



University of *Ljubljana*

Histogram Equalization with Cuda

**Final Project Report for Class:
Parallel and Distributed Systems and Algorithms**

Date: 05/02/2024

by Undergraduate Student

Irinej Slapal
Under the supervision of

Prof. Dr. Patricio Bulić

University of Ljubljana
Faculty of Computer and Information Science



Contents

1	Introduction	1
2	Parallel implementation with Cuda	2
2.1	Histogram Calculation	2
2.2	Comulative Distribution Function (CDF) Calculation	3
2.3	Equalization of the Image	9
3	Implementation testing	10
3.1	Image Sample 1	11
3.2	Image Sample 2	12
3.3	Image Sample 3	13
4	Results and Conclusion	14



1 Introduction

Histogram equalization is a fundamental technique in the field of digital image processing, used for improving the contrast of images. The primary goal of this method is to adjust the intensity distribution of an image, so that it spans the entire range of possible values more uniformly, so that the probabilities of all gray levels are somewhat close to equal.

Many images, especially those captured in low-light conditions or with limited dynamic range sensors, have poor contrast, resulting in concentration of pixel values in a narrow intensity range. This method is particularly useful for such images, adding contrast and bringing in details, invisible beforehand, to ones attention. The process of calculating a histogram is reading a color level (in this case, 8-bit grayscale level) from each pixel and track the appearance for each value. Finally the histogram shows how many pixels have a certain gray scale value.

The histogram H , we get from image, is the length of L :

$$H = \{n_0, n_1, n_2, \dots, n_{(L-1)}\}$$

Each value in histogram is equivalent to the number of pixels with the gray scale value l , $l \in [n_0, n_{L-1}]$. The probability of pixel having gray scale value of l , in an image, the size of $N \times M$ is calculated as:

$$p(l) = \frac{n_l}{N \times M}$$

The normalized histogram then contains the probabilities of the appearance of the different grey levels

$$H_{\text{norm}} = \{p(0), p(1), p(2), \dots, p(L - 1)\}$$

The cumulative distribution function tells us if $X \leq x$.

CDF corresponding to n_l is:

$$cdf(l) = p(X < l) = \sum_{i=0}^l p(i)$$

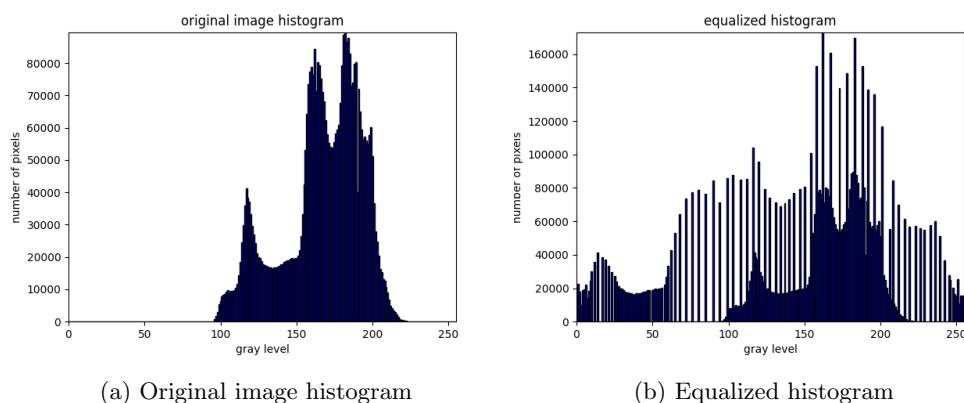


Figure 1: Comparison of original and equalized histograms, for image buca1.jpg

To sum it up: Histogram equalization redistributes the intensity values of an image, thereby enhancing its overall contrast and making details more visible.



