Getting started with CUDA Part 2 - Host view of GPU computation

Edwin Carlinet, Joseph Chazalon {firstname.lastname@lrde.epita.fr} April 2020

EPITA Research & Development Laboratory (LRDE)





Host view of GPU computation

You do not need to write kernels to run CUDA code: you can use kernels from a library, written by someone else.

This section is about how to properly launch CUDA kernels using their API only.

Sequential and parallel sections

We use the GPU(s) as co-processor(s).

Our program is made of a series of sequential and parallel sections.

Of course, CPU code can be multi-threaded too!

Sequential Execution			
Serial code	Host		
Parallel kernel	Device		
Kernel0<<<>>>()	Grid 0		
	Block (0, 0) Block (1, 0) Block (2, 0)		
	Block (0, 1) Block (1, 1) Block (2, 1)		
Serial code	Host		
	Device		
Parallel kernel Kernell<<<>>>()	Grid 1		
	Block (0, 0) Block (1, 0)		
	Block (0, 1) Block (1, 1)		
	Block (0, 2) Block (1, 2)		

Figure 1: Heterogeneous programming

We need to transfer inputs from the host to the device and outputs the other way around.



Figure 2: Computation on separate device

A proper kernel invocation

}

Let's fix this code!

```
#include <cuda.h>
void vecAdd(float *h_A, float *h_B, float *h_C, int n)
{
  int size = n * sizeof(float);
  float *d_A, *d_B, *d_C;
  // 1.1 Allocate device memory for A, B and C
  // 1.2 Copy A and B to device memory
 // 2. Launch kernel code - computation done on device
 // 3. Copy C (result) from device memory
  // Free device vectors
```

	1D	2D	3D
Allocate	cudaMalloc()	cudaMallocPitch()	cudaMalloc3D()
Сору	cudaMemcpy()	cudaMemcpy2D()	cudaMemcpy3D()
On-device init.	<pre>cudaMemset()</pre>	cudaMemset2D()	cudaMemset3D()
Reclaim		cudaFree()	

... plus many others detailed in the CUDA Runtime API documentation...

Why 2D and 3D variants?

- Strong alignment requirements in device memory
 - Enables correct loading of memory chunks to SM caches (correct bank alignment)
- Proper striding management in automated fashion

We just need the three following ones for now:

cudaError_t cudaMalloc (void** devPtr, size_t size_in_bytes)

Allocates space in the **device global** memory.

Asynchronous data transfer. cudaMemcpyKind pprox copy direction:

- cudaMemcpyHostToHost
- cudaMemcpyHostToDevice \leftarrow useful
- cudaMemcpyDeviceToHost \leftarrow useful
- cudaMemcpyDeviceToDevice
- cudaMemcpyDefault ← Direction inferred from pointer values. Requires unified virtual addressing.

cudaError_t cudaFree (void* devPtr)

```
#include <cuda.h>
void vecAdd(float *h A, float *h B, float *h C, int n)
{
  int size = n * sizeof(float);
  float *d A, *d B, *d C;
  // 1.1 Allocate device memory for A, B and C
  cudaMalloc((void **) &d_A, size); // TODO repeat for d_B and d_C
  // 1.2 Copy A and B to device memory
  cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
  cudaMemcpy(d_B, h_B, size, cudaMemcpyHostToDevice);
  // 2. Launch kernel code - computation done on device
  k_VecAdd<<<NB, NT>>>(d_A, d_B, d_C); // FIXME How to compute NB and NT?
  // 3. Copy C (result) from device memory
  cudaMemcpy(h C, d C, size, cudaMemcpyDeviceToHost);
  // Free device vectors
  cudaFree(d_A); // TODO repeat for d_B and d_C
}
```

```
In practice, we need to check for API errors
```

```
cudaError_t err = cudaMalloc((void **) &d_A, size);
if (err != cudaSuccess) {
    printf("%s in %s at line %d\n",
        cudaGetErrorString(err),
        __FILE__,
        __LINE__);
    exit(EXIT_FAILURE);
}
```

No: cudaMalloc(), cudaMemcpy() and cudaFree() shall be called from host only.

However, kernels may allocate, use and reclaim memory dynamically using regular malloc(), memset(), memcpy() and free() functions.

