

GPU Speed of Light Throughput

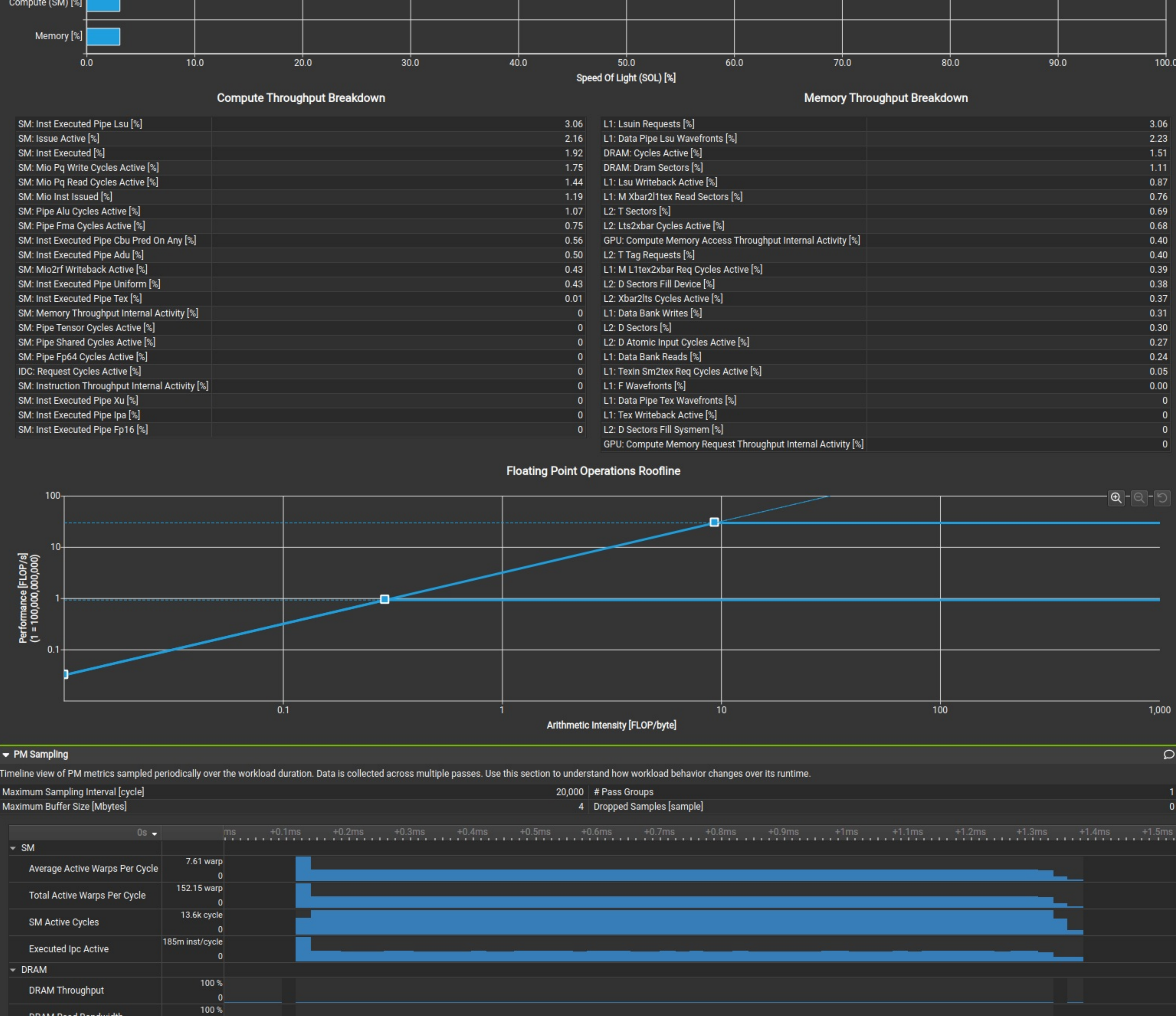
All

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 timeline chart.