2 Parallel implementation with Cuda

2.1 Histogram Calculation

Bellow is a sequential function for calculating Histogram from an image in C programming language. Notice, that the function iterates through every pixel in given image.

```

1 void CalculateHistogram(unsigned char* image, int width, int height, unsigned int *histogram
2 ) {
3     //Calculate histogram
4     for (int y=0; y<height; y++) {
5         for (int x=0; x<width; x++) {
6             histogram[image[y*width + x]]++;
7         }
8     }
}

```

There often is a case of image being too large, where the sequential approach would take too long.

There is a faster approach, in which threads are used to compute the histogram in parallel, as shown bellow:

```

1 __global__ void CalculateHistogramKernel(unsigned char* image, int width, int height,
2                                         unsigned int *histogram){
3
4     // calculate global x, y of pixel on image
5     int x = blockIdx.x * blockDim.x + threadIdx.x;
6     int y = blockIdx.y * blockDim.y + threadIdx.y;
7
8     // calculate local x, y of pixel in block
9     int lx = threadIdx.x;
10    int ly = threadIdx.y;
11
12    // allocate local memory for local histogram for each block
13    __shared__ unsigned long localHistogram[GRAYLEVELS];
14
15    // each thread sets its pixel in local (block) histogram to 0
16    localHistogram[blockDim.x * ly + lx] = 0;
17    __syncthreads();
18
19    //read value from image and increment local histogram according to its GRAYLEVEL
20    if (x < width && y < height) {
21        atomicAdd(&(localHistogram[image[y * width + x]]), 1);
22    }
23    __syncthreads();
24
25    //Finnaly coresponding local threads from each block sum their values to form a complete
26    //histogram
27    atomicAdd(&(histogram[ly * blockDim.x + lx]), localHistogram[ly * blockDim.x + lx]);
}

```

This kernel assigns each thread a unique pixel in the image, using the thread's global X and Y coordinates calculated from CUDA's built-in blockIdx, blockDim, and threadIdx variables. BlockSize must be 256 (16*16) here, so that histogram (which has length of 256) is correctly calculated.

GridSize is then calculated according to the size of input image, so that blocks "cover" the whole image. To optimize memory usage and access speed, shared memory is used within each block to store a local histogram of pixel intensities for each block in the grid.



Threads within a block read their corresponding pixel's intensity from the image stored in global memory and increment the appropriate entry in the local histogram.

Atomic operations are used, to avoid race conditions. After synchronizing all threads within the block to ensure the local histogram is complete, each thread then contributes its portion of the local histogram to a global histogram stored in global memory.

This approach utilizes the fast shared memory for intermediate calculations and optimizes global memory access.

2.2 Cumulative Distribution Function (CDF) Calculation

The Cumulative Distribution Function (CDF) is a statistical measure used to describe the distribution of a random variable.

In the context of image processing and analysis, the CDF of an image's pixel intensities provides a cumulative sum of the histogram counts, up to a certain intensity level. Essentially, for any given intensity value, the CDF indicates the proportion of pixels in the image that have an intensity less than or equal to that value.

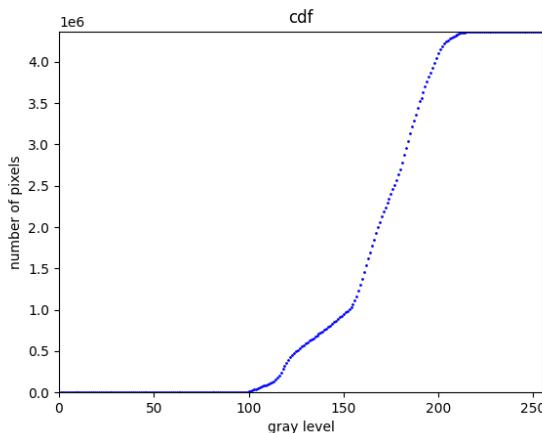


Figure 2: Comulative distribution function for image buca1.jpg.

