

Parallel programming, convolutions in CUDA

Christian Zanzi

E-mail address

christian.zanzi@edu.unifi.it

Abstract

In this report, we analyze the speedup achievable using a GPU to calculate a convolution on an image. First, we consider a sequential program written in C++ that works on a CPU, then we compare its execution time with a new program written in CUDA that works on NVIDIA GPUs. In this writing, some tricks are shown to make the convolution quicker on a GPU, knowing that the data transfer is the most burdensome operation.

Future Distribution Permission

The author(s) of this report give permission for this document to be distributed to Unifi-affiliated students taking future courses.

1. Introduction

In image processing and computer vision, image convolutions are widely used to find features of images or to modify them. It is a simple mathematical operation done between matrices or tensors of numeric values.

A filter is applied to each pixel of an image to obtain an output value. The **Figure 1** explains graphically the operation done on a single pixel. In practice is a sum of multiplication between the filter and the pixel values surrounding the pixel in which the filter is centered.

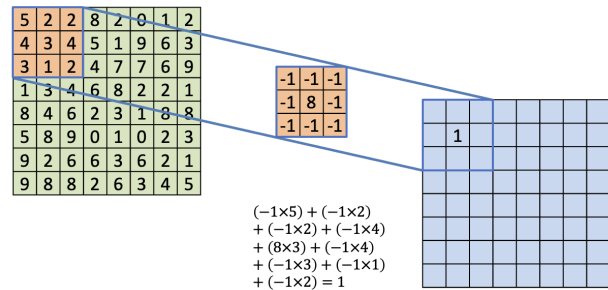


Figure 1. Filter applied on a single pixel.

1.1. Test image and filter

The tests are done using a 2048x2048 pixels image on which an edge detection filter is applied to obtain information on the edges inside the original image.

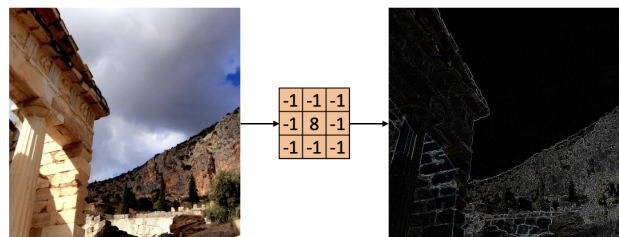


Figure 2. Filter applied on a single pixel.

1.2. Development methodologies

C++ is used as the programming language for the sequential program's code, and CUDA is used for the parallel version.

The sequential version just calculates linearly for each value of the matrix the convolution operation, while the CUDA program splits the execution across the streaming multiprocessors of a GPU.

Some tricks to speed up the convolution are used to store and load the image in the parallel program:

- use of the constant cache of the streaming multiprocessors to store the filter;
- use of the shared memory to store the tiles assigned from the GPU's scheduler to each streaming multiprocessor;
- 2 batch loading.

1.3. Performance evaluation

The speedup $\frac{t_{seq}}{t_{par}}$ measures the performance improvement obtained from the parallelization of the sequential code. Each program is executed 5 times, and the execution time is recorded with the *chrono* library. The sequential and parallel programs have been timed only considering the convolution operation without taking into account the time spent for the data transfer.

1.4. Test environment

The sequential program is runned on

Model name: MacBook Air

Chip: Apple M3

total number of Cores: 8 (prestazioni 4 ed efficienza 4)

Memory: 16 GB

Using the Mac terminal to run it:

```
1 clang++ -o sequential_convolution
  sequential_convolution.cpp
2 ./sequential_convolution input.jpg output.jpg
```

For the parallel program, we used the service offered by Amazon Web Services (AWS), called Elastic Compute Cloud (EC2), which allows scalable computational capability.

We launched a G4dn instance that includes the NVIDIA T4 GPU, useful for our problem testing.

