

	Result	Time	Cycles	Regs	GPU	SM Frequency	CC	Process		
Current	516 - ...	3,36 usecond	4.981	16	0 - NVIDIA GeForce RTX 3070	1,48 cycle/nsecond	8.6	[4796] test_4_256.out		

GPU Speed Of Light Throughput

All

High-level overview of the throughput for compute and memory resources of the GPU. For each unit, the throughput reports the achieved percentage of utilization with respect to the theoretical maximum. Breakdowns show the throughput for each individual sub-metric of Compute and Memory to clearly identify the highest contributor.

Compute (SM) Throughput [%]	0,16	Duration [usecond]	3,36
Memory Throughput [%]	1,19	Elapsed Cycles [cycle]	4.981
L1/TEX Cache Throughput [%]	0,72	SM Active Cycles [cycle]	1.368,46
L2 Cache Throughput [%]	1,03	SM Frequency [cycle/nsecond]	1,48
DRAM Throughput [%]	1,19	DRAM Frequency [cycle/nsecond]	6,70

Small Grid This kernel grid is too small to fill the available resources on this device, resulting in only 0.0 full waves across all SMs. Look at [Launch Statistics](#) for more details.

GPU Throughput



Compute Throughput Breakdown

SM: Inst Executed Pipe Lsu [%]	0,16
SM: Issue Active [%]	0,09
SM: Inst Executed [%]	0,08
SM: Mio Inst Issued [%]	0,06
SM: Mio2rf Writeback Active [%]	0,05
SM: Pipe Fmaheavy Cycles Active [%]	0,03
SM: Pipe Alu Cycles Active [%]	0,03
SM: Inst Executed Pipe Uniform [%]	0,03
SM: Inst Executed Pipe Adu [%]	0,02
SM: Pipe Fma Cycles Active [%]	0,02
SM: Inst Executed Pipe Cbu Pred On Any [%]	0,01
SM: Mio Pq Read Cycles Active [%]	0,00
SM: Mio Pq Write Cycles Active [%]	0,00
SM: Inst Executed Pipe Ipa [%]	0
SM: Inst Executed Pipe Tex [%]	0
SM: Inst Executed Pipe Xu [%]	0
IDC: Request Cycles Active [%]	0
SM: Pipe Fp64 Cycles Active [%]	0
SM: Pipe Tensor Cycles Active [%]	0

Memory Throughput Breakdown

DRAM: Cycles Active [%]	1,19
L2: T Sectors [%]	1,03
L2: Lts2xbar Cycles Active [%]	0,89
DRAM: Dram Sectors [%]	0,71
L2: T Tag Requests [%]	0,37
L2: Xbar2lts Cycles Active [%]	0,36
L2: D Sectors Fill Device [%]	0,34
L2: D Sectors [%]	0,25
L1: M Xbar2l1tex Read Sectors [%]	0,20
L1: Texin Sm2tex Req Cycles Active [%]	0,18
L1: Lsuin Requests [%]	0,16
L1: Data Pipe Lsu Wavefronts [%]	0,13
L1: M L1tex2xbar Req Cycles Active [%]	0,10
L1: Lsu Writeback Active [%]	0,09
L1: Data Bank Writes [%]	0,03
L1: Data Bank Reads [%]	0,02
L2: D Sectors Fill Sysmem [%]	0,00
L1: F Wavefronts [%]	0
L1: Tex Writeback Active [%]	0
L2: D Atomic Input Cycles Active [%]	0
L1: Data Pipe Tex Wavefronts [%]	0

Launch Statistics

Summary of the configuration used to launch the kernel. The launch configuration defines the size of the kernel grid, the division of the grid into blocks, and the GPU resources needed to execute the kernel. Choosing an efficient launch configuration maximizes device utilization.