In order to cumulatively sum histograms values (each value in resulting histogram, being the sum of all of its previous values) a sequential approach can be done like so:

```

1 void CalculateCDF(unsigned int *histogram, unsigned int *cdf){
2     // calculate cdf from histogram
3     cdf[0] = histogram[0];
4     for (int i=1; i<GRAYLEVELS; i++) {
5         cdf[i] = cdf[i-1] + histogram[i];
6     }
7 }
```

This doesn't cost as much as a sequential histogram calculation, because the number of iterations is always the same (256).



I wrote 3 different kernels for parallel cdf calculation:
The first one being the naive approach bellow:

```

1  __global__ void CalculateCDF_naive(unsigned int* histogram, unsigned int *cdf) {
2      __shared__ unsigned int temp[GRAYLEVELS*2];
3      int tid = threadIdx.x;
4
5      int pout = 0, pin = 1;
6
7      temp[tid] = histogram[tid];
8
9      __syncthreads();
10
11     for(int offset = 1; offset < GRAYLEVELS; offset <= 1) {
12         pout = 1 - pout;
13         pin = 1 - pout;
14         if (tid >= offset) {
15             temp[pout*GRAYLEVELS + tid] = temp[pin*GRAYLEVELS + tid] + temp[pin*GRAYLEVELS +
16             tid - offset];
17         } else {
18             temp[pout*GRAYLEVELS + tid] = temp[pin*GRAYLEVELS + tid];
19         }
20         __syncthreads();
21     }
22     cdf[tid] = temp[pout*GRAYLEVELS + tid];
}

```

This is a straightforward method of parallel prefix sum (scan) using shared memory and double buffering. In this method, each thread contributes by loading a value from the histogram into shared memory, and then iteratively updates its value based on the sum of preceding elements.

The kernel utilizes a loop that doubles the offset at each iteration, effectively summing pairs of elements spaced increasingly further apart. This iterative doubling reflects the essence of the naive approach, where the sum scan is built up through successive additions. Double buffering, achieved by toggling between two halves of a shared memory array, allows the kernel to read from one set of values while writing to another, preventing data overwrites during the summing process. This approach is very intuitive, but does not scale well with larger datasets.



Bellow is a better, work-efficient approach:

```

1  __global__ void CalculateCDFKernel_we(unsigned int* histogram, unsigned int*cdf) {
2      __shared__ unsigned int temp[GRAYLEVELS];
3
4      int tid = threadIdx.x; // 1block 1x128 threads, 128 threads
5      int offset = 1; // distance between elements in array that will be summed
6
7      // the sum of values, that each thread calculates in 1st step
8      temp[2*tid] = histogram[2*tid];
9      temp[2*tid+1] = histogram[2*tid+1];
10
11     for (int d = GRAYLEVELS >> 1; d > 0; d >>= 1) {
12         __syncthreads();
13
14         if (tid < d) {
15             int ai = offset*(2*tid+1)-1;
16             int bi = offset*(2*tid+2)-1;
17             temp[bi] += temp[ai];
18         }
19         offset *= 2;
20     }
21
22     if (tid == 0) {
23         temp[GRAYLEVELS - 1] = 0;
24     }
25
26     for (int d = 1; d < GRAYLEVELS; d *= 2) {
27         offset >>= 1;
28         __syncthreads();
29
30         if (tid < d) {
31             int ai = offset*(2*tid+1)-1;
32             int bi = offset*(2*tid+2)-1;
33
34             float t = temp[ai];
35             temp[ai] = temp[bi];
36             temp[bi] += t;
37         }
38     }
39     __syncthreads();
40     cdf[2*tid] = temp[2*tid];
41     cdf[2*tid+1] = temp[2*tid+1];
42 }
43

```

The implementation of this kernel is based on the following source

Based on the parallel prefix sum (scan) algorithm, detailed in GPU Gems 3[1], this kernel uses a single block with 128 threads to process the gray levels histogram f . Initially, each thread loads two consecutive histogram values into shared memory.

The kernel performs the scan in two phases: an **up-sweep** (reduction) phase, and a **down-sweep** phase.



