

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 [%]

62.34

Duration [usecond]

971.46

Memory Throughput [%]

62.34

Elapsed Cycles [cycle]

1.320.445

L1/TEX Cache Throughput [%]

66.97

SM Active Cycles [cycle]

1.229.096.43

L2 Cache Throughput [%]

23.41

SM Frequency [cycle/nsecond]

1.36

DRAM Throughput [%]

36.19

DRAM Frequency [cycle/nsecond]

6.01

Balanced Throughput

Compute and Memory are well-balanced. To reduce runtime, both computation and memory traffic must be reduced. Check both the [Compute Workload Analysis](#) and [Memory Workload Analysis](#) report sections.

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 6% of this device's fp32 peak performance and 0% of its fp64 peak performance.

GPU Throughput

Speed Of Light (SOL) [%]

Compute (SM) [%]

Memory [%]

0,0

10,0

20,0

30,0

40,0

50,0

60,0

70,0

80,0

90,0

100,0

Compute Throughput Breakdown

Memory Throughput Breakdown

SM: Inst Executed Pipe Lsu [%]	62,34	L1: Lsuin Requests [%]	62,34
SM: Issue Active [%]	25,51	DRAM: Cycles Active [%]	36,19
SM: Inst Executed [%]	25,51	L1: Data Pipe Lsu Wavefronts [%]	31,30
SM: Mio Inst Issued [%]	23,62	L2: T Sectors [%]	23,41
SM: Pipe Alu Cycles Active [%]	17,01	DRAM: Dram Sectors [%]	20,71
SM: Pipe Fma Cycles Active [%]	14,88	L1: Lsu Writeback Active [%]	19,92
SM: Mio2rf Writeback Active [%]	11,38	L2: Xbar2fts Cycles Active [%]	15,23
SM: Inst Executed Pipe Adu [%]	8,51	L2: Lts2xbar Cycles Active [%]	13,14
SM: Mio Pq Write Cycles Active [%]	5,75	L2: D Sectors [%]	11,02
SM: Mio Pq Read Cycles Active [%]	5,75	L2: D Sectors Fill Device [%]	10,58
SM: Inst Executed Pipe Cbu Pred On Any [%]	1,50	L1: M L1tex2ubar Req Cycles Active [%]	8,28
SM: Pipe Tensor Cycles Active [%]	0	L1: Data Bank Reads [%]	7,88
SM: Pipe Shared Cycles Active [%]	0	L2: T Tag Requests [%]	7,77
SM: Pipe Fp64 Cycles Active [%]	0	L1: M Xbar2f11tex Read Sectors [%]	7,25
IDC: Request Cycles Active [%]	0	L1: Data Bank Writes [%]	2,41
SM: Inst Executed Pipe Xu [%]	0	L1: Texin Sm2tex Req Cycles Active [%]	0,04
SM: Inst Executed Pipe Uniform [%]	0	L1: F Wavefronts [%]	0,00
SM: Inst Executed Pipe Tex [%]	0	L1: Data Pipe Tex Wavefronts [%]	0
SM: Inst Executed Pipe Ipa [%]	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

Performance [GOPS] (G = 100,000,000,000)

Arithmetic Intensity [FLOP/byte]

0,01

0,1

1

10

100

1.000

0,01

1

10

Compute Workload Analysis

All

Detailed analysis of the compute resources of the streaming multiprocessors (SM), including the achieved instructions per clock (IPC) and the utilization of each available pipeline. Pipelines with very high utilization might limit the overall performance.

Executed ipc Elapsed [inst/cycle]

1,02

SM Busy [%]

27,40

Executed ipc Active [inst/cycle]

1,10

Issue Slots Busy [%]

27,40

Issued ipc Active [inst/cycle]

1,10

High Utilization

LSU is the highest-utilized pipeline (67.0%). It executes load/store memory operations. The pipeline is well-utilized and might become a bottleneck if more work is added.

Pipe Utilization

LSU

ALU

FMA

ADU

CBU

FP16

FP64

TEX

Tensor (FP)

Tensor (INT)

Uniform

XU

0,0

25,0

50,0

75,0

100,0

Utilization [%]

Memory Workload Analysis

All

Detailed analysis of the memory resources of the GPU. Memory can become a limiting factor for the overall kernel performance when fully utilizing the involved hardware units (Mem Busy), exhausting the available communication bandwidth between those units (Max Bandwidth), or by reaching the maximum throughput of issuing memory instructions (Mem Pipes Busy). Detailed chart of the memory units. Detailed tables with data for each memory unit. Deprecated UI elements for backwards compatibility.