Compute GPU Throughput [%]	3.06	Elapsed Cycles [cycle]	732,570
Memory Throughput [%]	3.06	Elapsed Cycles [cycle]	732,570
L1/TEX Cache Throughput [%]	4.45	SM Active Cycles [cycle]	714,213.30
L2 Cache Throughput [%]	0.69	SM Frequency [Mhz]	584.97
DRAM Throughput [%]	1.51	DRAM Frequency [Ghz]	5.00

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 the [Kernel Profiling Guide](#) and [Warp State Statistics](#) for potential reasons.

Roofline Analysis The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 0% 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.



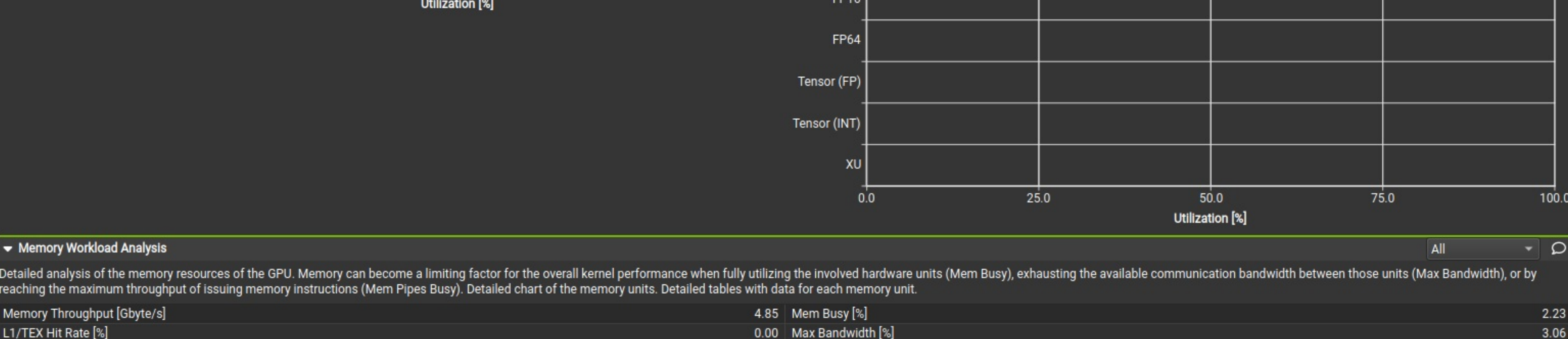
SM: Inst Executed Pipe Lru [%]	3.06	L1: Lru Requests [%]	3.06
SM: Issue Active [%]	2.16	L1: Data Pipe Lru Wavefronts [%]	2.23
SM: Inst Executed [%]	1.92	DRAM: Cycles Active [%]	1.51
SM: Mio Pq Write Cycles Active [%]	1.75	DRAM: Dram Sectors [%]	1.11
SM: Mio Pq Read Cycles Active [%]	1.44	L1: Lru Waveback Active [%]	0.87
SM: Mio Inst Issued [%]	1.19	L1: M Xbar211r Read Sectors [%]	0.76
SM: Pipe Au Cycles Active [%]	1.07	L2: T Sectors [%]	0.69
SM: Pipe Fma Cycles Active [%]	0.75	L2: L1x2zbar Cycles Active [%]	0.68
SM: Inst Executed Pipe Cpu Fired On Any [%]	0.56	GPU: Compute Memory Access Throughput Internal Activity [%]	0.40
SM: Inst Executed Pipe Au [%]	0.50	L2: T Tag Requests [%]	0.31
SM: Mio2rf Waveback Active [%]	0.43	L1: M L1x2zbar Req Cycles Active [%]	0.39
SM: Inst Executed Pipe Uniform [%]	0.43	L2: D Sectors Fill Device [%]	0.38
SM: Inst Executed Pipe Tex [%]	0.01	L2: Xbar211s Cycles Active [%]	0.37
SM: Memory Throughput Internal Activity [%]	0	L1: Data Bank Writes [%]	0.31
SM: Pipe Tensor Cycles Active [%]	0	L2: D Sectors [%]	0.30
SM: Pipe Shared Cycles Active [%]	0	L2: D Atomic Input Cycles Active [%]	0.27
SM: Pipe Pp4 Cycles Active [%]	0	L1: Data Bank Reads [%]	0.24
EOC: Request Cycles Active [%]	0	L1: Tex Sector Req Cycles Active [%]	0.05
SM: Instruction Throughput Internal Activity [%]	0	L1: F Wavefronts [%]	0.00
SM: Inst Executed Pipe Xu [%]	0	L1: Data Pipe Tex Wavefronts [%]	0
SM: Inst Executed Pipe Spa [%]	0	L1: Tex Waveback Active [%]	0
SM: Inst Executed Pipe Pp16 [%]	0	L2: D Sectors Fill System [%]	0
		GPU: Compute Memory Request Throughput Internal Activity [%]	0

