

# Concurrency and Parallelism Project

Guilherme Fernandes

Vladyslav Mikytiv

60045

60735

## ABSTRACT

This report provides the explanation and the implementation of parallelisation of the algorithm of equalization of images. Both an OpenMP and CUDA implementations.

## KEYWORDS

OpenMP Parallelism C++ CUDA

## 1 CRITICAL CODE ANALYSIS AND PARALLELIZATION

### 1.1 The first steps

At first we ran the code without doing any modifications and used the profiler in order to identify the hot spots doing a average value in each image of the time lost in functions from the profiler. The algorithm runs for 100 iterations.

| Function             | Computation Time Lost (%) |
|----------------------|---------------------------|
| correctColor&rescale | 57.98%                    |
| grayScale            | 23.17%                    |
| normalize            | 18.24%                    |

Table 1: Average Computation Time Lost in Functions

As seen in Listing 1 the zones of the code that took a long time to perform the computation and in overall slow down the algorithm are: the **normalization** of the image, the **correction of color**, the **conversion to the greyscale** and some other not so computational heavy functions.

To reduce the impact on processing time caused by these functions within the algorithm, we integrated OpenMP directives into these specific parts of the code to accelerate computation.

It's important to highlight the **thread management**. Since we are working with images we will have to do some simples calculus to determine the amount of work that each thread will be given. For that we must perform two calculations (one for RGB and one for the grey scale). For the RGB we will calculate  $WIDTH * HEIGHT * CHAN\_SIZE$  (the 3 RGB channels) and to obtain the  $CHUNK\_SIZE\_RGB$  we just divide  $WIDTH * HEIGHT * CHAN\_SIZE$  by  $N\_THREADS$ . For the grey scale it's similar but we don't multiply by  $CHAN\_SIZE$  getting the  $CHUNK\_SIZE$ . Now we will have the work that will be balanced between threads.

### 1.2 Function parallelization

The **normalization** function will be improved with the following code: **pragma omp parallel for** with the option **schedule(static, chunk\_size\_channels)** and **num\_threads(n\_threads)**. This code optimization can be seen in Listing 1.

```
1 void normalize(//omitting for space) {
2     #pragma omp parallel for schedule(static,
3     chunk_size_channels) num_threads(n_threads)
4     for (int i = 0; i < size_channels; i++)
5         uchar_image[i] = (unsigned char) (255 *
6         input_image_data[i]);
7 }
```

Listing 1: Normalization Function

The **grey scale conversion** will be improved in two ways. We will mix the **fill\_histogram** function with this one in order to do everything in the same function. Besides that we also apply the following OpenMP directives: **pragma omp parallel for** with the option **reduction(+: histogram)** and **num\_threads(n\_threads)**. This code optimization can be seen in Listing 2.

```
1 void extractGrayScale(//omitting for space) {
2     // filling the histogram with zeroes
3     #pragma omp parallel for reduction(+: histogram)
4     num_threads(n_threads)
5     for (int i = 0; i < size; i++){
6         auto r = uchar_image[num_chan * i];
7         auto g = uchar_image[num_chan * i + 1];
8         auto b = uchar_image[num_chan * i + 2];
9         gray_image[i] = // Calculated
10        histogram[gray_image[i]]++;
11    }
```

Listing 2: Grey Conversion Function

We can't parallelize the **CDF calculation** because it has true dependencies between iterations and also it's not worth due to the number of iterations (256). Another function that we managed to simplify was the **cdf\_min\_loop**. Since we minimum will always be on the first position we just return it and that's how we compute the minimum of the CDF.