During the **up-sweep**, starting with a stride of 1, the algorithm recursively doubles the offset, summing pairs of elements spaced apart by the offset and storing the result in the position of the second element in the pair.

This process constructs a sum tree in shared memory, where the final sum (which is total number of pixels) ends up at the last position of the array. However, to convert this into a CDF, a zero is placed at the end of the array (for the last element of the CDF to represent the total count) before the down-sweep phase.

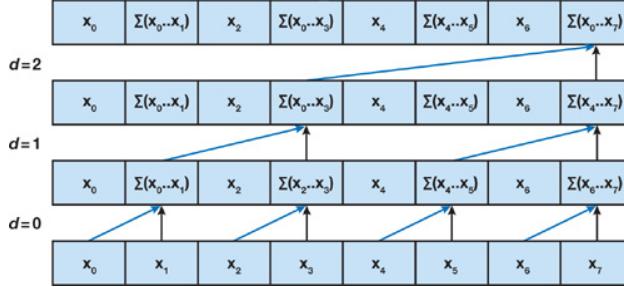


Figure 3: Illustration of up-sweep.

In the **down-sweep** phase, the kernel reverses the process, halving the offset at each step and using a simple swapping and summing technique to propagate the sums back through the array in a way that constructs the CDF from the sum tree. This phase starts by setting the last element to zero to ensure the correct initialization for the inclusive scan.

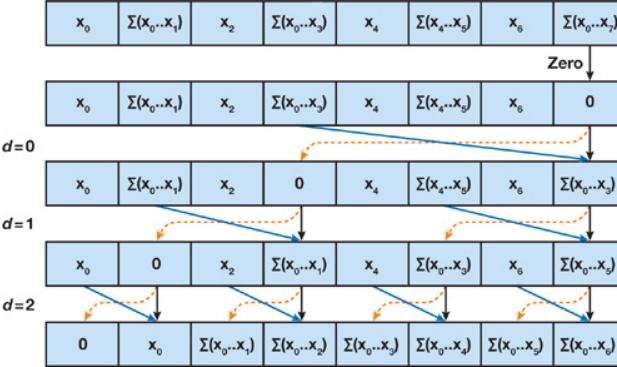


Figure 4: Illustration of down-sweep.

After completing both phases, each thread writes two elements of the computed CDF back to global memory.

However the efficient approach still has risk of memory bank conflicts.

Memory bank conflicts arise in GPU computing when multiple threads within the same warp access data from the same memory bank of shared memory simultaneously, and at least one of the operations is a write. This would occur after the 2 steps, when each thread writes two values to global memory.

Since each memory bank can service only one access per cycle, simultaneous accesses to the same bank cause serialization, leading to increased memory access times and reduced overall performance. These conflicts are a consequence of the parallel architecture of GPUs, where shared memory is divided into multiple banks to allow concurrent accesses by different threads.

