

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 routine chart.

Compute (SM) Throughput [%]

Memory Throughput [%]

L1/TEX Cache Throughput [%]

L2 Cache Throughput [%]

DRAM Throughput [%]

36.51

43.69

41.35

43.69

38.60

Duration [usecond]

Elapsed Cycles [cycle]

SM Active Cycles [cycle]

SM Frequency [cycle/nsecond]

DRAM Frequency [cycle/nsecond]

899.55

1,229,299

1,137,670.79

1.37

5.99

Latency Issue

This kernel 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.

Routine 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

Compute (SM) [%]

Memory [%]

Speed of Light (SOL) [%]

Compute Throughput Breakdown

Memory Throughput Breakdown

SM: Inst Executed Pipe Lsu [%]

SM: Issue Active [%]

SM: Inst Executed [%]

SM: Pipe Alu Cycles Active [%]

SM: Pipe Fma Cycles Active [%]

SM: Mio Inst Issued [%]

SM: Mio2rf Writeback Active [%]

SM: Inst Executed Pipe Xu [%]

SM: Inst Executed Pipe Cbu Pred On Any [%]

SM: Inst Executed Pipe Adu [%]

SM: Inst Executed Pipe Uniform [%]

SM: Mio Pq Read Cycles Active [%]

SM: Mio Pq Write Cycles Active [%]

SM: Pipe Tensor Cycles Active [%]

SM: Pipe Shared Cycles Active [%]

SM: Pipe Fp64 Cycles Active [%]

IDC: Request Cycles Active [%]

SM: Inst Executed Pipe Tex [%]

SM: Inst Executed Pipe Ipa [%]

SM: Inst Executed Pipe Fp16 [%]

L2: T Sectors [%]

DRAM: Cycles Active [%]

L1: Lsum Requests [%]

L2: Lts2xbar Cycles Active [%]

L2: Xbar2lts Cycles Active [%]

DRAM: Drm Sectors [%]

L1: Data Pipe Lsu Wavefronts [%]

L1: Lsu Writeback Active [%]

L1: M Xbar2l1tex Read Sectors [%]

L2: D Sectors [%]

L2: T Tag Requests [%]

L2: D Sectors Fill Device [%]

L1: M L1tex2xbar Req Cycles Active [%]

L1: Data Bank Reads [%]

L1: Data Bank Writes [%]

L1: F Wavefronts [%]

L1: Texin Sm2tex Req Cycles Active [%]

L1: Data Pipe Tex Wavefronts [%]

L1: Tex Writeback Active [%]

L2: D Atomic Input Cycles Active [%]

L2: D Sectors Fill System [%]

43.69

38.60

36.51

32.67

22.91

21.97

20.67

18.22

17.49

16.43

16.08

11.09

10.88

7.60

3.01

0.00

0.00

0

0

0

0

0

0

Floating Point Operations Routine

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

Arithmetic Intensity (FLOP/byte)

10

1

0.1

0.01

0.01

0.1

1

10

100

1,000

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]

Executed Ipc Active [inst/cycle]

Issued Ipc Active [inst/cycle]

1,19

1,28

1,28

SM Busy [%]

Issue Slots Busy [%]

32.06

32.06

Balanced

No pipeline is over-utilized.

Pipe Utilization

LSU

ALU

FMA

XU

ADU

Uniform

CBU

FP16

FP64

TEX

Tensor (FP)

Tensor (INT)

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. Depreciated (i) elements for backwards compatibility.

