

CUDA C PROGRAMMING QUICK REFERENCE¹

Function Qualifiers

<code>--global--</code>	called from host, executed on device
<code>--device--</code>	called from device, executed on device
<code>--host--</code>	called from host, executed on host
<code>--host-- --device--</code>	generates code for host and device
<code>--noinline--</code>	if possible, do not inline
<code>--forceinline--</code>	force compiler to inline

Variable Qualifiers (Device)

<code>--device--</code>	variable on device (Global Memory)
<code>--constant--</code>	variable in Constant Memory
<code>--shared--</code>	variable in Shared Memory
<code>--restrict--</code>	restricted pointers, assert to the compiler that pointers are not aliased (cf. aliased pointer)
– No Qualifier –	automatic variable, resides in Register or in Local Memory in some cases (local arrays, register spilling)

Built-in Variables (Device)

<code>dim3 gridDim</code>	dimensions of the current grid (<code>gridDim.x, ...</code>) (composed of independent blocks)
<code>dim3 blockDim</code>	dimensions of the current block (composed of threads) (total number of threads should be a multiple of warp size)
<code>dim3 blockIdx</code>	block location in the grid (<code>blockIdx.x, ...</code>)
<code>dim3 threadIdx</code>	thread location in the block (<code>threadIdx.x, ...</code>)

Shared Memory

Static allocation	<code>--shared-- int a[128]</code>
Dynamic allocation (at kernel launch)	<code>extern --shared-- float b[]</code>

Host / Device Memory

Allocate pinned / page-locked Memory on host	<code>cudaMallocHost(&dptr, size)</code> (for higher bandwidth, may degrade system performance)
Allocate Device Memory	<code>cudaMalloc(&devptr, size)</code>
Free Device Memory	<code>cudaFree(devptr)</code>
Transfer Memory	<code>cudaMemcpy(dst, src, size, cudaMemcpyKind kind)</code> kind = { <code>cudaMemcpyHostToDevice, ...</code> }
Nonblocking Transfer	<code>cudaMemcpyAsync(dst, src, size, kind[, stream])</code> (host memory must be page-locked)
Copy to constant or global memory	<code>cudaMemcpyToSymbol(symbol, src, size[, offset[, kind]])</code> kind= <code>cudaMemcpy[HostToDevice DeviceToDevice]</code>

Synchronizing

Synchronizing one Block	<code>--syncthreads()</code> (device call)
Synchronizing all Blocks	<code>cudaDeviceSynchronize()</code> (host call, CUDA Runtime API)

Kernel

Kernel Launch	<code>kernel<<<blocks, threadsPerBlock[, smem_size[, stream]]>>>(<.>)</code>
---------------	--

CUDA Device Management

Init device (context)	<code>cudaSetDevice(devID)</code>
Reset current device	<code>cudaDeviceReset()</code> (for profiler and flushing)

CUDA Runtime API Error Handling

CUDA Runtime API error as String	<code>cudaGetErrorString(cudaError_t err)</code>
Last CUDA error produced by any of the runtime calls	<code>cudaGetLastError()</code>

OpenGL Interoperability

Init device (within OpenGL context)	<code>cudaGLSetGLDevice(devID)</code> (mutually exclusive to <code>cudaSetDevice()</code>)
Register buffer object (must not be bound by OpenGL)	<code>cudaGraphicsGLRegisterBuffer(&res, id, flags)</code> res: <code>cudaGraphicsResource</code> pointer id: OpenGL Buffer Id flags: register flags (read/write access)
Register texture or render buffer	<code>cudaGraphicsGLRegisterImage(&res, id, target, flags)</code>

Graphics Interoperability

Unregister graphics resource	<code>cudaGraphicsUnregisterResource(res)</code>
Map graphics resources for access by CUDA	<code>cudaGraphicsMapResources(count, &res[, stream])</code>
Get device pointer (access a mapped graphics resource) (OpenGL: buffer object)	<code>cudaGraphicsResourceGetMappedPointer(&dptr, size, res)</code>
Get CUDA array of a mapped graphics resource (OpenGL: texture or renderbuffer)	<code>cudaGraphicsSubResourceGetMappedArray(&a, res, i, lvl)</code>
Unmap graphics resource	<code>cudaGraphicsUnmapResources(count, &res[, stream])</code>

CUDA Texture

Textures are read-only global memory, but cached on-chip, with texture interpolation

Declare texture (at file scope)	<code>texture<DataType, TexType, Mode> texRef</code>
Create channel descriptor	<code>cudaCreateChannelDesc<DataType>()</code>
Bind memory to texture	<code>cudaBindTexture(offset, texref, dptr, channelDesc, size)</code>
Unbind texture	<code>cudaUnbindTexture(texRef)</code>
Fetch Texel (texture pixel)	<code>tex1D(texRef, x)</code> <code>tex2D(texRef, x, y)</code> <code>tex3D(texRef, x, y, z)</code> <code>tex1DLayered(texRef, x, layer)</code> <code>tex2DLayered(texRef, x, y, layer)</code>

CUDA Streams (Concurrency Management)

Stream = instruction sequence. Streams may execute their commands out of order.