The **correct\_color\_loop** can be mixed with the **rescale** following the same strategy as we did with grey scale and fill histogram functions. The code optimization is the same as Listing 1 and can be seen in Listing 3.

```
1 void correct_color_loop_and_rescale(//omitting for
2 space) {
3     #pragma omp parallel for schedule(static,
4     chunk_size_channels) num_threads(n_threads)
5     for (int i = 0; i < size_channels; i++)
6     {
7         uchar_image[i] = correct_color(cdf[
8         uchar_image[i]], cdf_min);
9         output_image_data[i] = static_cast<float>(<
10        uchar_image[i]) / 255.0f;
11    }
```

Listing 3: Color Correction Function

### 1.3 Decisions

It's important to justify our reasoning. We used **static** in every single call of the OpenMP directives because the extra overhead of orchestrating the threads with a certain amount of work did not compensate. As we can see in the Table 2 **static** always shows better results and it's coherent to the theoretical analysis. For this simulation we used 16 threads and 100 iterations of the algorithm.

| Dynamic Speedup | Static Speedup | Image Name   |
|-----------------|----------------|--------------|
| 9.7             | 11.5           | borabora.ppm |
| 9.7             | 11.8           | input01.ppm  |
| 7.7             | 9.9            | sample.ppm   |

Table 2: Dynamic vs Static Mean Execution Times

## 2 METRIC ANALYSIS - OPENMP

Now, let's assess the impact of these changes on the runtime of our program. For that we will use **speedup** and **efficiency** and calculate the mean values of 5 runs for each thread.

The code was in a cluster executed with the following specifications: **2x Intel Xeon E5-2609 v4**, with **16 cores** and **NVIDIA Quadro M2000 GPU**.

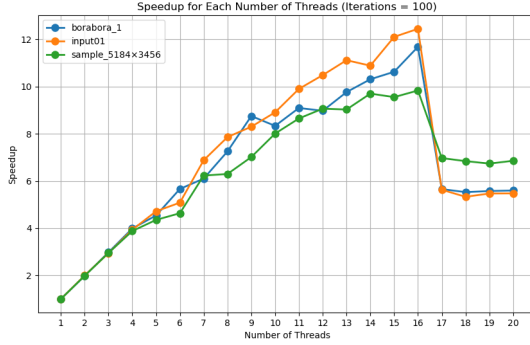


Figure 1: Speedup across different number of threads usage of different images

Figure 1 describes the speed up for the different images in our test set and Table 3 the standard deviations in the execution time. We can clearly see a rise on the speed up until a certain point. After that point the value of the speed up begins to decrease. The optimal NR\_THREADS to perform the computation is 16.

In order to calculate the efficiency we will use the best result that we had: the one with 16 threads. For that we define  $S_{16}$  which is the speed up with 16 threads and it will be equal to

$$S_{16} \approx 11$$

And the efficiency of the parallel implementation in Section 1 get's us an efficiency value of:

$$E_{16} = \frac{S_{16}}{16} \approx 69\%$$

| N° of Threads | borabora.ppm | input01.ppm | sample.ppm |
|---------------|--------------|-------------|------------|
| 1             | 21           | 18          | 18         |
| 2             | 16           | 21          | 59         |
| 3             | 9            | 10          | 249        |
| 4             | 3            | 15          | 230        |
| 5             | 19           | 14          | 403        |
| 6             | 15           | 15          | 77         |
| 7             | 28           | 4           | 98         |
| 8             | 16           | 3           | 313        |
| 9             | 3            | 9           | 874        |
| 10            | 20           | 9           | 452        |
| 11            | 5            | 10          | 58         |
| 12            | 2            | 8           | 251        |
| 13            | 1            | 9           | 84         |
| 14            | 3            | 6           | 59         |
| 15            | 5            | 7           | 37         |
| 16            | 3            | 11          | 46         |
| 17            | 4            | 5           | 47         |
| 18            | 3            | 2           | 48         |
| 19            | 4            | 2           | 128        |
| 20            | 2            | 1           | 70         |

Table 3: Standard Deviations of Execution Times

## 3 CUDA

We've attained speed enhancements leveraging the OpenMP library for parallelizing our code. However, there exists an avenue to further escalate our computational efficiency. By harnessing the GPU architecture through CUDA, we can unlock additional acceleration potential.

To leverage the GPU effectively, we need to define kernels for specific regions of our code that we wish to offload to the GPU for computation. This necessitates defining grids use in CUDA kernels.

Utilizing the GPU for optimization involves creating kernels to transfer data to the GPU, perform computations, and retrieve the results back to the CPU. Given the high cost associated with these operations, minimizing the number of kernels is imperative.