Memory Throughput [Gbyte/second]

69.59

Mem Busy [%]

31,30

L1/TEX Hit Rate [%]

10.99

Max Bandwidth [%]

62,34

L2 Hit Rate [%]

55.43

Mem Pipes Busy [%]

62,34

Memory Chart

Kernel

Global

Local

Texture

Surface

Shared

L1/TEX Cache

L2 Cache

System Memory

Device Memory

Peer Memory

1.06 M Inst

0.00 Inst

0.00 Inst

0.00 Inst

3.16 M Inst

800.77 K Req

261.89 K Req

0.00 Req

0.00 Req

0.00 Req

0.00 Req

40.90 MB

31.97 MB

32.96 MB

31.51 MB

0.00 B

0.00 B

0.00 B

Hit Rate: 10.99 %

Hit Rate: 55.43 %

0%

20%

40%

60%

80%

100%

% Peak

Shared Memory

	Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	2.356.992	2.356.992	2.356.992	12,75	0
Shared Load Matrix	0	0	0	0	0
Shared Store	800.769	800.769	800.769	4,33	0
Shared Atomic	0	0	0	0	0
Other	-	-	1.549.822	8,38	0
Total	3.157.761	3.157.761	4.707.583	25,47	0

L1/TEX Cache

	Instructions	Requests	Wavefronts	% Peak	Sectors	Sectors/Req	Hit Rate	Bytes	Sector Misses to L2	% Peak to L2	Returns to SM	% Peak to SM
Local Load	0	0	0	0	0	0	0	0	0	0	0	0
Global Load	800.769	800.769	800.769	4,33	1.634.817	2,04	18,03	52.314.144	1.340.073	7,25	800.769	4,33
Surface Load	0	0	0	0	0	0	0	0	0	0	0	0
Texture Load	0	0	0	0	0	0	0	0	0	0	0	0
Global Store	261.888	261.888	261.888	1,42	1.047.552	4	0	33.521.664	1.047.552	5,67	-	-
Local Store	0	0	0	0	0	0	0	0	0	0	-	-
Surface Store	0	0	0	0	0	0	0	0	0	0	-	-
Global Reduction	0	0	0	0	0	0	0	0	0	0	-	-
Surface Reduction	0	0	0	0	0	0	0	0	0	0	-	-
Global Atomic ALU	0	0	0	0	0	0	0	0	0	0	see above	see above
Global Atomic CAS	0	0	0	0	0	0	0	0	0	0	see above	see above
Surface Atomic ALU	0	0	0	0	0	0	0	0	0	0	see above	see above
Surface Atomic CAS	0	0	0	0	0	0	0	0	0	0	see above	see above
Loads	800.769	800.769	800.769	4,33	1.634.817	2,04	18,03	52.314.144	1.340.073	7,25	800.769	4,33
Stores	261.888	261.888	261.888	1,42	1.047.552	4	0	33.521.664	1.047.552	5,67	-	-
Total	1.062.657	1.062.657	1.062.657	5,75	2.682.369	2,52	10,99	85.835.808	2.387.625	12,92	800.769	4,33

L2 Cache

	Requests	Sectors	Sectors/Req	% Peak	Hit Rate	Bytes	Throughput	Sector Misses to Device	Sector Misses to System	Sector Misses to Peer
L1/TEX Load	530.669	1.340.081	2,53	13,13	20,34	42.882.592	44.142.598.326,64	1.079.228	0	0
L1/TEX Store	261.888	1.047.552	4	10,27	100	33.521.664	34.506.620.989,53	0	0	0
L1/TEX Atomic ALU	0	0	0	0	0	0	0	0	0	0
L1/TEX Atomic CAS	0	0	0	0	0	0	0	0	0	0
L1/TEX Reduction	0	0	0	0	0	0	0	0	0	0
L1/TEX Total	792.552	2.387.633	3,01	23,40	55,07	76.404.256	78.649.219.316,16	1.079.252	0	0
GPU Total	792.718	2.388.319	3,01	23,41	55,09	76.426.208	78.671.816.325,19	1.079.251	0	0

