, dimBlock>>>(d_inputImage,
45         d_grayImage, ImageSize);
46     cudaDeviceSynchronize();
47     kernelTime.stop();
48
49     kernelGpuTime[0]= kernelTime.getElapsed();
50
51     timer_setUp_Comm.start();
52     checkCudaCall(cudaMemcpy(grayImage, d_grayImage, ImageSize,
53         cudaMemcpyDeviceToHost));
54     timer_setUp_Comm.stop();
55
56     time_setup_comm[0] += timer_setUp_Comm.getElapsed();
57
58     checkCudaCall(cudaFree(d_grayImage));
59     checkCudaCall(cudaFree(d_inputImage));
60
61     cout << fixed << setprecision(6);
62     cout << "rgb2gray (gpu): \t\t" << kernelTime.getElapsed()
63         << " seconds." << endl;
64
65 }

```

## Appendix B

# Code for Histogram computation

```
1  __global__ void histogram1DCudaKernel(int ImageSize, unsigned int *
   device_histogram, unsigned char *d_grayImage)
2  {
3      // set the pointer to every element in d_grayImage
4      unsigned int index = blockIdx.x * blockDim.x + threadIdx.x;
5      if (index < ImageSize)
6      {
7          unsigned char Item = d_grayImage[index];
8          //use atomic operation to solve problem of memory
           conflict
9          atomicAdd (&(device_histogram[Item]), 1);
10     }
11 }
12
13
14
15 void histogram1DCuda(unsigned char *grayImage, unsigned char *
   histogramImage, const int width, const int height,
16                     unsigned int *histogram, const
17                     unsigned int HISTOGRAM_SIZE,
18                     const unsigned int BAR_WIDTH)
19 {
20     // set the number of threads in a single block
21     dim3 threadBlockSize(BLOCK_MAX_THREADSIZE);
22     unsigned int max = 0;
23     int ImageSize = width * height;
24     NSTimer kernelTime = NSTimer("kernelTime", false, false);
25     NSTimer timer_setUp_Comm = NSTimer("setUp_Communication_Timeee",
26                                       false, false);
27     timer_setUp_Comm.start();
```

```

28     memset(reinterpret_cast< void * >(histogram), 0,
29             HISTOGRAM_SIZE * sizeof(unsigned int));
30     // copy histogram to device_histogram
31     unsigned int* device_histogram = NULL;
32     checkCudaCall(cudaMalloc((void **) &device_histogram,
33                               HISTOGRAM_SIZE* sizeof(unsigned int)));
34     checkCudaCall(cudaMemcpy(device_histogram, histogram,
35                               HISTOGRAM_SIZE* sizeof(unsigned int),
36                               cudaMemcpyHostToDevice));
37     // copy grayImage to d_grayImage
38     unsigned char* d_grayImage = NULL;
39     checkCudaCall(cudaMalloc((void **) &d_grayImage, ImageSize)
40                     );
41     checkCudaCall(cudaMemcpy(d_grayImage, grayImage, ImageSize,
42                               cudaMemcpyHostToDevice));
43     timer_setUp_Comm.stop();
44
45     time_setup_comm[1]= timer_setUp_Comm.getElapsed();
46
47     kernelTime.start();
48     // set the number of blocks
49     dim3 BlockNum(width*height/threadBlockSize.x+1);
50     histogram1DCudaKernel<<<BlockNum, threadBlockSize>>>(
51         ImageSize, device_histogram, d_grayImage);
52     cudaDeviceSynchronize();
53     kernelTime.stop();
54
55     kernelGpuTime[1]= kernelTime.getElapsed();
56
57     timer_setUp_Comm.start();
58     checkCudaCall(cudaMemcpy(histogram, device_histogram,
59                               HISTOGRAM_SIZE* sizeof(unsigned int),
60                               cudaMemcpyDeviceToHost));
61     timer_setUp_Comm.stop();
62
63     time_setup_comm[1] += timer_setUp_Comm.getElapsed();
64
65     checkCudaCall(cudaFree(device_histogram));
66     // find the largest number in histogram
67     for ( unsigned int i = 0; i < HISTOGRAM_SIZE; i++ )
68     {
69         if ( histogram[i] > max )
70         {
71             max = histogram[i];
72         }
73     }

```