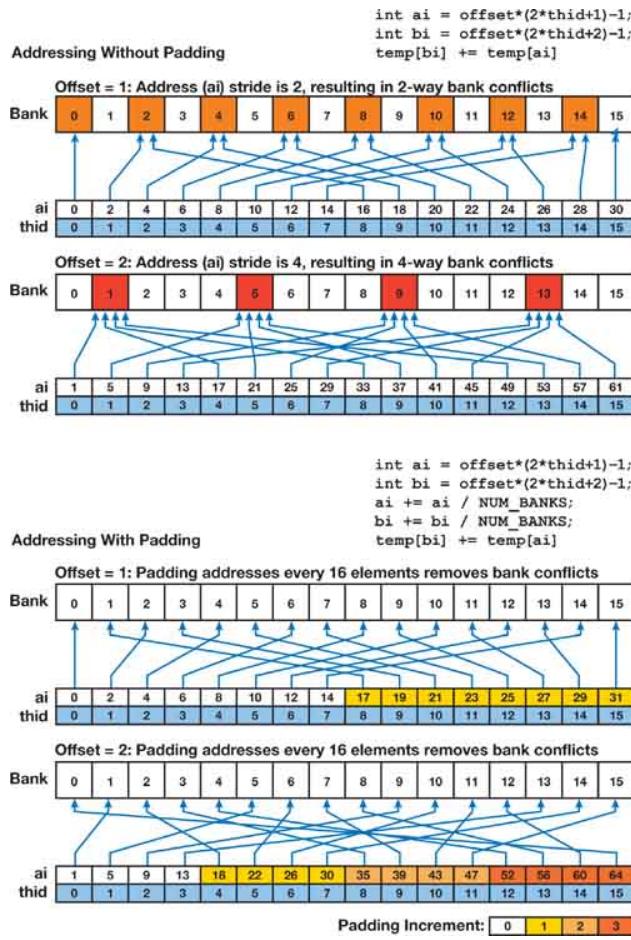


Figure 5: Illustration of memory bank conflicts.

With information given in GPU Gems 3, we can arrive to a strategy for mitigating the conflicts by padding, (intentionally inserting extra space between elements in shared memory) to ensure that consecutive threads access different memory banks. This improves memory access patterns and enhances the efficiency and speed of the kernel shown on next page.



```

1  __global__ void CalculateCDF_we_mbcf(unsigned int* histogram, unsigned int *cdf) {
2      __shared__ unsigned int temp[GRAYLEVELS + CONFLICT_FREE_OFFSET(GRAYLEVELS)];
3
4      int tid = threadIdx.x; // 1block 1x128 threads, 128 threads
5      int offset = 1; // distance between elements in array that will be summed
6
7      int ai= 2*tid;
8      int bi= tid + (GRAYLEVELS/2);
9      int bankOffsetA = CONFLICT_FREE_OFFSET(ai);
10     int bankOffsetB = CONFLICT_FREE_OFFSET(bi);
11
12     temp[ai + bankOffsetA] = histogram[ai];
13     temp[bi + bankOffsetB] = histogram[bi];
14
15     // the sum of values, that each thread calculates in 1st step
16
17     for (int d = GRAYLEVELS >> 1; d > 0; d >>= 1) {
18         __syncthreads();
19
20         if (tid < d) {
21             ai = offset*(2*tid+1)-1;
22             bi = offset*(2*tid+2)-1;
23             temp[bi] += temp[ai];
24         }
25         offset *= 2;
26     }
27
28     if (tid == 0) {
29         temp[GRAYLEVELS - 1 + CONFLICT_FREE_OFFSET(GRAYLEVELS - 1)] = 0;
30     }
31
32     for (int d = 1; d < GRAYLEVELS; d *= 2) {
33         offset >= 1;
34         __syncthreads();
35
36         if (tid < d) {
37             ai = offset*(2*tid+1)-1;
38             bi = offset*(2*tid+2)-1;
39
40             float t = temp[ai];
41             temp[ai] = temp[bi];
42             temp[bi] += t;
43         }
44     }
45     __syncthreads();
46     cdf[ai] = temp[ai + bankOffsetA];
47     cdf[bi] = temp[bi + bankOffsetB];
48 }

```



2.3 Equalization of the Image

The new (equalized) grayscale value is calculated for each pixel like so:

$$\text{new} = \frac{\text{cdf}(l) - \text{cdf}_{\min}}{(N \times M) - \text{cdf}_{\min}} \times (L - 1)$$

- l_{new} ... being the new value,
- $(N \times M)$... being the size of an image,
- $\text{cdf}(l)$... being the number of pixels on input image with grayscale value l ,
- cdf_{\min} ... being the smallest **non-zero** value in CDF histogram, 0-s represent no accumulation in the histogram
- L ... in this case 256, which is the number of gray levels