Grid Size	25	Function Cache Configuration	cudaFuncCachePreferNone
Registers Per Thread [register/thread]	16	Static Shared Memory Per Block [byte/block]	0
Block Size	32	Dynamic Shared Memory Per Block [byte/block]	0
Threads [thread]	800	Driver Shared Memory Per Block [Kbyte/block]	1,02
Waves Per SM	0,03	Shared Memory Configuration Size [Kbyte]	16,38

Small Grid The grid for this launch is configured to execute only 25 blocks, which is less than the GPU's 46 multiprocessors. This can underutilize some multiprocessors. If you do not intend to execute this kernel concurrently with other workloads, consider reducing the block size to have at least one block per multiprocessor or increase the size of the grid to fully utilize the available hardware resources. See the [Hardware Model](#) description for more details on launch configurations.

Occupancy

Occupancy is the ratio of the number of active warps per multiprocessor to the maximum number of possible active warps. Another way to view occupancy is the percentage of the hardware's ability to process warps that is actively in use. Higher occupancy does not always result in higher performance, however, low occupancy always reduces the ability to hide latencies, resulting in overall performance degradation. Large discrepancies between the theoretical and the achieved occupancy during execution typically indicates highly imbalanced workloads.

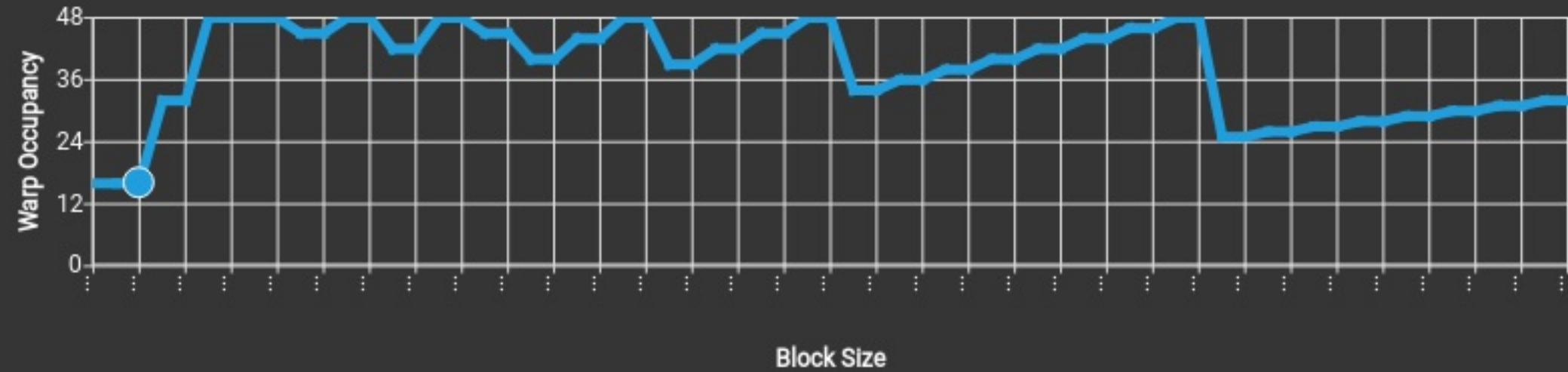
Theoretical Occupancy [%]	33,33	Block Limit Registers [block]	128
Theoretical Active Warps per SM [warp]	16	Block Limit Shared Mem [block]	16
Achieved Occupancy [%]	2,12	Block Limit Warps [block]	48
Achieved Active Warps Per SM [warp]	1,02	Block Limit SM [block]	16

Occupancy Limiters This kernel's theoretical occupancy (33.3%) is limited by the required amount of shared memory This kernel's theoretical occupancy (33.3%) is limited by the number of blocks that can fit on the SM The difference between calculated theoretical (33.3%) and measured achieved occupancy (2.1%) can be the result of warp scheduling overheads or workload imbalances during the kernel execution. Load imbalances can occur between warps within a block as well as across blocks of the same kernel. See the [CUDA Best Practices Guide](#) for more details on optimizing occupancy.

Impact of Varying Register Count Per Thread



Impact of Varying Block Size



Impact of Varying Shared Memory Usage Per Block