Instance name: g4dn.2xlarge

Chip: NVIDIA T4

total number of Cuda Cores: 2560

Memory: 32 GiB

To run the program, we used the instance terminal:

```
1 nvcc -o cuda_convolution cuda_convolution.cu
2 ./cuda_convolution input.jpg output.jpg
```

2. Programs comparison

In this section, we will discuss the difference between the programs and their implementation. We used the *stbi* library to load the image and save it after the convolution. In the two programs, the code used for image loading is unaltered

```
1 int width, height, channels;
2 unsigned char *input_image = stbi_load(argv[1], &
  width, &height, &channels, 0);
3 if (!input_image) {
4     std::cerr << "Failed to load image\n";
5     return 1;
6 }
7
8 size_t image_size = width * height * channels;
9 float *host_input = new float[image_size];
10 float *host_output = new float[image_size];
11 for (size_t i = 0; i < image_size; ++i)
12     host_input[i] = input_image[i] / 255.0f;
```

With *stbi_load()* we are obtaining an unsigned char array with each pixel's value. It is common to normalize these values for:

- transforming them to floating points, since GPUs work better with floats;
- increasing the numeric precision and limiting the risk of overflow;
- controlling the results obtained from the convolution with filters with small values.

2.1. Sequential program

The convolution function of the sequential program loops over channels, height rows, and width columns to apply the filter convolution on every pixel.

Before making the sums, the code checks if the index on which the operations should be done is correct.

```

1 void convolution(float *input, float *output,
2               float *filter, int channels, int width, int
3               height) {
4     // image matrix memorized in row-major order,
5     // linear memory access to consecutive addresses
6     for (int k = 0; k < channels; ++k) {
7         for (int y = 0; y < height; ++y) {
8             for (int x = 0; x < width; ++x) {
9                 float sum = 0;
10
11                 for (int fy = - FILTER_RADIUS; fy <=
12                     FILTER_RADIUS; ++fy) {
13                     for (int fx = - FILTER_RADIUS; fx
14                         <= FILTER_RADIUS; ++fx) {
15                         // setting the index position
16                         // of the image tile with respect to the filter
17                         int nx = x + fx;
18                         int ny = y + fy;
19                         // border handling with
20                         // replicate padding for the ghost cells
21                         if (nx < 0) nx = 0;
22                         if (ny < 0) ny = 0;
23                         if (nx >= width) nx = width -
24                             1;
25                         if (ny >= height) ny = height
26                             -1;
27
28                         const int idx = (ny * width +
29                             nx) * channels + k;
30                         sum += input[idx] * filter[(fx
31                             + FILTER_RADIUS) * FILTER_WIDTH + (fy +
32                             FILTER_RADIUS)];
33                     }
34                     int out_idx = (y * width + x) *
35                         channels + k;
36                     output[out_idx] = sum;
37                 }
38             }
39         }
40     }
41 }

```

The program prints the execution time at the end and saves the convoluted image.

2.2. Parallel program

The parallel program has an initial definition of constant variables that do not change during execution.

```

1 #define FILTER_WIDTH 3
2 #define FILTER_RADIUS (FILTER_WIDTH / 2)
3 #define TILE_WIDTH 32
4 #define w (TILE_WIDTH + FILTER_WIDTH - 1)
5 #define clamp(x) (min(max((x), 0.0f), 1.0f))
6 __constant__ float device_filter[FILTER_WIDTH *
7     FILTER_WIDTH];

```

At the beginning, the variable that will be stored in the GPU's constant memory is also defined.

2.2.1 launch_convolution()

To launch the kernel on the GPU, we have defined a function that receives the image width, height, and channels and specifies the space on the device memory (*cudaMalloc*). It transfers the image data from the host to the device with the *cudaMemcpy()*. The filter has to be stored in the constant memory because it does not change during execution and is the same for every convolution operation, and for that reason is used *cudaMemcpyToSymbol()*.