A kernel for finding the smallest non-zero value in an array, using principle of reduction:

```

1  __global__ void findMinKernel(unsigned int* cdf, unsigned int*d_cdfmin) {
2      // Allocate shared memory
3      __shared__ unsigned int partial_mins[256];
4
5      // Calculate thread ID
6      int tid = threadIdx.x;
7      // Load elements into shared memory
8      // we are looking for the smallest NON-ZERO value in CDF so we can UINT_MAX all the
9      // zeros
10     if (tid < 128) {
11         partial_mins[tid] = cdf[tid] == 0 ? UINT_MAX : cdf[tid];
12         partial_mins[tid + 128] = cdf[tid + 128] == 0 ? UINT_MAX : cdf[tid + 128];
13     }
14     // Start at 1/2 block stride and divide by two each iteration
15     for (int s = GRAYLEVELS/2; s > 0; s >>= 1) {
16         __syncthreads();
17         // Each thread does work unless it is further than the stride
18         if (tid < s) {
19             partial_mins[tid] = min(partial_mins[tid], partial_mins[tid + s]);
20         }
21     }
22     __syncthreads();
23     if (threadIdx.x == 0) {
24         *d_cdfmin = partial_mins[0];
25     }
}

```

```

1  __device__ unsigned char scale(unsigned int cdf, unsigned int cdfmin, unsigned int imageSize
2      ) {
3      float scale;
4      scale = (float)(cdf - cdfmin) / (float)(imageSize - cdfmin);
5      scale = round(scale * (float)(GRAYLEVELS-1));
6      return (int)scale;
}

```

Device function scale calculates the new value of pixel, according to the equation before.



Kernel below builds the new equalized image in parallel, by calling scale function for each pixel with each thread for unique pixel in image. BlockSize here is 256 (16*16) and GridSize is calculated the same way as for the histogram calculation kernel.

```

1  __global__ void EqualizeKernel(unsigned char * image_in, unsigned char * image_out, int
2   width, int height, unsigned int *cdf, unsigned int *cdfmin) {
3   unsigned int imageSize = width * height;
4   int x = blockIdx.x * blockDim.x + threadIdx.x;
5   int y = blockIdx.y * blockDim.y + threadIdx.y;
6
6   //Equalize
7   if (x < width && y < height){
8       image_out[(y*width + x)] = scale(cdf[image_in[y*width + x]], *cdfmin, imageSize);
9   }
10 }
```

3 Implementation testing

I have tested my implementation on 7 different images of varying sizes. I used cudaEvent_t sturct for timing the execution. I timed 3 different executions:

- One where all 3 steps are in parallel (Cuda Parallel)
- One where all 3 are sequential (Sequential)

- and One, where I implemented the kernels for sequential execution

and ran only 1 thread (Kernel Sequential) For normalization purposes each implementation was run 10 times.

Image Name	Size of Image ($h \times w$)	Cuda Parallel	Sequential	Kernel Sequential
flower2.jpg	280 x 180	0.1968 ms	1.2869 ms	20.8876 ms
pa3cio.jpg	400 x 250	0.1918 ms	2.5469 ms	52.0048 ms
ograja.jpg	400 x 600	0.2159 ms	6.1029 ms	230.8772 ms
ferari.jpg	600 x 338	0.2082 ms	5.1372 ms	197.1191 ms
puscava.jpg	1024 x 683	0.3530 ms	17.7262 ms	1513.1527 ms
flower1.jpg	1280 x 720	0.4385 ms	23.3550 ms	2380.4128 ms
buca1.jpg	2560 x 1707	1.1806 ms	110.7117 ms	49685.0430 ms

Table 1: Times for parallel and sequential processing.

