

▼ GPU Speed of Light Throughput

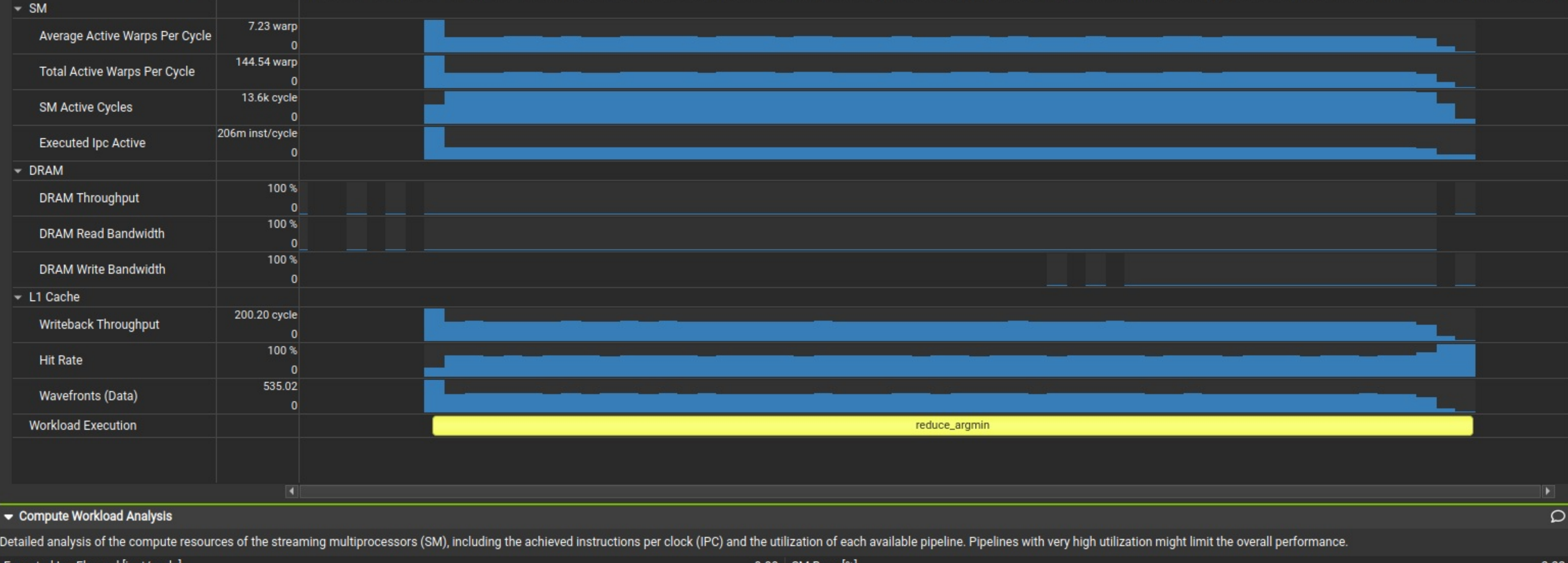
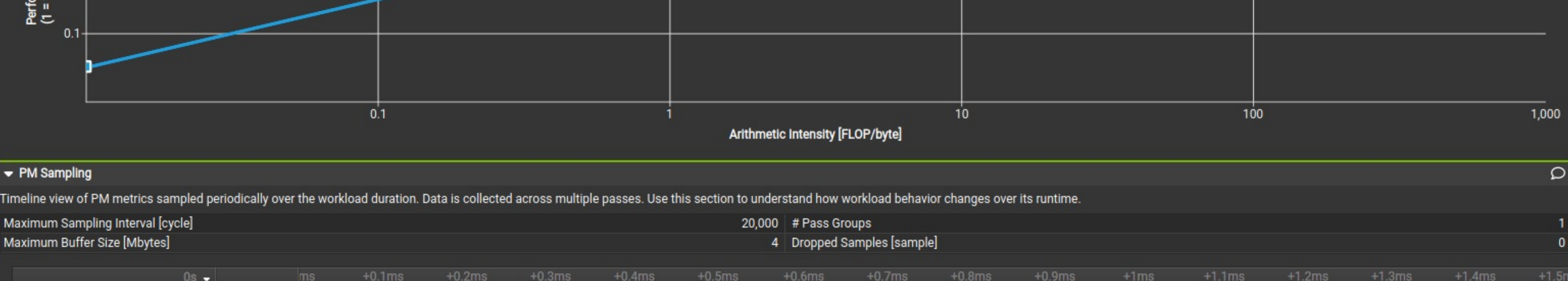
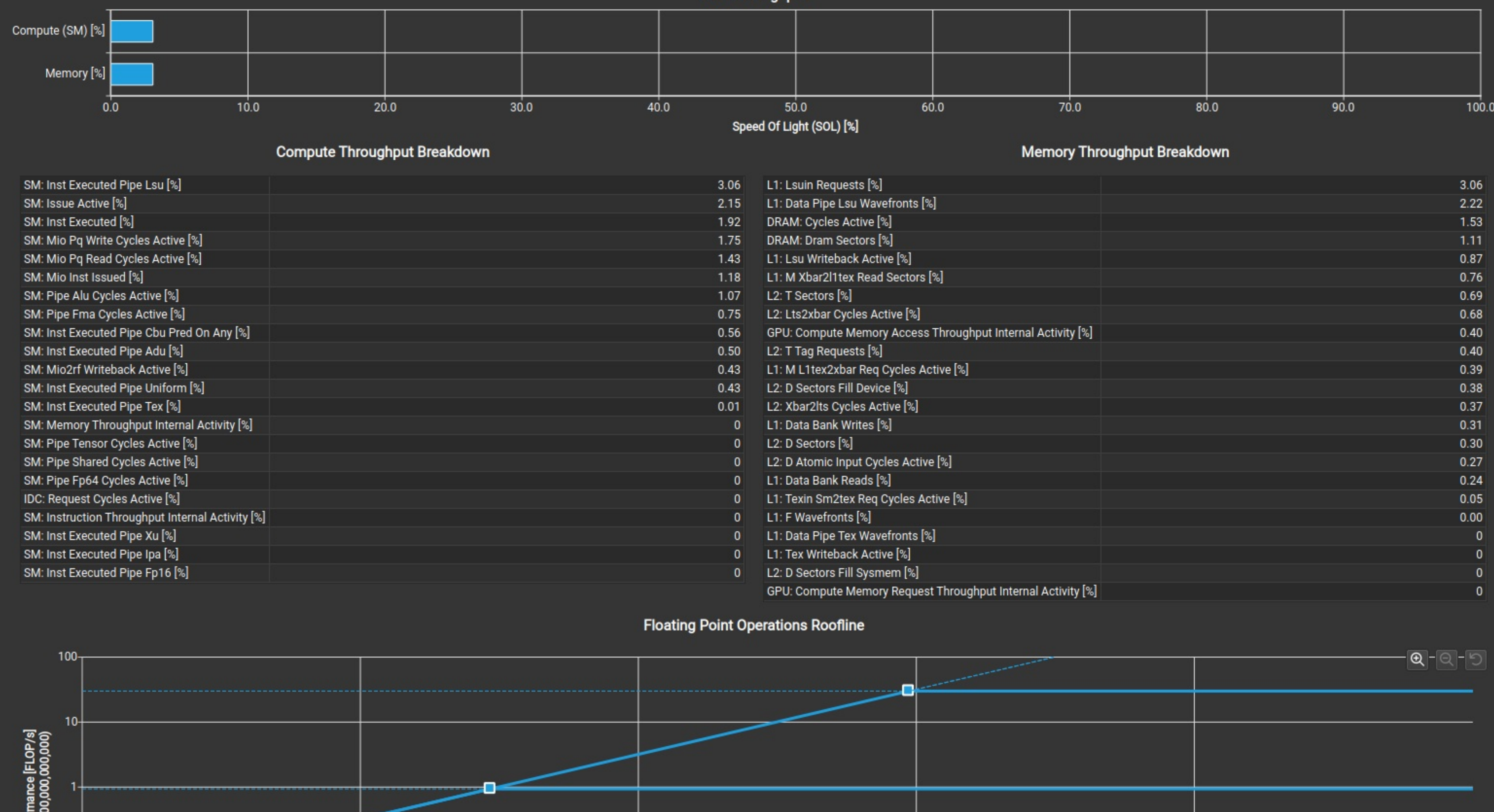
All

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 [%]	3.06	Duration [ms]	1.25
Memory Throughput [%]	3.06	Elapsed Cycles [cycle]	731,325
L1/TEX Cache Throughput [%]	4.45	SM Active Cycles [cycle]	712,002.62
L2 Cache Throughput [%]	0.69	SM Frequency [Mhz]	584.96
DRAM Throughput [%]	1.53	DRAM Frequency [Ghz]	4.99

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.

