

CUDA learning

Alessio Cimma

April 12, 2025

Contents

1	Project structure	2
2	CUDA architecture	3
2.1	Hierarchy of CUDA architecture	3
2.2	How to launch a CUDA kernel	3
2.3	Best parameters (block_size)	4
2.4	Best parameters (grid_size)	4
3	Exercise 1	5
4	Exercise 2	6
5	Exercise 3	7

1 Project structure

The program will contain a **Bridge** class, which will help call different functions during training. The idea is to have a comfortable environment to learn new things. I'll add in near future new tools to make debugging, plotting, benchmarking and comparisons easier. Each exercise will be contained in a different `self.execute_exerciseX()`. At the beginning of the program you will be asked to choose which CUDA kernel and function to load. Change it each time you want test different exercises. Possibly, each kernel will contain only one exercise, if the exercise needs more than one function they will all be contained in a single `.cu` file.

```
1 # Leggi il file contenente il kernel CUDA
2 with open(f'{kernel_name}.cu', 'r') as f:
3     kernel_code = f.read()
```

Listing 1: Loading CUDA kernel (kernel.cu)

2 CUDA architecture

2.1 Hierarchy of CUDA architecture

We have a grid, which determines the number of blocks, which are a collection of threads. A grid 2x2 contains 4 blocks. If each block is 32 threads, we will have 128 threads total.

To determine the optimal amount of grid and blocks, organize the number of blocks and then derive the size of the grid using the following subsections.

2.2 How to launch a CUDA kernel

Remember that when launching the kernel, the first 2 arguments need to be passed as tuple (grid_size and block_size), even if it's just one number, to do it use the syntax (x,):

```
1 # 1D version
2 self.imported_kernel((grid_size,), (block_size,), (x, y, z, n))
3
4 # 2D version
5 self.imported_kernel(grid_size, block_size, (x, y, z, n))
```

Listing 2: Launch CUDA kernel

2.3 Best parameters (block_size)

The following statistics should help choosing the best block-size to use in future projects, remember that you can't exceed the **MaxThreadsPerBlock** value (1024 on my laptop):

- Memory-bound kernels (loads/stores): 256-1024 \rightarrow *Hide latency via massive parallelism*
- Compute-bound kernels: 128-512 \rightarrow *More registers per thread may be needed*
- Shared memory usage per block: 128-256 \rightarrow *Avoid limiting number of resident blocks*
- Register-heavy kernels: 128 or lower \rightarrow *Prevent spilling and reduce pressure*
- Small data sizes ($n < 10k$): 64-256 \rightarrow *Larger blocks may cause underutilization*
- Thread divergence (if-else logic): 32-64 \rightarrow *Keep warps smaller to minimize divergence waste*

2.4 Best parameters (grid_size)

The best grid size you can use is the following, it maximizes the usage of the GPU by dividing the problem equally across all blocks:

```
1 # 1D version
2 grid_size = (n + block_size - 1) // block_size
3
4 # 2D version
5 grid_size = (
6     (W_delta + block_size[0] - 1) // block_size[0],
7     (W_delta + block_size[1] - 1) // block_size[1]
8 )
```

Listing 3: Determine grid size

3 Exercise 1

Objective: Creating a kernel that sums two 1D arrays element-wise.

What I learned:

- To access the correct index at which I am working, I need to use:

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

Listing 4: Index access

This way, I enter the correct block position multiplied by the size of the block (measured in threads) and then add the thread index of its parent block.

- I should the following command to cover the edge case where i have a potential thread that can do some work, but there's no more work to be found (due to the size of the work not being a power of two):

```
1 if (idx < n)
```

Listing 5: Index edge-case

4 Exercise 2

Objective: Creating a kernel that performs a convolution (morphological erosion).

What I learned:

- If the grid is 2D I can access the indices X and Y of the block and thread:

```
1 int i = blockIdx.y * blockDim.y + threadIdx.y;  
2 int j = blockIdx.x * blockDim.x + threadIdx.x;
```

Listing 6: Index access

- Now to check the new boundaries:

```
1 if (i >= Size_x || j >= Size_y) return;
```

Listing 7: Index edge-case

- I can obtain the value of infinity without importing header files using:

```
1 float min_val = 1.0f / 0.0f;
```

Listing 8: Infinity

5 Exercise 3

Objective: Trying to speeding up the execution time of the Exercise 2 using shared memory.

Execution Time	NO shared-memory	shared-memory
	75.3s	84.7s

Table 1: Performance comparison

What I learned: No need for shared memory for data that needs to be accessed from ALL threads just once. Useful when each thread needs to access a certain value multiple time during the execution. The best case is when an entire block needs a batch of data to use multiple times, otherwise the memory overhead and transfer time makes it slower than just accessing the data from global memory.

Cooperative-loading: Interesting way of loading memory using all threads:

```
1 for (int idx = local_thread_idx; idx < total_tile_size; idx +=
  threads_per_block) {
2   int ki = idx / TILE_KW;
3   int kj = idx % TILE_KW;
4
5   int global_ki = tile_offset_i + ki;
6   int global_kj = tile_offset_j + kj;
7
8   float val = (global_ki < KH && global_kj < KW)
9     ? ker[global_ki * KW + global_kj] : 0.0f;
10
11   ker_tile[ki * (TILE_KW + 1) + kj] = val;
12 }
```

Listing 9: Index edge-case