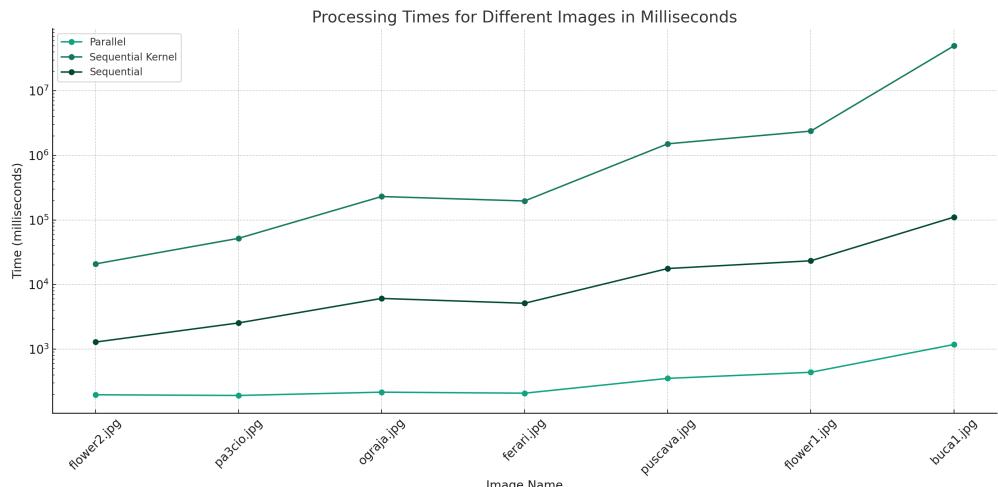


Figure 6: Graph of times for parallel and sequential processing.