In CUDA, when a kernel is called, the work is divided between multiple blocks of multiple threads. The dimension of the block defines how many threads operate in it. The grid dimension specifies how many blocks the GPU scheduler has to manage and assign to the streaming multiprocessors.

```

1 void launch_convolution(float *host_input, float
2     *host_filter, float *host_output, int
3     channels, int width, int height) {
4     float *device_input, *device_output;
5     size_t image_size = width * height * channels
6         * sizeof(float);
7     size_t filter_size = FILTER_WIDTH *
8         FILTER_WIDTH * sizeof(float);
9
10    cudaMalloc(&device_input, image_size);
11    cudaMalloc(&device_output, image_size);
12
13    cudaMemcpy(device_input, host_input,
14        image_size, cudaMemcpyHostToDevice);
15    cudaMemcpyToSymbol(device_filter, host_filter,
16        filter_size);
17
18    dim3 dimBlock(TILE_WIDTH, TILE_WIDTH);
19    dim3 dimGrid((width + TILE_WIDTH - 1) /
20        TILE_WIDTH, (height + TILE_WIDTH - 1) /
21        TILE_WIDTH);
22
23    std::chrono::time_point<std::chrono::
24        system_clock> start, end;
25    start = std::chrono::system_clock::now();
26    convolution<<<dimGrid, dimBlock>>>(&
27        device_input, device_output, channels, width,
28        height);
29
30    cudaDeviceSynchronize();
31    end = std::chrono::system_clock::now();
32    std::chrono::duration<double> milliseconds =
33        end - start;
34    std::cout << "Total convolution time: " <<
35        milliseconds.count() << " ms" << std::endl;
36
37    cudaError_t err = cudaGetLastError();
38    if (err != cudaSuccess) {
39        std::cerr << "CUDA error: " <<
40            cudaGetErrorString(err) << std::endl;
41    }
42 }

```

```

27 }
28
29 cudaMemcpy(host_output, device_output,
30            image_size, cudaMemcpyDeviceToHost);
31 cudaFree(device_input);
32 cudaFree(device_output);
33 }

```

The function transfers the data back to the host when the convolution is applied by each thread on the original image using *cudaMemcpy()*, and then frees the allocation on the GPU.

2.2.2 convolution()

The kernel called from the *launch_convolution()* contains the instruction that each thread of every block should execute. The first operation of the block is the instantiation of the variable that will memorize the values of the tile assigned to it in the shared memory. We have used the shared memory to reduce the time spent to access the global memory, since it is outside the multiprocessor, indeed, it requires more time to retrieve the data.

```

1 __shared__ float N_s[w][w];

```

Using the shared memory is a good choice not only for the access time reduction but also because we can load into it the values outside the tile's border that are needed to compute the convolution and are assigned and used from other blocks for its tile. To load every pixel value that we need, we have implemented a 2-batch loading technique.

The shared memory with size $(TILE_WIDTH + FILTER_WIDTH - 1)^2$ must be loaded by the threads with the values of the image inside the global memory before doing the convolution. Since the threads of a block are not able to fully load the shared memory, we need to let them bring more than just one value from the global memory. Fortunately, the number of uploads required for each thread, in our case, where $TILE_WIDTH \gg FILTER_WIDTH$, is at most two, and usually not for all of them. so for each thread we compute the first loading

```

1 int dest = threadIdx.y * TILE_WIDTH +
  threadIdx.x,
2   destY = dest / w, destX = dest % w,
3   srcY = blockIdx.y * TILE_WIDTH + destY
  - FILTER_RADIUS,
4   srcX = blockIdx.x * TILE_WIDTH + destX
  - FILTER_RADIUS,
5   src = (srcY * width + srcX) * channels
  + k;
6
7 if (srcY >= 0 && srcY < height && srcX >= 0
  && srcX < width)
8   N_s[destY][destX] = input[src];
9 else
10  N_s[destY][destX] = 0.0f;