```
66
67 for ( int x = 0; x < HISTOGRAM_SIZE * BAR_WIDTH; x += BAR_WIDTH )
68 {
69     unsigned int value = HISTOGRAM_SIZE - ((histogram[x / BAR_WIDTH
70         ] * HISTOGRAM_SIZE) / max);
71
72     for ( unsigned int y = 0; y < value; y++ )
73     {
74         for ( unsigned int i = 0; i < BAR_WIDTH; i++ )
75         {
76             histogramImage[(y * HISTOGRAM_SIZE * BAR_WIDTH) + x + i] =
77                 0;
78         }
79     }
80     for ( unsigned int y = value; y < HISTOGRAM_SIZE; y++ )
81     {
82         for ( unsigned int i = 0; i < BAR_WIDTH; i++ )
83         {
84             histogramImage[(y * HISTOGRAM_SIZE * BAR_WIDTH) + x + i] =
85                 255;
86         }
87     }
88 }
89
90 cout << fixed << setprecision(6);
91 cout << "histogram1D (gpu): \t\t" << kernelTime.getElapsed() << "
92     seconds." << endl;
```

## Appendix C

### Code for Contrast enhancement

```
1  __global__ void contrast1DKernel(unsigned char *grayImage, const
    int width, const int height, int min, int max, int diff, int
    grayImageSize)
2  {
3
4      unsigned int index = blockIdx.x * blockDim.x + threadIdx.x;
5
6      //ensure we dont use more threads than image size
7      if(index < grayImageSize)
8      {
9
10         unsigned char pixel = grayImage[index];
11
12         if ( pixel < min )
13         {
14             pixel = 0;
15         }
16         else if ( pixel > max )
17         {
18             pixel = 255;
19         }
20         else
21         {
22             pixel = static_cast< unsigned char >(255.0f * (pixel -
                min) / diff);
23         }
24         grayImage[index] = pixel;
25     }
26 }
27
28 void contrast1DCuda(unsigned char *grayImage, const int width,
    const int height,
29     unsigned int *histogram, const unsigned int HISTOGRAM_SIZE,
```



```

30         const unsigned int CONTRAST_THRESHOLD)
31 {
32
33     unsigned int i = 0;
34     NSTimer kernelTime = NSTimer("kernelTime", false, false);
35     NSTimer timer_setUp_Comm = NSTimer("setUp_Communication_Timeee",
36         false, false);
37
38     while ( (i < HISTOGRAM_SIZE) && (histogram[i] <
39         CONTRAST_THRESHOLD) )
40     {
41         i++;
42     }
43     unsigned int min = i;
44
45     i = HISTOGRAM_SIZE - 1;
46     while ( (i > min) && (histogram[i] < CONTRAST_THRESHOLD) )
47     {
48         i--;
49     }
50     unsigned int max = i;
51     float diff = max - min;
52
53     int grayImageSize= width * height;
54
55     timer_setUp_Comm.start();
56     // Allocate device memory for grayImage
57     unsigned char *d_grayImage;
58     checkCudaCall(cudaMalloc((void **)&d_grayImage,
59         grayImageSize));
60
61     // Copy host memory to device
62     checkCudaCall(cudaMemcpy(d_grayImage, grayImage,
63         grayImageSize, cudaMemcpyHostToDevice));
64
65     // Setup execution parameters
66     dim3 threads(BLOCK_MAX_THREADSIZE);
67     dim3 grid(grayImageSize/threads.x);
68     timer_setUp_Comm.stop();
69
70     time_setup_comm[2]= timer_setUp_Comm.getElapsed();
71
72     // Kernel launch
73     kernelTime.start();
74     contrast1DKernel<<<grid, threads>>>(d_grayImage, width, height, min,
75         max, diff, grayImageSize);

```

```
72     cudaDeviceSynchronize();
73     kernelTime.stop();
74
75     kernelGpuTime[2]= kernelTime.getElapsed();
76
77     timer_setUp_Comm.start();
78     // Copy result from device to host
79     checkCudaCall(cudaMemcpy(grayImage,d_grayImage,
80                             grayImageSize,cudaMemcpyDeviceToHost));
81     timer_setUp_Comm.stop();
82     time_setup_comm[2] += timer_setUp_Comm.getElapsed();
83
84
85     cout << fixed << setprecision(6);
86     cout << "contrast1DCUDA (gpu): \t\t" << kernelTime.getElapsed()
87           << " seconds." << endl;
88
89     // clean device memory
90     cudaFree(d_grayImage);
91 }
```

## Appendix D

# Code for Image Smoothing