SM: Average Active Warps Per Cycle	7.61 warp		
Total Active Warps Per Cycle	152.15 warp		
SM Active Cycles	13.6k cycle		
Executed Ipc Active	185m inst/cycle		

DRAM Throughput	100 %		
DRAM Read Bandwidth	100 %		
DRAM Write Bandwidth	100 %		
Writeback Throughput	210.43 cycle		
Hit Rate	100 %		
Wavefronts (Data)	561.70		
Workload Execution	0		

Compute Workload Analysis			
Executed Ipc Elapsed [inst/cycle]	0.08	SM Busy [%]	2.21
Executed Ipc Active [inst/cycle]	0.08	Issue Slots Busy [%]	2.21
Issued Ipc Active [inst/cycle]	0.09		

Low Utilization Est. Local Speedup: 98.90% All compute pipelines are underutilized. Either this kernel is very small or it doesn't issue enough warps per scheduler. Check the [Launch Statistics](#) and [Scheduler Statistics](#) sections for further details.



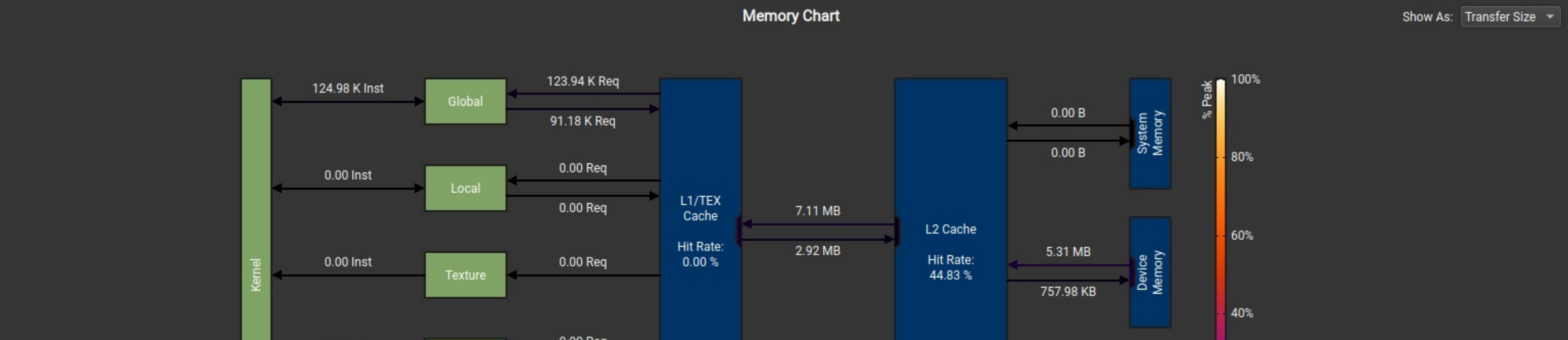
Memory Throughput [cbyte/s]	4.85	Mem Busy [%]	2.23
L1/TEX Hit Rate [%]	0.00	Max Bandwidth [%]	3.06
L2 Hit Rate [%]	44.63	Mem Pipes Busy [%]	3.06

Memory L2 Compression The optional metric lru_average_gcomp_input_sector_success_rate_pct could not be found. Collecting it as an additional metric could enable the rule to provide more guidance.

L1/TEX Global Load Access Pattern Est. Speedup: 0.03% The memory access pattern for global loads from L1/TEX might not be optimal. On average, only 31.8 of the 32 bytes transmitted per sector are utilized by each thread. This could possibly be caused by a stride between threads. Check the [Source Locations](#) section for uncoalesced global loads.

L1/TEX Global Store Access Pattern Est. Speedup: 3.89% The memory access pattern for global stores to L1/TEX might not be optimal. On average, only 4.0 of the 32 bytes transmitted per sector are utilized by each thread. This could possibly be caused by a stride between threads. Check the [Source Locations](#) section for uncoalesced global stores.

