



PDC Summer School 2016

CUDA developer tools

2015-08-19

Michael Schliephake

Szilárd Páll

KTH – CSC – HPCViz



Compiler

- Compiler toolchain
- Host – device compilation
- Device code generation
 - Multi-pass
 - Allows fallback code-path
 - Embed multiple code versions (virtual vs physical architecture)



Memory checker

- Cuda-memcheck
- Checks for:
 - Misaligned/out of bounds access
 - Device malloc()/free()
 - API errors, hardware errors
 - Memory leaks
 - races



Debugger

- Cuda-gdb
- Nsight: VS or Eclipse plugin
- Parallel debuggers:
 - Totalview
 - Allinea DDT

Profilers

- nvprof
 - Command line tool for device profiling
 - API tracing
 - Can generate data for nvvp
- Visual Profiler (nvvp)
 - Linux/OS X
- Nsight (VS)
 - Visual Studio plugin
- 3rd party tools: TAU, Vampirtrace, PAPI



GPU management & monitoring

- nvidia-smi
 - Monitor GPU at runtime
 - Manage behavior (e.g. clocks, ECC)
- NVML
 - Programmatically do the same (and more)



CUDA Occupancy Calculator

http://developer.download.nvidia.com/compute/cuda/CUDA_Occupancy_calculator.xls

	A	B	C	D	E	F	G	H	I	J	K	L
1	CUDA GPU Occupancy Calculator											
2												
3												
4	Just follow steps 1, 2, and 3 below! (or click here for help)											
5												
6	1.) Select Compute Capability (click):	2,0	(Help)									
7	1.b) Select Shared Memory Size Config (bytes)	49152										
8												
9	2.) Enter your resource usage:											
10	Threads Per Block	256	(Help)									
11	Registers Per Thread	16										
12	Shared Memory Per Block (bytes)	4096										
13												
14	(Don't edit anything below this line)											
15												
16	3.) GPU Occupancy Data is displayed here and in the graphs:											
17	Active Threads per Multiprocessor	1536	(Help)									
18	Active Warps per Multiprocessor	48										
19	Active Thread Blocks per Multiprocessor	6										
20	Occupancy of each Multiprocessor	100%										
21												
22												
23	Physical Limits for GPU Compute Capability:	2,0										
24	Threads per Warp	32										
25	Warps per Multiprocessor	48										
26	Threads per Multiprocessor	1536										
27	Thread Blocks per Multiprocessor	8										

Click Here for detailed instructions on how to use this occupancy calculator.
 For more information on NVIDIA CUDA, visit <http://developer.nvidia.com/cuda>

Your chosen resource usage is indicated by the red triangle on the graphs. The other data points represent the range of possible block sizes, register counts, and shared memory allocation.

Impact of Varying Block Size

Threads Per Block	Multiprocessor or Warp Occupancy (# warps)
64	8
128	24
192	48
256	40 (My Block Size 256)
320	44
384	40
448	44
512	48
576	32
640	36
704	40
768	44
832	24
896	28
960	32
1024	36