```

and the second one

```

1 dest = threadIdx.y * TILE_WIDTH + threadIdx
  .x + TILE_WIDTH * TILE_WIDTH;
2 destY = dest / w; destX = dest % w;
3 srcY = blockIdx.y * TILE_WIDTH + destY -
  FILTER_RADIUS;
4 srcX = blockIdx.x * TILE_WIDTH + destX -
  FILTER_RADIUS;
5 src = (srcY * width + srcX) * channels + k;
6
7 if (destY < w) {
8   if (srcY >= 0 && srcY < height && srcX
  >= 0 && srcX < width)
9     N_s[destY][destX] = input[src];
10  else
11    N_s[destY][destX] = 0.0f;
12 }
13 __syncthreads();

```

After each thread of a block finishes uploading the pixel values, the convolution can be computed by applying the filter to the pixel assigned to the thread and its neighbours.

```

1 float accum = 0.0f;
2 for (int fy = 0; fy < FILTER_WIDTH; fy++)
3   for (int fx = 0; fx < FILTER_WIDTH; fx
  ++){
4     accum += N_s[threadIdx.y + fy][
  threadIdx.x + fx] * device_filter[fy *
  FILTER_WIDTH + fx];
5
6   int y = blockIdx.y * TILE_WIDTH + threadIdx
  .y;
7   int x = blockIdx.x * TILE_WIDTH + threadIdx
  .x;
8   if (y < height && x < width)
9     output[(y * width + x) * channels + k] =
  clamp(accum);
10  __syncthreads();

```

Before ending, the normalized result is stored in the global memory. When copied to the host memory, the value can be reconverted to the unsigned char type to obtain the new image.

3. Sperimental results

The execution times of the programs have been collected, considering only the time spent on the convolution without considering the time that we need to move the data from and to the host. We do this to focus our attention on the actual operation of convolution. Practically, when we consider a huge dataset of images, the data transmission cost is much less relevant. The results are shown in *Table 1*.

We can notice how the execution time required for the CUDA convolution is 424 ($Speedup = \frac{T_{seq}}{T_{par}} = \frac{0.399}{0.00094} = 424,052$) times less than the time spent on the sequential operation. The uploading of the filter into the constant cache and the block's tile into the shared memory is fast, thanks to the 2-batch loading and the leverage of constant memory. The data transfer for just one image is the bottleneck that slows down the execution of the entire CUDA program; the uploading to the GPU takes between 20 and 35 milliseconds, reducing the speedup obtained from the program.

In *Table 2*, we show the time the program takes to upload to the GPU and download the data from it. For just an image is quite burdensome, but especially in machine learning, where there are more images to work with, using batches to load them reduces the total data transfer cost and improves the performance.

mented several techniques that enhance execution on GPU architecture; using constant cache and shared memory reduces the multiple accesses of streaming multiprocessors to global memory. In a convolution context, 2-batch loading makes filling shared memory extremely beneficial for the operation. The upload cost for one image is the most computationally expensive part of the CUDA program and slows down the process. In future implementation, we should test the programs with more than one image, and perform other tests on the dimension of the block and the number of threads defined for them.

Version	t1	t2	t3	t4	t5	t6	mean
Sequential(ms)	0.433983	0.423375	0.380958	0.378117	0.375455	0.402566	0.399075667
Parallel(ms)	0.000800096	0.00078254	0.00118689	0.0011847	0.000970656	0.000721927	0.000941135

Table 1. Speedup results.

	t1	t2	t3	t4	t5	t6
Upload(ms)	0.325928	0.24965	0.24803	0.246484	0.248367	0.296436
Download(ms)	0.0398763	0.0399647	0.0405561	0.0400545	0.0398323	0.0401356

Table 2. GPU data tranfer.

4. Conclusion

We have demonstrated that CUDA convolution is a hundred times faster than convolution performed on a CPU. Additionally, we have imple-