Shared Load Bank Conflicts Est. Speedup: 0.57% The memory access pattern for shared loads might not be optimal and causes on average a 1.2 - way bank conflict across all 94504 shared load requests. This results in 14219 bank conflicts, which represent 12.70% of the overall 111953 wavefronts for shared loads. Check the [Source Locations](#) section for uncoalesced shared loads.



Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	94,504	94,504	111,953	0.38
Shared Load Matrix	0	0	14,219	
Shared Store	102,992	102,992	133,209	0.52
Shared Atomic	0	0	0	
Other	-	-	139,393	0.48
Total	197,496	197,496	404,555	1.38
				14,219

Instructions	Requests	Wavefronts	% Peak	Sectors	Sectors/Req	Hit Rate	Bytes	Sector Misses to L2	% Peak to L2	Returns to SM
Global Load	33,792	33,792	33,792	3,931	0.12	132,096	3,931	0.00	4,227,072	132,096
Surface Load	0	0	0	0	0.00	0	0	0	0	0
Store	0	0	1,031	0.00	1,031	1.00	0.19	32,928	0	0
Global Store	1,031	1,031	1,031	0.00	1,031	1.00	0.19	32,928	0	0
Local Store	0	0	0	0	0.00	0	0	0	1,029	0.00
Surface Store	0	0	0	0	0.00	0	0	0	0	
Global Reduction	0	0	0	0	0.00	0	0	0	0	
DSMEM Reduction	0	0	0	0	0.00	0	0	0	0	
Surface Reduction	0	0	0	0	0.00	0	0	0	0	
Global Atomic ALU	0	0	90,086	0.31	90,291	1.00	0	2,889,312	90,291	0.02
Global Atomic CAS	90,145	90,145	90,086	0.31	90,291	1.00	0	2,889,312	90,291	0.02
Surface Atomic ALU	0	0	0	0	0.00	0	0	0	0	
Surface Atomic CAS	0	0	0	0	0.00	0	0	0	0	
Loads	33,792	33,792	33,792	0.12	132,096	3,931	0.00	4,227,072	132,096	0.45
Stores	1,031	1,031	1,031	0.00	1,031	1.00	0.19	32,928	1,029	0.00
Atoms & Reductions	90,145	90,145	90,086	0.31	90,291	1.00	0	2,889,312	90,291	0.02
Total	124,976	124,968	124,909	0.43	223,416	1.79	0.00	7,149,312	223,415	0.47

Requests	Sectors	Sectors/Req	% Peak	Hit Rate	Bytes	Throughput	Sector Misses to Device	Sector Misses to System	Sector Misses to Peer
L1/TEX Load	33,791	132,095	3,931	0.39	0.77	4,227,040	3,375,367,318.26	131,072	0
L1/TEX Store	1,035	1,035	1	0.00	99.81	33,120	26,446,914.53	2	0
L1/TEX Atomic ALU	-17	108	-6.35	0.00	100	3,456	2,759,678.04	0	0
L1/TEX Atomic CAS	90,289	90,194	1.00	0.26	100.00	2,886,208	2,304,688,897.41	2	0
L1/TEX Reduction	0	0	0	0	0	0	0	0	0
L1/TEX Total	125,159	223,318	1.78	0.65	41.34	7,146,176	5,706,349,814.74	131,076	0
EOC Total	-	10	-	0.00	-	320	255,525.74	10	-
GPU Total	136,065	234,826	1.73	0.69	44.38	7,514,432	6,000,408,841.19	131,090	0

Sectors	% Peak	Bytes	Throughput
Load	146,042	1.33	5,213,244
Store	23,687	0.19	757,984
Total	169,729	1.51	6,071,328
			4,848,064,392.49

