

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 [%]	83.72	Duration [ms]	1.15
Memory Throughput [%]	58.68	Elapsed Cycles [cycle]	675,416
L1/TEX Cache Throughput [%]	61.75	SM Active Cycles [cycle]	660,293.62
L2 Cache Throughput [%]	5.32	SM Frequency [Mhz]	584.97
DRAM Throughput [%]	1.87	DRAM Frequency [Ghz]	4.99

High Throughput

The kernel is utilizing greater than 80.0% of the available compute or memory performance of the device. To further improve performance, work will likely need to be shifted from the most utilized to another unit. Start by analyzing workloads in the [Compute Workload Analysis](#) section.

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 19% 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 (SM) [%]

Memory [%]

Speed Of Light (SOL) [%]

Compute Throughput Breakdown

Memory Throughput Breakdown

SM: Issue Active [%]	83.72	L1: Lsuin Requests [%]	58.68
SM: Inst Executed [%]	83.72	L1: Data Pipe Lsu Wavefronts [%]	30.88
SM: Pipe Alu Cycles Active [%]	68.08	L1: Lsu Writeback Active [%]	26.19
SM: Inst Executed Pipe Lsu [%]	58.68	L1: Data Bank Reads [%]	15.04
SM: Inst Executed Pipe Xu [%]	53.35	L1: M L1tex2xbar Req Cycles Active [%]	6.12
SM: Inst Executed Pipe Adu [%]	52.65	L2: Xbar2lts Cycles Active [%]	5.32
IDC: Request Cycles Active [%]	51.19	L2: T Sectors [%]	5.08
SM: Pipe Fma Cycles Active [%]	49.16	DRAM: Cycles Active [%]	1.87
SM: Mio Inst Issued [%]	37.11	L2: D Sectors [%]	1.44
SM: Mio Pq Write Cycles Active [%]	28.02	L1: Data Bank Writes [%]	1.38
SM: Mio Pq Read Cycles Active [%]	28.01	DRAM: Dram Sectors [%]	1.37
SM: Mio2rf Writeback Active [%]	25.95	GPU: Compute Memory Access Throughput Internal Activity [%]	1.28
SM: Inst Executed Pipe Cbu Pred On Any [%]	0.12	L2: T Tag Requests [%]	1.26
SM: Pipe Tensor Cycles Active [%]	0	L1: M Xbar2l1tex Read Sectors [%]	0.87
SM: Pipe Shared Cycles Active [%]	0	L2: Lts2xbar Cycles Active [%]	0.75
SM: Pipe Fp64 Cycles Active [%]	0	L2: D Sectors Fill Device [%]	0.08
SM: Memory Throughput Internal Activity [%]	0	L1: Texin Sm2tex Req Cycles Active [%]	0.00
SM: Instruction Throughput Internal Activity [%]	0	L1: F Wavefronts [%]	0.00
SM: Inst Executed Pipe Uniform [%]	0	L2: D Atomic Input Cycles Active [%]	0
SM: Inst Executed Pipe Tex [%]	0	L2: D Sectors Fill Sysmem [%]	0
SM: Inst Executed Pipe Ipa [%]	0	L1: Tex Writeback Active [%]	0
SM: Inst Executed Pipe Fp16 [%]	0	GPU: Compute Memory Request Throughput Internal Activity [%]	0
		L1: Data Pipe Tex Wavefronts [%]	0

Floating Point Operations Roofline

Performance [FLOP/s] (1 = 100,000,000,000)

Arithmetic Intensity [FLOP/byte]

Floating Point Operations Roofline (Double Precision)

Performance [FLOP/s] (1 = 100,000,000,000)

Arithmetic Intensity [FLOP/byte]

Floating Point Operations Roofline (Half Precision)

Performance [FLOP/s] (1 = 1E+12)

Arithmetic Intensity [FLOP/byte]

Floating Point Operations Roofline (Single Precision)

Performance [FLOP/s] (1 = 1E+12)

Arithmetic Intensity [FLOP/byte]

Tensor Core Operations Roofline

Performance [OP/s]

Arithmetic Intensity [OP/byte]

	# Operations	# Operations / Cycle	# Operations / s	Peak %	Peak Operations / Cycle	Peak Operations / s
Src:fp16,bf16,tf32 Dst:fp32	0	0	0	0	40,960	23,970.00
Src:fp16 Dst:fp16	0	0	0	0	40,960	23,970.00
Src:int1	0	0	0	0	81,920	47,940.01
Src:int4	0	0	0	0	81,920	47,940.01
Src:int8	0	0	0	0	81,920	47,940.01

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]	660,293.62	Average L1 Active Cycles [cycle]	660,293.62
Average L2 Active Cycles [cycle]	282,723.47	Average SMSP Active Cycles [cycle]	657,002.41
Average DRAM Active Cycles [cycle]	107,818.50	Total SM Elapsed Cycles [cycle]	27,026,928
Total L1 Elapsed Cycles [cycle]	27,026,928	Total L2 Elapsed Cycles [cycle]	31,587,680
Total SMSP Elapsed Cycles [cycle]	108,107,712	Total DRAM Elapsed Cycles [cycle]	46,054,400

Workload Distribution

	Average	Min	Max	Sum
SM Active Cycles	660,293.62	645,474	670,542	26,411,745
SMSP Active Cycles	657,002.41	640,699	668,310	105,120,385
L1 Active Cycles	660,293.62	645,474	670,542	26,411,745
L2 Active Cycles	282,723.47	254,979	311,067	9,047,151
DRAM Active Cycles	107,818.50	100,316	120,912	862,548