```
1 __global__ void triangularSmoothKernel(unsigned char *grayImage,
2   unsigned char *smoothImage, const int width, const int height,
3   const float *filter)
4 {
5   // The total number of threads per block times the number of
6   // blocks
7   unsigned int i = blockDim.x * blockIdx.x + threadIdx.x; // read-
8   only variable i
9   unsigned int j = blockDim.y * blockIdx.y + threadIdx.y; // read-
10  only variable j
11
12  if(i < width && j < height)
13  {
14    unsigned int filterItem = 0;
15    float filterSum = 0.0f;
16    float smoothPix = 0.0f;
17
18    for ( int fy = j - 2; fy < j + 3; fy++ )
19    {
20      for ( int fx = i - 2; fx < i + 3; fx++ )
21      {
22        if ( ((fy < 0) || (fy >= height)) || ((fx < 0) || (fx >=
23          width)) )
24        {
25          filterItem++;
26          continue;
27        }
28        smoothPix += grayImage[(fy * width) + fx] * filter[
29          filterItem];
30        filterSum += filter[filterItem];
31        filterItem++;
32      }
33    }
34  }
```

```

28     smoothPix /= filterSum;
29     smoothImage[(j * width) + i] = static_cast< unsigned char > (
        smoothPix);
30 }
31 }
32
33 void triangularSmoothCuda(unsigned char *grayImage, unsigned char *
        smoothImage, const int width, const int height,
34         const float *filter)
35 {
36     NSTimer kernelTime = NSTimer("kernelTime", false, false);
37
38     // 2a. Allocate the memory on the GPU
39     unsigned char *d_grayImage;
40     unsigned char *d_smoothImage;
41     float *d_filter;
42     checkCudaCall(cudaMalloc((void **)&d_grayImage, width * height));
43     checkCudaCall(cudaMalloc((void **)&d_smoothImage, width * height)
        );
44     checkCudaCall(cudaMalloc((void **)&d_filter, sizeof(filter)/
        sizeof(const float)));
45
46     // 2b. Move data over (host memory to device memory) for the
        function to execute
47     // cudaMemcpy(void *dst, void *src, size_t nbytes, enum
        cudaMemcpyKind direction);
48     checkCudaCall(cudaMemcpy((void *)d_grayImage, (void *)grayImage,
        (cudaMemcpyKind)width*height, cudaMemcpyHostToDevice));
49     checkCudaCall(cudaMemcpy((void *)d_smoothImage, (void *)
        smoothImage, (cudaMemcpyKind)width*height,
        cudaMemcpyHostToDevice));
50     checkCudaCall(cudaMemcpy(d_filter, filter, (cudaMemcpyKind)sizeof
        (filter)/sizeof(const float), cudaMemcpyHostToDevice));
51
52     // 3. Modify the function call in order to enable it to launch on
        the GPU
53     // Kernel invocation with one block of width * height * 1 threads
54     dim3 threadsPerBlock(16, 16);
55     dim3 numberOfBlocks(width/threadsPerBlock.x, height/
        threadsPerBlock.y);
56
57     kernelTime.start();
58     //Kernel invocation with N threads that executes it
59     triangularSmoothKernel<<<numberOfBlocks, threadsPerBlock>>> (
        d_grayImage, d_smoothImage, width, height, d_filter);
60     cudaDeviceSynchronize();
61     // /Kernel

```

```
62     kernelTime.stop();
63
64     // 4. Move data back over (device memory to host memory)
65     checkCudaCall(cudaMemcpy((void *)grayImage, (void *)d_grayImage,
66                               (cudaMemcpyKind)width*height, cudaMemcpyDeviceToHost));
67     checkCudaCall(cudaMemcpy((void *)smoothImage, (void *)
68                               d_smoothImage, (cudaMemcpyKind)width*height,
69                               cudaMemcpyDeviceToHost));
70
71     cout << fixed << setprecision(6);
72     cout << "triangularSmooth (cpu): \t" << kernelTime.getElapsed()
73           << " seconds." << endl;
74
75     // Free up device memory
76     cudaFree(d_grayImage);
77     cudaFree(d_smoothImage);
78     cudaFree(d_filter);
79 }
```