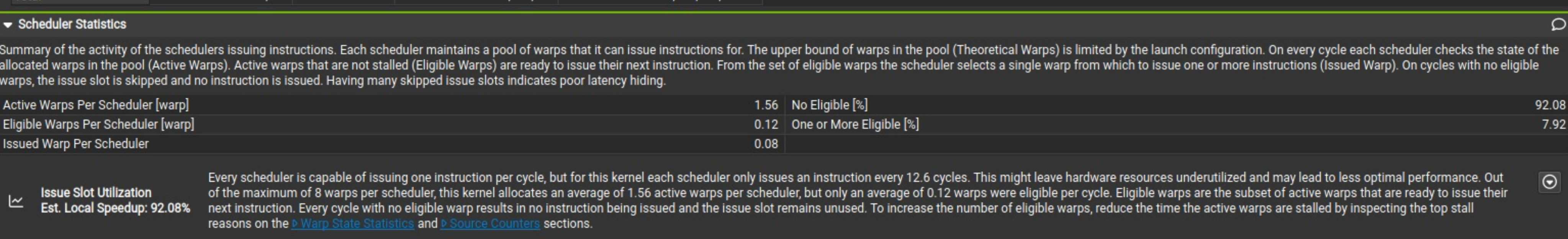
Scheduler Statistics Summary of the activity of the schedulers issuing instructions. Each scheduler maintains a pool of warps that it can issue instructions for. The upper bound of warps in the pool (Theoretical Warps) is limited by the launch configuration. On every cycle each scheduler checks the state 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] 1.56 No Eligible [%] 92.08

Issued Warps Per Scheduler [warp] 0.12 One or More Eligible [%] 7.92

Issued Warp Per Scheduler 0.08

Issue Slot Utilization Est. Local Speedup: 12.08% Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 12.6 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 1.56 active warps per scheduler, but only an average of 0.12 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, reduce the time the active warps are stalled by inspecting the top stall reasons on the [Warp State Statistics](#) and [Source Locations](#) sections.



GPU Maximum Warps Per Scheduler	12.08
Theoretical Warps Per Scheduler	12.08
Active Warps Per Scheduler	1.56
Eligible Warps Per Scheduler	0.12
Issued Warp Per Scheduler	0.08

Warp State Statistics 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 warps are stalled. Lower, the less time a warp spends in a particular state. This kernel achieves an average of 18.5 threads being active per cycle. This is further reduced to 15.5 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.

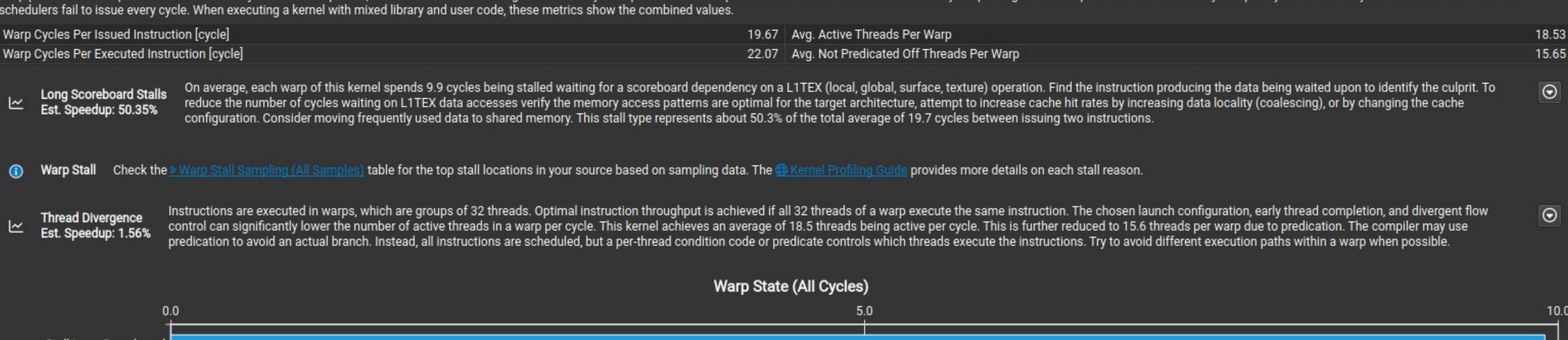
Warp Cycles Per Issued Instruction [cycle] 19.67 Avg. Active Threads Per Warp 18.53

Warp Cycles Per Executed Instruction [cycle] 22.07 Avg. Not Predicated Off Threads Per Warp 15.65

Long Scoreboard Stalls Est. Speedup: 50.35% On average, each warp of this kernel spends 9.9 cycles being stalled waiting for a scoreboard dependency on a L1/TEX (local, global, surface, texture) operation. Find the instruction producing the data being waited upon to identify the culprit. 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 (coalescing), or by changing the cache configuration. Consider moving frequently used data to shared memory. This stall type represents about 50.3% of the total average of 19.7 cycles between issuing two instructions.

