

## Compute terminology translation guide

CPU (with an x86 bias)	NVidia CUDA	Khronos (OpenCL, SYCL, Vulkan...)
SIMD lane	CUDA core	Processing Element
SIMD vector	Warp	Subgroup
Simultaneous Multi-Threading / Hyper-Threading	Hardware multithreading	<i>Not named, hidden inside PE/CU model</i>
Advanced Matrix eXtensions, Neural Processing Unit	Tensor core	<i>No standard name (only usable via Vulkan vendor extensions)</i>
<i>No equivalent, ray tracing is not special</i>	RTX core	<i>Not named, exposed via VK_KHR extensions</i>
Intel Data Streaming Accelerator	Tensor Memory Accelerator (kind of [1])	<i>Not named, exposed via OpenCL async copy API</i>
Core	Streaming Multiprocessor	Compute Unit
Core complex / Sub-NUMA cluster	Thread block cluster (kind of [2])	<i>Not exposed yet</i>
NUMA node	Device	Device
SIMD registers, with manual [3] spill to caches/RAM	Local memory, with automatic spill to global memory	Private memory, with automatic spill to global memory
L1/L2 cache, with automatic spill to lower caches/RAM	Shared memory	Local memory
<i>No equivalent, all memory is writable</i>	Constant memory	Constant memory
<i>No equivalent, images are not special</i>	Texture memory	<i>Not named, exposed via image API</i>
Cache hierarchy + RAM	Global memory	Global memory
SIMD lane	Threads within the same warp	Work-items within the same subgroup
SIMD vector	Warp	Subgroup
Thread	Threads within the same thread block	Work-items within the same work-group
<i>No equivalent, all threads can synchronize with each other</i>	Thread block	Work-group
<i>No equivalent, all threads can synchronize with each other</i>	Thread block cluster	<i>Not exposed yet</i>
<i>No equivalent, threads are independent from each other</i>	Grid	NDrange
<i>No equivalent, can spawn threads anytime</i>	Stream	Command queue
<i>No equivalent, can spawn threads anytime</i>	Graph	Command buffer ( <i>only exposed in Vulkan</i> )
Buffer	Buffer	Buffer
<i>No equivalent, images are not special</i>	Texture	Image
<i>No equivalent, images are not special</i>	Texture	Sampler
RAM	Memory	Device memory
<i>No equivalent, all useful memory is reachable</i>	Unified memory	Shared memory
<i>No equivalent, all useful memory is reachable</i>	Host/pinned memory	Host memory

[1] Both TMA and DSA allow offloading memory copy work from the compute cores, but Intel DSA is more focused on device DMA and Nvidia TMA on VRAM-cache copies.

[2] On a hardware level, it's the same idea, but CUDA only allows some forms of synchronization when threads belong to the same thread block cluster.

[3] Speaking from the perspective of machine code here. Of course, the compiler of most programming languages will automate it for you.