L1/TEX Hit Rate [%]

L2 Hit Rate [%]

74.04

66.04

79.51

Mem Busy [%]

Max Bandwidth [%]

Mem Pipes Busy [%]

43.69

38.60

36.51

Memory Chart

Kernel

Global

Local

Texture

Surface

Shared

L1/TEX Cache

L2 Cache

System Memory

Device Memory

Peer Memory

2.62 M Inst

0.00 Inst

0.00 Inst

0.00 Inst

0.00 Inst

0.00 Inst

2.35 M Req

261.63 K Req

0.00 Req

0.00 Req

0.00 Req

0.00 Req

0.00 Req

0.00 Req

0.00 Req

0.00 Req

94.46 MB

31.94 MB

0.00 B

0.00 B

32.00 MB

31.51 MB

0.00 B

0.00 B

Hit Rate: 66.04 %

Hit Rate: 79.51 %

% Peak

100%

80%

60%

40%

20%

0%

	Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	0	0	0	0	0
Shared Load Matrix	0	0	0	0	0
Shared Store	0	0	0	0	0
Shared Atomic	0	0	0	0	0
Other	-	-	524,288	3,05	0
Total	0	0	524,288	3,05	0

	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	3,100,644	18,02	2,608,911	15,16
Global Load	2,354,688	2,354,688	2,458,932	14,29	10,896,403	4,63	72,39	348,684,896	0	0	0	0
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	261,632	261,632	261,632	1,52	1,046,528	4	0	33,488,896	1,046,528	6,08	-	-
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	-	-
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	-	-
Loads	2,354,688	2,354,688	2,458,932	14,29	10,896,403	4,63	72,39	348,684,896	3,100,644	18,02	2,608,911	15,16
Stores	261,632	261,632	261,632	1,52	1,046,528	4	0	33,488,896	1,046,528	6,08	-	-
Total	2,616,320	2,616,320	2,720,564	15,81	11,942,931	4,56	66,05	382,173,792	4,147,172	24,11	2,608,911	15,16

	Requests	Sectors	Sectors/Req	% Peak	Hit Rate	Bytes	Throughput	Sector Misses to Device	Sector Misses to System	Sector Misses to Peer
L1/TEX Load	1,263,370	3,094,212	2,45	32,67	68,09	99,014,784	110,071,217,672,80	1,048,576	0	0
L1/TEX Store	261,632	1,046,528	4	11,05	100	33,488,896	37,228,415,922,59	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	1,525,102	4,137,366	2,71	43,69	75,84	132,395,712	147,179,609,405,57	1,048,576	0	0
GPU Total	1,525,145	4,137,664	2,71	43,69	75,83	132,405,248	147,190,210,237,99	1,048,585	0	0

	Sectors	% Peak	Bytes	Throughput
Load	1,048,691	19,45	33,558,112	37,305,360,890,75
Store	1,022,522	19,15	33,040,704	36,790,176,799,12
Total	2,081,213	38,60	66,598,816	74,035,537,689,87

	Instructions	L1/TEX Requests	% Peak	Hit Rate	L2 Requests	% Peak	L2 Returns	% Peak	L1/TEX Returns	% Peak
Global Load Cached	2,354,688	2,354,688	13,69	72,39	-	-	-	-	-	-
Global Load Uncached	-	-	-	-	-	-	3,100,644	18,02	2,608,911	15,16
Local Load Cached	-	0	0	0	-	-	-	-	-	-
Local Load Uncached	-	-	-	-	-	-	-	-	-	-
Surface Load	0	0	0	0	-	-	0	0	0	0
Texture Load	0	0	0	0	-	-	0	0	0	0
Global Store	261,632	261,632	1,52	0	1,046,528	6,08	-	-	-	-
Local Store	0	0	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	0	0	see above	see above
Surface Atomic	0	0	0	0	0	0	0	0	see above	see above
Surface Atomic Cas	0	0	0	0	0	0	0	0	see above	see above
Loads	2,354,688	2,354,688	13,69	72,39	-	-	3,100,644	18,02	2,608,911	15,16
Stores	261,632	261,632	1,52	0	1,046,528	6,08	-	-	-	-
Total	2,616,320	2,616,320	15,21	66,05	1,046,528	6,08	3,100,644	18,02	2,608,911	15,16

