

CUDA Occupancy Calculator

Compute Capability version

Threads per block

Registers per thread

Shared memory per block

GPU Occupancy Data is displayed here and in the graphs

Active Threads per Multiprocessor	2048
Active Warps per Multiprocessor	64
Active Thread Blocks per Multiprocessor	8
Occupancy of each Multiprocessor	1

Physical Limits for GPU Compute Capability

Version	6.1
Threads per Warp	32
Warps per Multiprocessor	64
Threads per Multiprocessor	2048
Thread Blocks per Multiprocessor	32
Total # of 32-bit registers per Multiprocessor	65536
Register allocation unit size	256
Register allocation granularity	warp
Max registers per Block	65536
Max registers per thread	255
Shared Memory per Multiprocessor (bytes)	98304
Shared Memory Allocation unit size	256
Warp allocation granularity (for register allocation)	4
Max thread block size	1024

Allocation Per Thread Block

Warps	8
Registers	4096
Shared Memory	1024

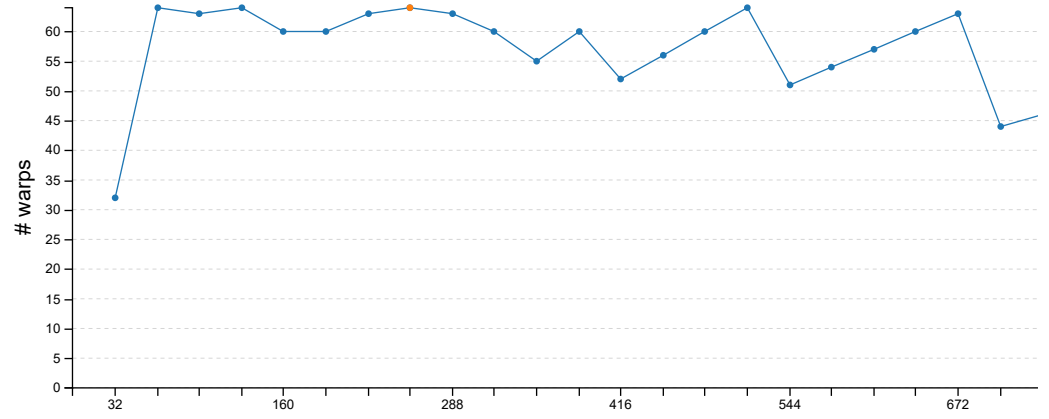
Maximum Thread Blocks Per Multiprocessor

Limited by Max Warps / Blocks per Multiprocessor	8
Limited by Registers per Multiprocessor	16

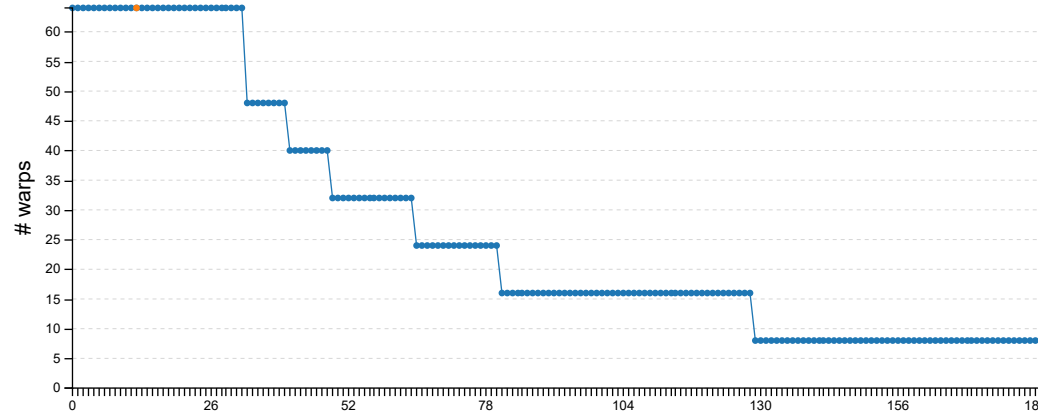
Limited by Shared Memory per Multiprocessor

96

Impact of Varying Block Size



Impact of Varying Register Count Per Thread



Impact of Varying Shared Memory Usage Per Block