Device Memory

	Sectors	% Peak	Bytes	Throughput
Load	1.079.969	18,50	34.559.008	35.574.444.956,85
Store	1.032.659	17,69	33.045.088	34.016.041.899,99
Total	2.112.628	36,19	67.604.096	69.590.486.856,84

L1/TEX Cache (Deprecated)

	Instructions	L1/TEX Requests	% Peak	Hit Rate	L2 Requests	% Peak	L2 Returns	% Peak	L1/TEX Returns	% Peak
Global Load Cached	800.769	800.769	4,33	18,03	-	-	-	-	-	-
Global Load Uncached	-	-	-	-	-	-	-	-	-	-
Local Load Cached	0	0	0	0	-	-	1.340.073	7,25	800.769	4,33
Local Load Uncached	-	-	-	-	-	-	-	-	-	-
Surface Load	0	0	0	0	-	-	0	0	-	-
Texture Load	0	0	0	0	-	-	0	0	0	0
Global Store	261.888	261.888	1,42	0	1.047.552	5,67	-	-	-	-
Local Store	0	0	0	0	-	-	-	-	-	-
Surface Store	0	0	0	0	0	0	-	-	-	-
Global Reduction	0	0	0	0	0	0	-	-	-	-
Surface Reduction	0	0	0	0	0	0	-	-	-	-
Global Atomic	0	0	0	0	0	0	-	-	-	-
Global Atomic Cas	0	0	0	0	0	0	-	-	see above	see above
Surface Atomic	0	0	0	0	0	0	-	-	see above	see above
Surface Atomic Cas	0	0	0	0	0	0	-	-	see above	see above
Loads	800.769	800.769	4,33	18,03	-	-	1.340.073	7,25	800.769	4,33
Stores	261.888	261.888	1,42	0	1.047.552	5,67	-	-	-	-
Total	1.062.657	1.062.657	5,75	10,99	1.047.552	5,67	1.340.073	7,25	800.769	4,33

L2 Cache (Deprecated)

	L2 Requests	% Peak	L2 Returns	% Peak	Total Bytes	Total Throughput
Global Load Cached	-	-	-	-	42.882.336	44.142.334.804,66
Global Load Uncached	-	-	-	-	-	-
Local Load Cached	-	-	1.340.073	7,25	-	-
Local Load Uncached	-	-	-	-	-	-
Surface Load	-	-	0	0	0	0
Texture Load	-	-	0	0	0	0
Global Store	1.047.552	5,67	-	-	33.521.664	34.506.620.989,53
Local Store	0	0	-	-	0	0
Surface Store	0	0	-	-	0	0
Global Reduction	0	0	-	-	0	0
Surface Reduction	0	0	-	-	0	0
Global Atomic	0	0	-	-	0	0
Global Atomic Cas	0	0	-	-	0	0
Surface Atomic	0	0	-	-	0	0
Surface Atomic Cas	0	0	-	-	0	0
Loads	-	-	1.340.073	7,25	42.882.336	44.142.334.804,66
Stores	1.047.552	5,67	-	-	33.521.664	34.506.620.989,53
Total	1.047.552	5,67	1.340.073	7,25	76.404.000	78.648.955.794,19

Scheduler Statistics

All

Summary of the activity of the schedulers issuing instructions. Each scheduler maintains a pool of warps that it can issue instructions for. The upper bound of warps in the pool (Theoretical Warps) is limited by the launch configuration. On every cycle each scheduler checks the state of the allocated warps in the pool (Active Warps). Active warps that are not stalled (Eligible Warps) are ready to issue their next instruction. From the set of eligible warps the scheduler selects a single warp from which to issue one or more instructions (Issued Warp). On cycles with no eligible warps, the issue slot is skipped and no instruction is issued. Having many skipped issue slots indicates poor latency hiding.