The initial kernel handles both the normalization and grayscale operations. We can merge these two operations in one kernel since we don't need any prior computation before to calculate it. We don't need to create an extra overhead with two kernels in here. The Listing 4 shows this kernel implementation.

```

1 __global__ void normalize_kernel(//omitting for space) {
2     int ii = blockIdx.y * blockDim.y + threadIdx.y;
3     int jj = blockIdx.x * blockDim.x + threadIdx.x;
4     int idx = (ii * width + jj)*3;
5     if (ii < height && jj < width) {
6         for(int i = 0; i < num_chan; i++)
7             // Normalize logic
8             auto r = uchar_image[idx];
9             auto g = uchar_image[idx + 1];
10            auto b = uchar_image[idx + 2];
11            idx = ii * width + jj;
12            gray_image[idx] = // Calculation
13        }
14    }

```

Listing 4: Grey Conversion and Normalization

Before proceeding to the second kernel, it's essential to calculate the histogram. While we could create a custom kernel for this task, a more efficient approach involves leveraging the **CUB** library. By utilizing its operations, we can significantly enhance the performance of our code.

Using `cub::DeviceHistogram::HistogramEven` we are able to calculate the histogram in the GPU without making the kernel ourselves.

After having the histogram we can calculate the probabilities in order to use them in another call to the **CUB** library. For that we create a second kernel, which will seem more clear why ahead. This one will differ from the first one in the `numBlocks` and `blockSize`. The `numBlocks` will be 256 and the `blockSize` will be  $(\text{HISTOGRAM\_LEN} + \text{blockSize} - 1) / \text{blockSize}$ , that's why we couldn't re-use the first kernel due to this size and the CUB library usage.

Once we have computed this probability array, we can further utilize it in another function provided by the **CUB** library.

The `cub::DeviceScan::InclusiveSum` function facilitates the calculation of the final cumulative distribution function (CDF) values. This not only streamlines the programming process but also optimizes computation time, offering significant efficiency gains. We used the `InclusiveSum` because it was not possible to use `cub::DeviceScan::InclusiveScan` due to restrictions in the library (there was no way to pass the probabilities to this call). That's why we must calculate the probabilities beforehand and use the `InclusiveSum`.

Ultimately, we employ our final kernel to run the last part of the algorithm. The code can be seen in Listing 5.

```
1 __global__ void correct_kernel(//omitting for space) {
2     int ii = blockIdx.y * blockDim.y + threadIdx.y;
3     int jj = blockIdx.x * blockDim.x + threadIdx.x;
4     int idx = (ii * width + jj)*3;
5
6     if (ii < height && jj < width) {
7         for (int i = 0; i < num_chan; i++) {
8             auto cdf_val = d_cdf[uchar_image[idx+i]];
9             float cdf_min = d_cdf[0];
10
11             uchar_image[idx+i] = // Logic
12
13             output_image_data[idx+i] = // Logic
14         }
15     }
16 }
```

Listing 5: Correct and Rescale

These kernels constitute the entirety of the GPU-accelerated operations. Due to space constraints, significant portions of the code have been omitted. For the complete implementation, please refer to the accessible code at <https://github.com/Gui28F/CP-Project>.

## 4 METRIC ANALYSIS - CUDA

Utilizing the capabilities of the GPU necessitates defining the grid dimensions. To determine the optimal `TILE_WIDTH` value, we conducted an experiment, the results of which are depicted in Figure 2. It was discerned that the most favorable value was 16. Notably, we are constrained from exceeding 32 due to the limitation imposed by the maximum number of thread per block, capped at  $32 * 32 = 1024$ .

