

PMC Lecture 17

Gianmaria Romano

10 December 2025

Compute Unified Device Architecture

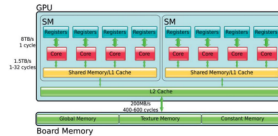
Remember that, when performing memory allocation operations, data allocated on host memory cannot be seen from the GPU, and viceversa, therefore requiring explicit data transfers during communications.

Memory Copy Kind

The `cudaMemcpyKind` parameter of `cudaMemcpy` is also an enumerated type. The `kind` parameter can take one of the following values:

- `cudaMemcpyHostToHost` = 0, Host to Host
- `cudaMemcpyHostToDevice` = 1, Host to Device
- `cudaMemcpyDeviceToHost` = 2, Device to Host
- `cudaMemcpyDeviceToDevice` = 3, Device to Device (for multi-GPU configurations)
- `cudaMemcpyDefault` = 4, used when Unified Virtual Address space capability is available (see [Section 6.7](#))

Memory Types



Variable declaration	Memory	Scope	Lifetime
Automatic variables other than arrays	Register	Thread	Kernel
Automatic array variables	Global	Thread	Kernel
<code>__device__ __shared__ int SharedVar;</code>	Shared	Block	Kernel
<code>__device__ int GlobalVar;</code>	Global	Grid	Application
<code>__device__ __constant__ int ConstVar;</code>	Constant	Grid	Application

6.1.1 The global memory

The global memory is an off-chip memory component that can store large amounts of data but is physically distant from the processors.

For this reason, accessing data from the global memory requires higher latency compared to on-chip memory components.

6.1.2 Registers

Similarly to CPUs, GPU registers are used to store a thread's local variables. Generally speaking, the available registers on a streaming multiprocessor are evenly distributed among all resident threads, but, whenever the register demand of a kernel exceeds the number of available registers per thread, some local variables will be spilled to the off-chip global memory, which, however, is much slower to access.

Since register usage affects how many threads can be active on a streaming multiprocessor, it will also influence the number of resident warps that can be supported by the streaming multiprocessor.

For this reason, it is possible to define occupancy as the ratio between the number of resident warps and the maximum number of warps that can be supported by the streaming multiprocessor.

Ideally, occupancy should be close to 1 so that the hardware has enough resources to effectively hide latencies.

6.1.3 On-chip shared memory

A shared memory is an on-chip memory component that is shared among threads within a block.

Despite having smaller capacity, the shared memory has a much lower latency compared to other (off-chip) components, making it useful for storing frequently accessed data that, otherwise, would take much longer to fetch.

Example: The following code provides an implementation of a stencil, which is a geometric tool that is often used for solving mathematical problems:

```
1  #include <stdio.h>
2  #include <cuda.h>
3
4  __global__ void stencil_1d(int *in, int *out) {
5      __shared__ int temp[BLOCK_SIZE + 2 * RADIUS]; // RADIUS adds a padding on both sides.
6      int gindex = threadIdx.x + blockIdx.x * blockDim.x; // Global index in the grid.
7      int lindex = threadIdx.x + radius; // Local index in the shared memory.
8
9      // Read input elements into the shared memory.
10     temp[lindex] = in[gindex];
11     if (threadIdx.x < RADIUS) {
12         temp[lindex - RADIUS] = in[gindex - RADIUS];
13         temp[lindex + BLOCK_SIZE] = in[gindex + BLOCK_SIZE];
14     }
15
16     __syncthreads(); // Synchronize all threads within a block to avoid dependencies.
17
18     // Apply the stencil to perform computations.
19     int result = 0;
20     for (int offset = -RADIUS; offset <= RADIUS; offset++) {
21         result += temp[lindex + offset];
22     }
23
24     // Store the result.
25     out[gindex] = result;
26 }
```

6.1.4 The constant memory

A constant memory is a read-only memory component that is used to store constant values that do not change during execution.

In particular, the constant memory provides performance benefits because it is cached and supports broadcasting when threads within a warp need to access the same memory location.

```
__constant__ type variable_name; // static
cudaMemcpyToSymbol(variable_name, host_ptr, sizeof(type), cudaMemcpyHostToDevice);
// warning: cannot be dynamically allocated
```

- data will reside in the constant memory address space
- has static storage duration (persists until the application ends)
- readable from all threads of a kernel

6.2 Performance estimation for GPUs

Generally speaking, GPU performance is measured in terms of the number of floating-point operations it can perform in one second.

In practice, however, the performance of a program is limited by the data transfer rate, causing the application to be "memory-bound".

For this reason, it can come in handy to define the compute-to-global-memory-access ratio, which measures the number of floating-point operations that are executed per global memory access within a certain region of the program.