

GPU Speed Of Light Throughput

GPU Throughput Chart

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 [%]	5,16	Duration [ms]	5,18
Memory Throughput [%]	39,94	Elapsed Cycles [cycle]	9.945.315
L1/TEX Cache Throughput [%]	28,44	SM Active Cycles [cycle]	9.819.361,08
L2 Cache Throughput [%]	30,10	SM Frequency [Ghz]	1,92
DRAM Throughput [%]	39,94	DRAM Frequency [Ghz]	7,99

Latency Issue

This workload exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of this device. Achieved compute throughput and/or memory bandwidth below 60.0% of peak typically indicate latency issues. Look at [Scheduler Statistics](#) and [Warp State Statistics](#) for potential reasons.

Key Performance Indicators

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	11.719	Function Cache Configuration	CachePreferNone
Registers Per Thread [register/thread]	18	Static Shared Memory Per Block [byte/block]	0
Block Size	256	Dynamic Shared Memory Per Block [byte/block]	0
Threads [thread]	3.000.064	Driver Shared Memory Per Block [Kbyte/block]	1,02
Waves Per SM	81,38	Shared Memory Configuration Size [Kbyte]	16,38
Uses Green Context	0	Stack Size	1.024
# SMs [SM]	24	# TPCs	12
Enabled TPC IDs	all	-	-

Occupancy

% Occupancy Graphs

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 [%]	100	Block Limit Registers [block]	10
Theoretical Active Warps per SM [warp]	48	Block Limit Shared Mem [block]	16
Achieved Occupancy [%]	88,46	Block Limit Warps [block]	6
Achieved Active Warps Per SM [warp]	42,46	Block Limit SM [block]	24

Achieved Occupancy

Est. Local Speedup: 11.54%

The difference between calculated theoretical (100.0%) and measured achieved occupancy (88.5%) 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.

Key Performance Indicators

GPU and Memory Workload Distribution

Analysis of workload distribution in active cycles of SM, SMP, SMSP, L1 & L2 caches, and DRAM

Average SM Active Cycles [cycle]	9.819.361,08	Average L1 Active Cycles [cycle]	9.819.361,08
Average L2 Active Cycles [cycle]	8.892.463,44	Average SMSP Active Cycles [cycle]	9.818.121,91
Average DRAM Active Cycles [cycle]	16.535.316	Total SM Elapsed Cycles [cycle]	240.571.584
Total L1 Elapsed Cycles [cycle]	240.571.584	Total L2 Elapsed Cycles [cycle]	144.091.376
Total SMSP Elapsed Cycles [cycle]	962.286.336	Total DRAM Elapsed Cycles [cycle]	165.611.520

Missing Roofline and Memory Charts? Profile again with the detailed metric set to collect all necessary metrics. Also consider collecting all available sections with the full metric set.