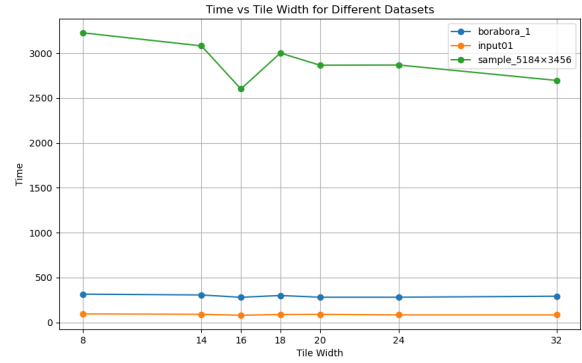


Figure 2: Tiles Width and Execution Time

To effectively evaluate the performance of the GPU-accelerated version, it's important to compare it against the most optimized configuration achieved with OpenMP (specifically, the one utilizing 16 threads). Figure 3 illustrates the differences in speedup compared to the optimal OpenMP configuration (and the others also).

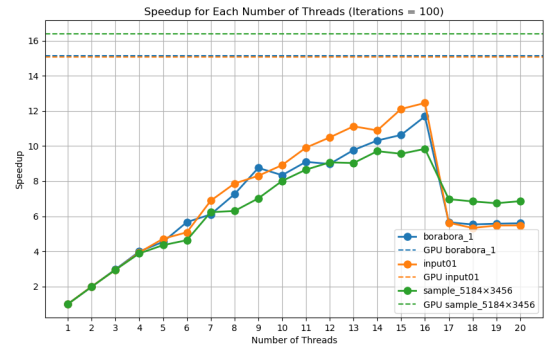


Figure 3: OpenMP vs CUDA

## 5 FINAL RESULTS

After implementing optimizations using **OpenMP** and **CUDA**, we now have a clearer understanding of the overall performance. The table presents execution times for the purely sequential implementation, the **OpenMP** implementation, and the **CUDA** implementation. These results, obtained from 100 iterations of the algorithm, represent the mean calculated values derived from running the algorithm five times. The names were shortened in the table <sup>1</sup>.

Each generation of results was derived from the average of five separate runs the algorithm with 100 iterations. By calculating the mean of these five separate executions, we aimed to obtain more reliable and trustworthy results. This approach ensured that our findings were not skewed by anomalies in any single run, providing a more accurate representation of the algorithm's performance.

<sup>1</sup> **Description:** SQT – Sequential Time, OMPT – OpenMP Time, CUDAT – CUDA Time

| SQT (ms) | OMPT (ms) | CUDAT (ms) | Image        |
|----------|-----------|------------|--------------|
| 1140     | 98        | 76         | borabora.ppm |
| 956      | 76        | 63         | input01.ppm  |
| 43593    | 4371      | 2635       | sample.ppm   |

**Table 4: Final Execution Times**

| OpenMP Speedup | CUDA Speedup | Image        |
|----------------|--------------|--------------|
| 11.5           | 15           | borabora.ppm |
| 11.8           | 15.2         | input01.ppm  |
| 9.9            | 16.5         | sample.ppm   |

**Table 5: Final Speedup Values**

As observed in all the images of Table 4 and Table 5, even with the best OpenMP configuration, CUDA consistently outperforms it in terms of speed. The superior capabilities of the GPU significantly accelerate our computations, and this performance difference becomes increasingly noticeable with larger image sizes.

## 6 TESTS

During development, we utilized GoogleTest to guide our implementations for both OpenMP and CUDA. These tests ensured the accuracy of the pixel values in the resulting images, verifying correctness beyond mere visual inspection. Additionally, the test files

used are available in our GitHub repository, providing transparency and allowing others to validate our results.

## 7 CONCLUSIONS

Whether through the integration of OpenMP directives or harnessing the power of GPU acceleration, we consistently observe improved results compared to the raw sequential implementation. In all our tests, CUDA consistently delivers superior performance compared to OpenMP. The enhanced capabilities of GPUs allow CUDA to handle computations more efficiently, resulting in faster processing times across various image sizes and complexities. This makes CUDA the preferred choice for achieving optimal performance in our implementations.

## 8 CODE

All the code mentioned above can be consulted in <https://github.com/Gui28F/CP-Project>.

## 9 INDIVIDUAL CONTRIBUTION AND COMMENTS

The project was a joint effort, with both team members sharing the workload equally. We utilized pair programming sessions on Discord calls and in the lab to accomplish our tasks effectively. We are grateful to Professor Hervé for his invaluable support and constant availability to assist us throughout the project.