Note that if some device code allocates some memory, it must free it.

We want to fix this line:

```
k_VecAdd<<<NB, NT>>>(d_A, d_B, d_C);
```

Kernel invocation syntax:

kernel<<<blocks, threads_per_block, shmem, stream>>>(param1, param2, ...);

- blocks: number of blocks in the grid;
- threads_per_block: number of threads for each block;
- shmem: (opt.) amount of shared memory to allocate (in bytes);
- stream: (opt.) CUDA stream (not discussed in this course, see the documentation).

Lvl 0: Naive trial with as many threads as possible

k_VecAdd<<<1, size>>>(d_A, d_B, d_C);

LvI 0: Naive trial with as many threads as possible

k_VecAdd<<<1, size>>>(d_A, d_B, d_C);

Will fail with large vectors.

Hardware limitation on the maximum number of threads per block (1024 for Compute Capability 3.0-7.5).

Will fail with vectors of size which is not a multiple of warp size.

Lvl 1: It works with just enough blocks

```
// Get max threads per block
int devId = 0; // There may be more devices!
cudaDeviceProp deviceProp;
cudaGetDeviceProperties(&deviceProp, devId);
printf("Maximum grid dimensions: %d x %d x %d\n",
        deviceProp.maxGridSize[0],
        deviceProp.maxGridSize[1],
        deviceProp.maxGridSize[2]);
printf("Maximum block dimensions: %d x %d x %d\n",
        deviceProp.maxThreadsDim[0],
        deviceProp.maxThreadsDim[1],
        deviceProp.maxThreadsDim[2]);
// Compute the number of blocks
int xThreads = deviceProp.maxThreadsDim[0];
dim3 DimBlock(xThreads, 1, 1); // 1D VecAdd
int xBlocks = (int) ceil(n/xThreads);
dim3 DimGrid(xBlocks, 1, 1);
// Launch the kernel
k VecAdd<<<DimGrid, DimBlock>>>(d A, d B, d C);
```

Lvl 2: Tune block size given kernel requirements and hardware constraints It is important to understand the difference between:

the logical decomposition of your program:

 $\mathsf{problem}\approx\mathsf{grid}\to\mathsf{blocks}\to\mathsf{threads}$

- the scheduling of the computation on the hardware:
 - assignment of each block to a Streaming Multiprocessor (SM)
 - groups threads into warps
 - run groups of warps concurrently

The hardware constraints are different between each *Compute Capability* version. See the CUDA C programming manual, Appendix H for details about each hardware version. In particular, the amount of memory available on each SM may limit the number of threads one would actually want to launch...

But this depends on the kernel code!

The **CUDA Occupancy Calculator** APIs are designed to assist programmers in choosing the best number of threads per block based on register and shared memory requirements of a given kernel.

```
#include <stdio.h>
__global__ void print_kernel() {
    printf("Hello!\n");
}
```

This code prints nothing!

```
int main() {
    print_kernel<<<1, 1>>>();
```

}

16

```
#include <stdio.h>
__global__ void print_kernel() {
    printf("Hello!\n");
}
```

```
int main() {
    print_kernel<<<1, 1>>>();
    cudaDeviceSynchronize();
}
```

Host code synchronization requires cudaDeviceSynchronize() because **kernel invocation is asynchronous** from host perspective.

On the device, kernel invocations are strictly sequential (unless you schedule them on different *streams*).

Yes: This is the basis of *dynamic parallelism*.

Some restrictions over the stack size apply.

Remember that the device runtime is a functional subset of the host runtime, ie you can perform device management, kernel launching, device memcpy, etc., but with some restrictions (see the documentation for details).

The compiler may inline some of those calls, though.

A host-only view of the computation is sufficient for most of the cases:

- 1. upload input data to the device
- 2. fire a kernel
- 3. download output data from the device

Advanced CUDA requires to make sure we saturate the SMs, and may imply some kernel study to determine the best:

- amount of threads per blocks
- amount of blocks per grid
- work per thread (if applicable)
- • • •

This depends on:

- hardware specifications: maximum gridDim and blockDim, etc.
- kernel code: amount of register and shared memory used by each thread