3.1 Image Sample 1

bucal1.jpg

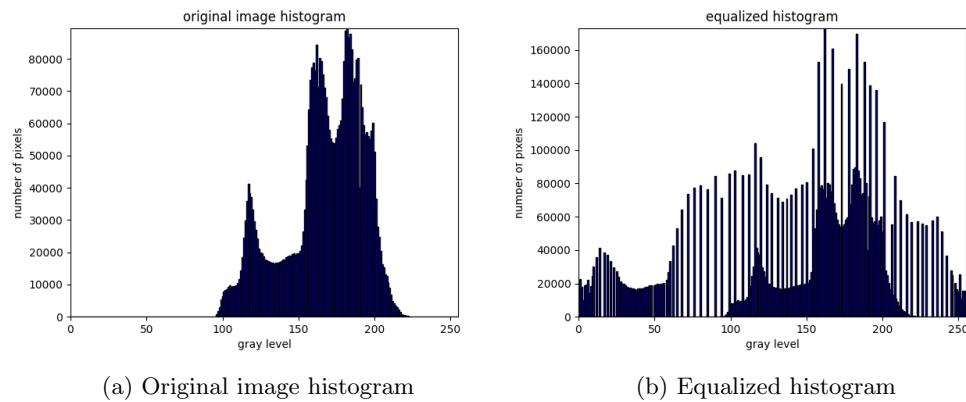


Figure 7: Comparison of original and equalized histograms, for image buca1.jpg



Figure 8: Comparison of original and equalized histograms, for image buca1.jpg



3.2 Image Sample 2

pa3cio.jpg

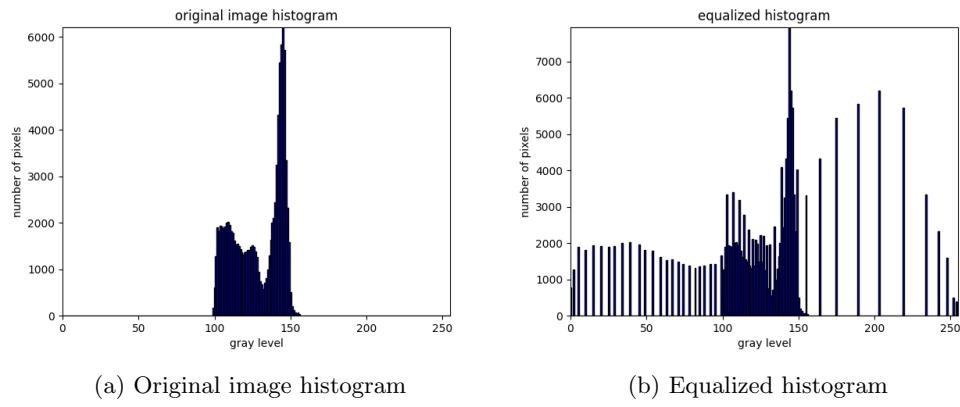


Figure 9: Comparison of original and equalized histograms, for image pa3cio.jpg



Figure 10: Comparison of original and equalized images, for image pa3cio.jpg

3.3 Image Sample 3

flower1.jpg

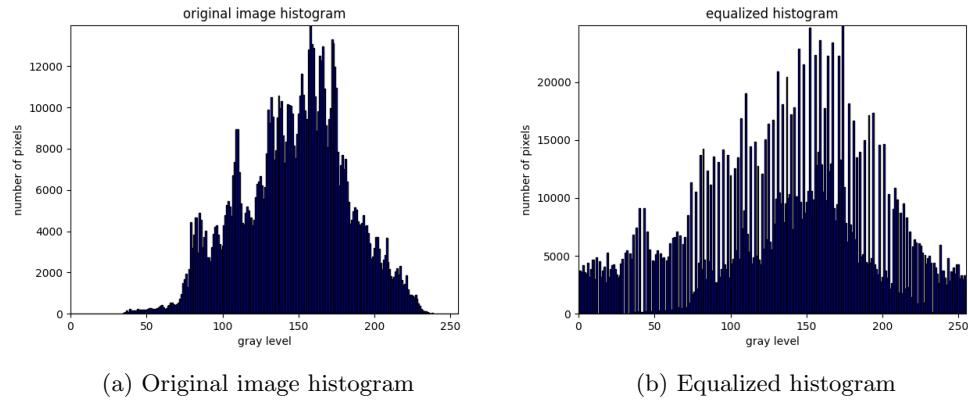


Figure 11: Comparison of original and equalized histograms, for image buca1.jpg



Figure 12: Comparison of original and equalized histograms, for image flower1.jpg



4 Results and Conclusion

In conclusion of this report, the comparative analysis between parallel and sequential processing methodologies for the task of gray-scale image histogram equalization has shown a substantial difference between each. The results unambiguously support the hypothesis that parallel computing, particularly when executed on Graphics Processing Units (GPUs), is superior in performance relative to its sequential counterpart running on Central Processing Units (CPUs).

Throughout the range of sample sizes subjected to testing, it has been consistently observed that the parallel implementation outpaces sequential processing. The distinction in performance, while noticeable even at smaller scales, becomes stark as the dimensions of the image increase. This trend can be attributed to the intrinsic nature of parallel computing, which, unlike sequential processing, can simultaneously handle multiple data points. GPUs, with their multitude of cores designed for concurrent operations, are particularly adept at managing large sets of pixel data inherent in high-resolution images.

Parallel processing is not always the optimal approach due to several limiting factors, most notably described by Amdahl's Law. Amdahl's Law provides a theoretical maximum improvement in performance using parallel processing and highlights the importance of sequential execution time in the overall computation. The law can be expressed by the formula:

$$\text{latency}(s) = \frac{1}{(1-p) + \frac{p}{s}} \quad (1)$$

where

- S_{latency} is the theoretical speedup of the execution of the whole task,
- s is the speedup of the part of the task that benefits from improved system resources,
- p is the proportion of the task that can be parallelized (i.e., the proportion of execution time that the part benefiting from improved resources originally occupied).

According to Amdahl's Law, the speedup of a program using multiple processors in parallel computing is limited by the time needed for the sequential fraction of the program. For instance, if only 50% of a program can be parallelized, the maximum speedup using parallel computing would be 2x, no matter how many processors are used. This implies that there is a point of diminishing returns where adding more processing power yields minimal improvements. This is because the non-parallelizable part of the program (the sequential component) becomes a bottleneck.



References

- [1] Shubhabrata Sengupta, Mark Harris, and John D. Owens. Chapter 39. parallel prefix sum (scan) with cuda.