

## Summary

Now that the veil has been lifted, we can see that GPUs are really just multi-threaded SIMD processors, although they have more processors, more lanes per processor, and more multithreading hardware than do traditional multicore computers. For example, the Fermi GTX 480 has 15 SIMD processors with 16 lanes per processor and hardware support for 32 SIMD threads. Fermi even embraces instruction-level parallelism by issuing instructions from two SIMD threads to two sets of SIMD lanes. They also have less cache memory—Fermi's L2 cache is 0.75 megabyte—and it is not coherent with the distant scalar processor.

| Type                 | More descriptive name used in this book | Official CUDA/NVIDIA term | Book definition and AMD and OpenCL terms   | Official CUDA/NVIDIA definition  |
|----------------------|---|---------------------------|--|--|
| Program abstractions | Vectorizable loop                       | Grid                      | A vectorizable loop, executed on the GPU, made up of one or more “Thread Blocks” (or bodies of vectorized loop) that can execute in parallel. OpenCL name is “index range.” AMD name is “NDRange”.           | A grid is an array of thread blocks that can execute concurrently, sequentially, or a mixture.   |
|                      | Body of Vectorized loop                 | Thread Block              | A vectorized loop executed on a multithreaded SIMD Processor, made up of one or more threads of SIMD instructions. These SIMD Threads can communicate via Local Memory. AMD and OpenCL name is “work group”. | A thread block is an array of CUDA Threads that execute concurrently together and can cooperate and communicate via Shared Memory and barrier synchronization. A Thread Block has a Thread Block ID within its Grid. |
|                      | Sequence of SIMD Lane operations        | CUDA Thread               | A vertical cut of a thread of SIMD instructions corresponding to one element executed by one SIMD Lane. Result is stored depending on mask. AMD and OpenCL call a CUDA Thread a “work item.”                 | A CUDA Thread is a lightweight thread that executes a sequential program and can cooperate with other CUDA Threads executing in the same Thread Block. A CUDA Thread has a thread ID within its Thread Block.        |
| Machine object       | A Thread of SIMD instructions           | Warp                      | A traditional thread, but it contains just SIMD instructions that are executed on a multithreaded SIMD Processor. Results are stored depending on a per-element mask. AMD name is “wavefront.”               | A warp is a set of parallel CUDA Threads (e.g., 32) that execute the same instruction together in a multithreaded SIMT/SIMD Processor.   |
|                      | SIMD instruction                        | PTX instruction           | A single SIMD instruction executed across the SIMD Lanes. AMD name is “AMDIL” or “FSAIL” instruction.  | A PTX instruction specifies an instruction executed by a CUDA Thread.  |

**Figure 4.24** Conversion from terms used in this chapter to official NVIDIA/CUDA and AMD jargon. OpenCL names are given in the book definition.

| Type                | More descriptive name used in this book | Official CUDA/NVIDIA term | Book definition and AMD and OpenCL terms  | Official CUDA/NVIDIA definition  |
|---------------------|---|---------------------------|---|--|
| Processing hardware | Multithreaded SIMD processor            | Streaming multi-processor | Multithreaded SIMD Processor that executes thread of SIMD instructions, independent of other SIMD Processors. Both AMD and OpenCL call it a “compute unit.” However, the CUDA Programmer writes program for one lane rather than for a “vector” of multiple SIMD Lanes. | A streaming multiprocessor (SM) is a multithreaded SIMT/ SIMD Processor that executes warps of CUDA Threads. A SIMT program specifies the execution of one CUDA Thread, rather than a vector of multiple SIMD Lanes. |
|                     | Thread block scheduler                  | Giga thread engine        | Assigns multiple bodies of vectorized loop to multithreaded SIMD Processors. AMD name is “Ultra-Threaded Dispatch Engine”.  | Distributes and schedules thread blocks of a grid to streaming multiprocessors as resources become available.  |
|                     | SIMD Thread scheduler                   | Warp scheduler            | Hardware unit that schedules and issues threads of SIMD instructions when they are ready to execute; includes a scoreboard to track SIMD Thread execution. AMD name is “Work Group Scheduler”.  | A warp scheduler in a streaming multiprocessor schedules warps for execution when their next instruction is ready to execute.  |
|                     | SIMD Lane                               | Thread processor          | Hardware SIMD Lane that executes the operations in a thread of SIMD instructions on a single element. Results are stored depending on mask. OpenCL calls it a “processing element.” AMD name is also “SIMD Lane”.   | A thread processor is a datapath and register file portion of a streaming multiprocessor that executes operations for one or more lanes of a warp.   |
| Memory hardware     | GPU Memory                              | Global Memory             | DRAM memory accessible by all multithreaded SIMD Processors in a GPU. OpenCL calls it “Global Memory.”  | Global memory is accessible by all CUDA Threads in any thread block in any grid; implemented as a region of DRAM, and may be cached.   |
|                     | Private Memory                          | Local Memory              | Portion of DRAM memory private to each SIMD Lane. Both AMD and OpenCL call it “Private Memory.”   | Private “thread-local” memory for a CUDA Thread; implemented as a cached region of DRAM.   |
|                     | Local Memory                            | Shared Memory             | Fast local SRAM for one multithreaded SIMD Processor, unavailable to other SIMD Processors. OpenCL calls it “Local Memory.” AMD calls it “Group Memory”.  | Fast SRAM memory shared by the CUDA Threads composing a thread block, and private to that thread block. Used for communication among CUDA Threads in a thread block at barrier synchronization points.               |
|                     | SIMD Lane registers                     | Registers                 | Registers in a single SIMD Lane allocated across body of vectorized loop. AMD also calls them “Registers”.  | Private registers for a CUDA Thread; implemented as multithreaded register file for certain lanes of several warps for each thread processor.  |

**Figure 4.25** Conversion from terms used in this chapter to official NVIDIA/CUDA and AMD jargon. Note that our descriptive terms “Local Memory” and “Private Memory” use the OpenCL terminology. NVIDIA uses SIMT, single-instruction multiple-thread, rather than SIMD, to describe a streaming multiprocessor. SIMT is preferred over SIMD because the per-thread branching and control flow are unlike any SIMD machine.