

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. High-level overview of the utilization for compute and memory resources of the GPU presented as a roofline chart.

Compute (SM) Throughput [%]	74.89	Duration [ms]	1.43
Memory Throughput [%]	16.89	Elapsed Cycles [cycle]	834,310
L1/TEX Cache Throughput [%]	18.68	SM Active Cycles [cycle]	811,024.80
L2 Cache Throughput [%]	0.34	SM Frequency [Mhz]	584.98
DRAM Throughput [%]	0.88	DRAM Frequency [Ghz]	4.99

High Compute Throughput

Compute is more heavily utilized than Memory: Look at the [Compute Workload Analysis](#) section to see what the compute pipelines are spending their time doing. Also, consider whether any computation is redundant and could be reduced or moved to look-up tables.

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 1% of this device's fp32 peak performance and 0% of its fp64 peak performance. See the [Kernel Profiling Guide](#) for more details on roofline analysis.

GPU Throughput

Compute Throughput Breakdown

Memory Throughput Breakdown

SM: Pipe Alu Cycles Active [%]	74.89	L1: Lsuin Requests [%]	16.89
SM: Issue Active [%]	64.39	L1: Data Pipe Lsu Wavefronts [%]	9.34
SM: Inst Executed [%]	64.38	L1: Lsu Writeback Active [%]	4.07
SM: Inst Executed Pipe Adu [%]	36.41	L1: Data Bank Reads [%]	0.97
IDC: Request Cycles Active [%]	24.54	DRAM: Cycles Active [%]	0.88
SM: Pipe Fma Cycles Active [%]	22.55	DRAM: Dram Sectors [%]	0.64
SM: Mio Inst Issued [%]	17.77	L1: Data Bank Writes [%]	0.53
SM: Inst Executed Pipe Lsu [%]	16.89	L1: M L1tex2xbar Req Cycles Active [%]	0.39
SM: Mio Pq Write Cycles Active [%]	16.20	L2: T Sectors [%]	0.34
SM: Mio Pq Read Cycles Active [%]	16.20	L2: Xbar2lts Cycles Active [%]	0.34
SM: Inst Executed Pipe Cbu Pred On Any [%]	12.77	L2: D Sectors [%]	0.15
SM: Mio2rf Writeback Active [%]	8.22	GPU: Compute Memory Access Throughput Internal Activity [%]	0.09
SM: Inst Executed Pipe Uniform [%]	0.05	L2: T Tag Requests [%]	0.09
SM: Memory Throughput Internal Activity [%]	0	L2: Lts2xbar Cycles Active [%]	0.00
SM: Pipe Tensor Cycles Active [%]	0	L2: D Sectors Fill Device [%]	0.00
SM: Pipe Shared Cycles Active [%]	0	L1: F Wavefronts [%]	0.00
SM: Pipe Fp64 Cycles Active [%]	0	L1: Texin Sm2tex Req Cycles Active [%]	0.00
SM: Instruction Throughput Internal Activity [%]	0	L1: Data Pipe Tex Wavefronts [%]	0
SM: Inst Executed Pipe Xu [%]	0	GPU: Compute Memory Request Throughput Internal Activity [%]	0
SM: Inst Executed Pipe Tex [%]	0	L1: M Xbar2l1tex Read Sectors [%]	0
SM: Inst Executed Pipe lpa [%]	0	L1: Tex Writeback Active [%]	0
SM: Inst Executed Pipe Fp16 [%]	0	L2: D Atomic Input Cycles Active [%]	0
		L2: D Sectors Fill Sysmem [%]	0

Floating Point Operations Roofline

Floating Point Operations Roofline (Double Precision)

Floating Point Operations Roofline (Half Precision)

Floating Point Operations Roofline (Single Precision)

Tensor Core Operations Roofline

	# Operations	# Operations / Cycle	# Operations / s	Peak %	Peak Operations / Cycle	Peak Operations / s
Src:fp16,bf16,tf32 Dst:fp32	0	0	0	0	40,960	23,960.44
Src:fp16 Dst:fp16	0	0	0	0	40,960	23,960.44
Src:int1	0	0	0	0	81,920	47,920.89
Src:int4	0	0	0	0	81,920	47,920.89
Src:int8	0	0	0	0	81,920	47,920.89

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]	811,024.80	Average L1 Active Cycles [cycle]	811,024.80
Average L2 Active Cycles [cycle]	46,334.31	Average SMSP Active Cycles [cycle]	731,784.65
Average DRAM Active Cycles [cycle]	62,439	Total SM Elapsed Cycles [cycle]	33,371,656
Total L1 Elapsed Cycles [cycle]	33,371,656	Total L2 Elapsed Cycles [cycle]	39,019,360
Total SMSP Elapsed Cycles [cycle]	133,486,624	Total DRAM Elapsed Cycles [cycle]	56,980,480

Workload Distribution

	Average	Min	Max	Sum
SM Active Cycles	811,024.80	793,217	828,930	32,440,992
SMSP Active Cycles	731,784.65	698,764	768,585	117,085,544
L1 Active Cycles	811,024.80	793,217	828,930	32,440,992
L2 Active Cycles	46,334.31	43,854	54,709	1,482,698
DRAM Active Cycles	62,439	61,860	63,064	499,512