Warp Stall Check the [Warp Stall Sampling \(All Samples\)](#) table for the top stall locations in your source based on sampling data. The [Kernel Profiling Guide](#) provides more details on each stall reason.

Thread Divergence Est. Speedup: 1.56% 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 can significantly lower the number of active threads in a warp per cycle. This kernel achieves an average of 18.5 threads being active per cycle. This is further reduced to 15.5 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.



Opcodes	Executed Warp-Level Instructions/Opcode
IMAD	380,000.0
ISETP	200,000.0
BRA	100,000.0
STS	100,000.0
LDS	100,000.0
LEA	100,000.0
ULDC	100,000.0
IADD3	100,000.0
ATOMG	100,000.0
YIELD	100,000.0
VOTEU	100,000.0
VOTE	100,000.0
LDP3	100,000.0
BAR	100,000.0
EXIT	100,000.0
WARPSYNC	100,000.0
FSFETP	100,000.0
SHF	100,000.0
STG	100,000.0
LDG	100,000.0
USHF	100,000.0
UIMAD	100,000.0
UIADD3	100,000.0
S2UR	100,000.0
S2R	100,000.0
MEMBAR	100,000.0
ERRBAR	100,000.0
CTL	100,000.0

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] 714,213.30 Average L1 Active Cycles [cycle] 714,213.30

Average L2 Active Cycles [cycle] 112,371.41 Average SMSP Active Cycles [cycle] 199,618.64

Average DRAM Active Cycles [cycle] 94,864.50 Total SM Elapsed Cycles [cycle] 29,316,624

Total L1 Elapsed Cycles [cycle] 29,316,624 Total L2 Elapsed Cycles [cycle] 34,251,696

Total SMSP Elapsed Cycles [cycle] 117,266,496 Total DRAM Elapsed Cycles [cycle] 50,122,752

SMSP Workload Imbalance Est. Speedup: 12.08% One or more SMSPs have a much lower number of active cycles than the average number of active cycles. Maximum instance value is 44.34% above the average, while the minimum instance value is 71.05% below the average.

L2 Slices Workload Imbalance Est. Speedup: 9.34% One or more L2 Slices have a much higher number of active cycles than the average number of active cycles. Maximum instance value is 89.00% above the average, while the minimum instance value is 29.75% below the average.

Workload Distribution	Average	Min	Max	Sum
SM Active Cycles	714,213.30	699,701	730,224	28,848,532
SMSP Active Cycles	199,618.64	67,795	388,071	31,938,883
L1 Active Cycles	714,213.30	699,701	730,224	28,568,532
L2 Active Cycles	112,371.41	78,941	1,021,835	3,595,885
DRAM Active Cycles	94,864.50	94,284	95,652	758,916

Source Counters Source metrics, including branch efficiency and sampled warp stall reasons. Warp Stall Sampling 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] 415,461 Branch Efficiency [%] 15.31%

Branch Instructions Ratio [%] 0.18 Avg. Divergent Branches 32

Uncoalesced Shared Accesses Est. Speedup: 11.25% This kernel has uncoalesced shared accesses resulting in a total of 96622 excessive wavefronts (12% of the total 265162 wavefronts). Check the L1 Wavefronts Shared Excessive table for the primary source locations. The [CUDA doc](#) has an example on optimizing shared memory accesses.

Location	Value	Value (%)
kernel on 172 (0x7a6075937e40 in reduce_argmin)	15,911	56
kernel on 186 (0x7a6075937e40 in reduce_argmin)	15,311	58
kernel on 186 (0x7a6075937e40 in reduce_argmin)	0	0
kernel on 177 (0x7a6075937e40 in reduce_argmin)	0	0
kernel on 177 (0x7a6075937e40 in reduce_argmin)	0	0

Follow the rules outputs to get guidance on how to navigate through the report and quickly discover performance bottlenecks in this kernel. You could also disable [warp state](#) to focus on selected performance aspects and make profiling faster.