Create CUDA Stream	<code>cudaStreamCreate(cudaStream_t &stream)</code>
Destroy CUDA Stream	<code>cudaStreamDestroy(stream)</code>
Synchronize Stream	<code>cudaStreamSynchronize(stream)</code>
Stream completed?	<code>cudaStreamQuery(stream)</code>

¹November 12, 2019. Cf. Complete Reference: "NVIDIA CUDA C Programming Guide" – Note, that functions may have optional arguments not listed here

	TESLA				FERMI		KEPLER			MAXWELL	
COMPUTE CAPABILITY	1.0	1.1	1.2	1.3	2.0	2.1	3.0	3.5	3.7	5.0	5.2
Max. dimensionality of grid		2						3			
Max. dimensionality of block							3				
Max. x-,y- or z-dimension of a grid		$2^{16}-1$						$2^{32}-1$			
Max. x- or y-dimension of a block		512						1024			
Max. z-dimension of a block							64				
Max. threads per block		512						1024			
Warp Size							32				
Max. resident blocks per SM			8					16		32	
Max. resident warps per SM	24		32		48				64		
Max. resident threads per SM	768		1024		1536				2048		
Number of 32-bit registers per SM	8 K		16 K		32 K		64 K		128 K		64 K
Max. registers per thread		124			63				255		
Max. shared memory per SM (≥ 2.0 : configurable L1 Cache)		16 KB			48 KB				112 KB	64 KB	96 KB
Number of shared memory banks		16						32			
Local memory per thread		16 KB						512 KB			
Constant memory size					64 KB						
Cache working set per SM for constant					8 KB					10 KB	
Cache working set per SM for texture		6-8 KB			12 KB			12 KB-48 KB		24 KB	48 KB
Max. instructions per kernel		2 million						512 million			
Max. width for 1D texture (array)		8192						65 536			
Max. width 1D texture (linear)							2^{27}				
Max. width×layers for 1D texture		8192×512						16 384×2048			
Max. textures bound to kernel			128						256		
Max. width×layers for 1D surface		N/A						65 536×2048			
Max. surfaces bound to kernel		N/A			8				16		
ARCHITECTURE SPECIFICATIONS											
Number of cores (with FPU and ALU)		8			32	48		192		128	
Number of special function units		2			4	8			32		
Number of texture units		2			4	8		16		8	
Number of warp schedulers		1			2				4		
Number of instructions issued by scheduler			1						2		
COMPUTE CAPABILITY	1.0	1.1	1.2	1.3	2.0	2.1	3.0	3.5	3.7	5.0	5.2

CC	GPU _s	Features
FERMI		
2.0	GF100, GF110	ECC, Better Caches (L1 and L2), dual warp scheduler, concurrent kernel execution, better atomics, int64 shared memory atomics, float32-atomicAdd, unified address space, ballot, thread-fence, surface, FMA, int32 ALU
2.1	GF104, ...	—
KEPLER (Focus: perf/watt, doubles, HPC)		
3.0	GK104, GK106, GK107	Polymorph Engine 2.0, GPU Boost, TXAA, warp shuffle, bindless textures, h.264 encoder NVENC, adaptive VSync, PCIe 3.0
3.5	GK110, GK208	Dynamic Parallelism, Hyper-Q, Grid Management Unit, GPUDirect (RDMA), funnel shift
3.7	GK210 (K80)	
MAXWELL (Focus: perf/watt, perf/area, single-precision, Gaming)		
5.0	GM107, GM108	
5.2	GM200, GM204, GM206	Polymorph Engine 3.0, VXGI (Global Illumination), H.265 encoding

GPU	GTX 580	GTX 680	GTX 780	GTX 980
Launch	Nov 2010	Mar 2012	May 2013	Sep 2014
Model	GF110	GK104	GK110	GM204
Core Clock (MHz)	772	1006	863	1126
Shader Clock (Mhz)	1544	–	–	–
Boost Clock (MHz)	–	1058	900	1216
PCIe Bus Support	2.0	3.0	3.0	3.0
CUDA Cores	512	1536	2304	2048
Memory Bandwidth (GB/sec)	192.4	192.2	288.4	224
Memory Clock (Mhz)	4008	6008	6008	7010
Memory Interface Width (bit)	384	256	384	256
Standard Memory Config MiB	1536	2048	3072	4096
Texture Fill Rate (billion/sec)	49.4	128.8	160.5	144
Max Temp. °C	97	98	95	98
Max Power W	244	195	250	165
SM count	16	8	12	16
Transistors	3×10^9	3.5×10^9	7×10^9	5.2×10^9
GFLOPS/s (SP/DP)	1581 / 198	3090 / 128	3977 / 166	4612 / 144
Fab-Size nm	40	28	28	28
Die Size mm ²	520	294	561	398
L2 Cache Size KiB	768	512	1536	2048
ROPs	48	32	48	64
Texture Units	64	128	192	128

Memory	Location	Cached	Access	Scope	Lifetime
Register	On-chip	N/A	R/W	Thread	Thread
Local	Off-chip	Yes	R/W	Thread	Thread
Shared	On-chip	N/A	R/W	Block	Block
Global	Off-chip	Yes	R/W	Global	Application
Constant	Off-chip	Yes	R	Global	Application
Texture	Off-chip	Yes	R	Global	Application
Surface	Off-chip	Yes	R/W	Global	Application