	L2 Requests	% Peak	L2 Returns	% Peak	Total Bytes	Total Throughput
Global Load Cached	-	-	-	-	99,220,608	110,300,024,901,28
Global Load Uncached	-	-	-	-	-	-
Local Load Cached	-	-	3,100,644	18,02	-	-
Local Load Uncached	-	-	-	-	-	-
Surface Load	-	-	0	0	0	0
Texture Load	-	-	0	0	0	0
Global Store	1,046,528	6,08	-	-	33,488,896	37,228,415,922,59
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	0	0
Surface Atomic	0	0	0	0	0	0
Surface Atomic Cas	0	0	0	0	0	0
Loads	-	-	3,100,644	18,02	99,220,608	110,300,024,901,28
Stores	1,046,528	6,08	-	-	33,488,896	37,228,415,922,59
Scheduler	1,046,528	6,08	3,100,644	18,02	132,709,504	147,528,440,823,88

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]

Eligible Warps Per Scheduler [warp]

Issued Warp Per Scheduler

7,12

0,95

0,33

No Eligible [%]

One or More Eligible [%]

67,35

32,65

Issue Slot Utilization

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 3.1 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 7.12 active warps per scheduler, but only an average of 0.95 warps were eligible per cycle. Eligible warps are the subset of active warps that are ready to issue their next instruction. Every cycle with no eligible warp results in no instruction being issued and the issue slot remains unused. To increase the number of eligible warps, avoid possible load imbalances due to highly different execution durations per warp. Reducing stalls indicated on the [Warp State Statistics](#) and [Warp Cycles](#) sections can help, too.

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]

Warp Cycles Per Executed Instruction [cycle]

21,82

21,82

Avg. Active Threads Per Warp

Avg. Not Predicated Off Threads Per Warp

32,00

29,95

long_scoreboard

On average, each warp of this kernel spends 8.8 cycles being stalled waiting for a scoreboard dependency on a L1/TEX (local, global, surface, texture) operation. This represents about 40.3% of the total average of 21.8 cycles between issuing two instructions. To reduce the number of cycles waiting on L1/TEX data accesses verify the memory access patterns are optimal for the target architecture, attempt to increase cache hit rates by increasing data locality or by changing the cache configuration, and consider moving frequently used data to shared memory.

Warp Stall

Check the [Source Counters](#) section for the top stall locations in your source based on sampling data.

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]

Issued Instructions [inst]

20,423,168

20,426,136

Avg. Executed Instructions Per Scheduler [inst]

Avg. Issued Instructions Per Scheduler [inst]

364,699,43

364,752,43

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

Logical Links

0

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

Block Size

Threads [thread]

Waves Per SM

Function Cache Configuration

8,192

1,024

8,388,608

595,14

cudaFuncCachePrefNone

Registers Per Thread [register/thread]

Static Shared Memory Per Block [byte/block]

Dynamic Shared Memory Per Block [byte/block]

Driver Shared Memory Per Block [byte/block]

Shared Memory Configuration Size [kbyte]

28

0

0

0

32,77

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 [%]	100	Block Limit Registers [block]	2
Theoretical Active Warps per SM [warp]	32	Block Limit Shared Mem [block]	16
Achieved Occupancy [%]	88,86	Block Limit Warps [block]	1
Achieved Active Warps Per SM [warp]	28,34	Block Limit SM [block]	16

Occupancy Limiters

This kernel's theoretical occupancy is not impacted by any block limit. The difference between calculated theoretical (100.0%) and measured achieved occupancy (88.6%) 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.

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]

Branch Instructions Ratio [%]

1,048,064

0,05

Branch Efficiency [%]

Avg. Divergent Branches

100

0

Uncoalesced Global Accesses

Uncoalesced global access, expected 1046528 sectors, got 1307138 (1.25x) at PC [0x7f9856f9b390](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 1046528 sectors, got 1307138 (1.25x) at PC [0x7f9856f9b370](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 1046528 sectors, got 1307138 (1.25x) at PC [0x7f9856f9b3a0](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 1046528 sectors, got 1307138 (1.25x) at PC [0x7f9856f9b3b0](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 1046528 sectors, got 1307138 (1.25x) at PC [0x7f9856f9b3c0](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 1046528 sectors, got 1307138 (1.25x) at PC [0x7f9856f9b400](#)