Active Warps Per Scheduler [warp]

7,71

No Eligible [%]

72,22

Eligible Warps Per Scheduler [warp]

1,05

One or More Eligible [%]

27,78

Issued Warp Per Scheduler

0,28

Issue Slot Utilization

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 3.6 cycles. This might leave hardware resources underutilized and may lead to less optimal performance. Out of the maximum of 8 warps per scheduler, this kernel allocates an average of 7.71 active warps per scheduler, but only an average of 1.05 warps were eligible per cycle. This is further reduced to 23.8 threads per warp due to predication. The compiler may use predication to avoid an actual branch. Instead, all instructions are scheduled, but a per-thread condition code or predicate controls which threads execute the instructions. Try to avoid different execution paths within a warp when possible. In addition, ensure your kernel makes use of Independent Thread Scheduling, which allows a warp to reconverge after a data-dependent conditional block by explicitly calling __syncwarp().

Warp State Statistics

All

Analysis of the states in which all warps spent cycles during the kernel execution. The warp states describe a warp's readiness or inability to issue its next instruction. The warp cycles per instruction define the latency between two consecutive instructions. The higher the value, the more warp parallelism is required to hide this latency. For each warp state, the chart shows the average number of cycles spent in that state per issued instruction. Stalls are not always impacting the overall performance nor are they completely avoidable. Only focus on stall reasons if the schedulers fail to issue every cycle. When executing a kernel with mixed library and user code, these metrics show the combined values.

Warp Cycles Per Issued Instruction [cycle]

27,75

Avg. Active Threads Per Warp

31,99

Warp Cycles Per Executed Instruction [cycle]

27,75

Avg. Not Predicated Off Threads Per Warp

23,80

Thread Divergence

Instructions are executed in warps, which are groups of 32 threads. Optimal instruction throughput is achieved if all 32 threads of a warp execute the same instruction. The chosen launch configuration, early thread completion, and divergent flow control can significantly lower the number of active threads in a warp per cycle. This kernel achieves an average of 32.0 threads being active per cycle. This is further reduced to 23.8 threads per warp due to predication. The compiler may use predication to avoid an actual branch. Instead, all instructions are scheduled, but a per-thread condition code or predicate controls which threads execute the instructions. Try to avoid different execution paths within a warp when possible. In addition, ensure your kernel makes use of Independent Thread Scheduling, which allows a warp to reconverge after a data-dependent conditional block by explicitly calling __syncwarp().

Instruction Statistics

All

Statistics of the executed low-level assembly instructions (SASS). The instruction mix provides insight into the types and frequency of the executed instructions. A narrow mix of instruction types implies a dependency on few instruction pipelines, while others remain unused. Using multiple pipelines allows hiding latencies and enables parallel execution. Note that 'Instructions/Opcode' and 'Executed Instructions' are measured differently and can diverge if cycles are spent in system calls.

Executed Instructions [inst]

18.860.032

Avg. Executed Instructions Per Scheduler [inst]

336.786,29

Issued Instructions [inst]

18.862.582

Avg. Issued Instructions Per Scheduler [inst]

336.831,82

NVLink

All

High-level summary of NVLink utilization. It shows the total received and transmitted (sent) memory, as well as the overall link peak utilization. NVLink Topology. Detailed tables with properties for each NVLink.

Physical Links

0

Logical Links

0

Launch Statistics

All

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

8.192

Registers Per Thread [register/thread]

26

Block Size

1.024

Static Shared Memory Per Block [kbyte/block]

4,62

Threads [thread]

8.388.608

Dynamic Shared Memory Per Block [byte/block]

0

Waves Per SM

585,14

Driver Shared Memory Per Block [byte/block]

0

Function Cache Configuration

cudaFuncCachePreferNone

Shared Memory Configuration Size [kbyte]

32,77

Occupancy

All

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]

2

Theoretical Active Warps per SM [warp]

32

Block Limit Shared Mem [block]

13

Achieved Occupancy [%]

95,13

Block Limit Warps [block]

1

Achieved Active Warps per SM [warp]

30,76

Block Limit SM [block]

16

Occupancy Limiters

This kernel's theoretical occupancy is not impacted by any block limit.

Source Counters

All

Source metrics, including branch efficiency and sampled warp stall reasons. Sampling Data metrics are periodically sampled over the kernel runtime. They indicate when warps were stalled and couldn't be scheduled. See the documentation for a description of all stall reasons. Only focus on stalls if the schedulers fail to issue every cycle.

Branch Instructions [inst]

524.032

Branch Efficiency [%]

0

Branch Instructions Ratio [%]

0,03

Avg. Divergent Branches

0