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.



▼ PM Sampling

Timeline view of PM metrics sampled periodically over the workload duration. Data is collected across multiple passes. Use this section to understand how workload behavior changes over its runtime.

Maximum Sampling Interval [cycle]	20,000	# Pass Groups	1
Maximum Buffer Size [Mbytes]	4	Dropped Samples [sample]	0



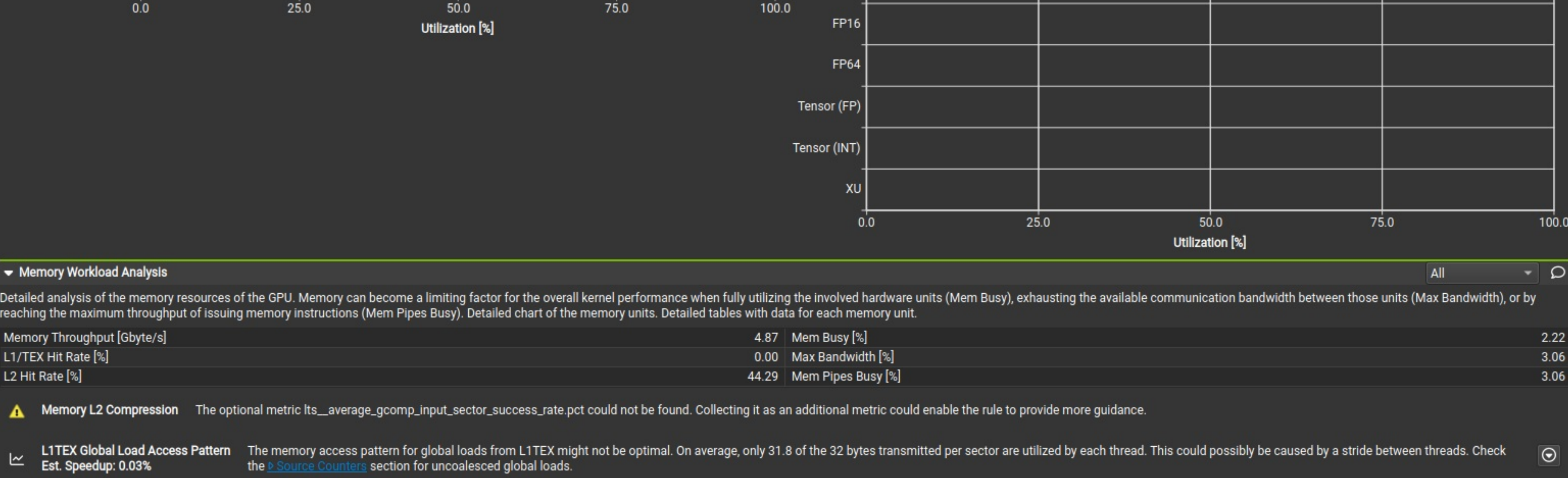
▼ Compute Workload Analysis

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]	0.08	SM Busy [%]	2.22
Executed Ipc Active [inst/cycle]	0.08	Issue Slots Busy [%]	2.22
Issued Ipc Active [inst/cycle]	0.09		

Low Utilization All compute pipelines are under-utilized. 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.

Est. Local Speedup: 98.90%



▼ 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.

Memory Throughput [Gbyte/s]	4.87	Mem Busy [%]	2.22
L1/TEX Hit Rate [%]	0.00	Max Bandwidth [%]	3.06
L2 Hit Rate [%]	44.29	Mem Pipes Busy [%]	3.06

Memory L2 Compression The optional metric `lvs_avg_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 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 Counters](#) section for uncoalesced global loads.

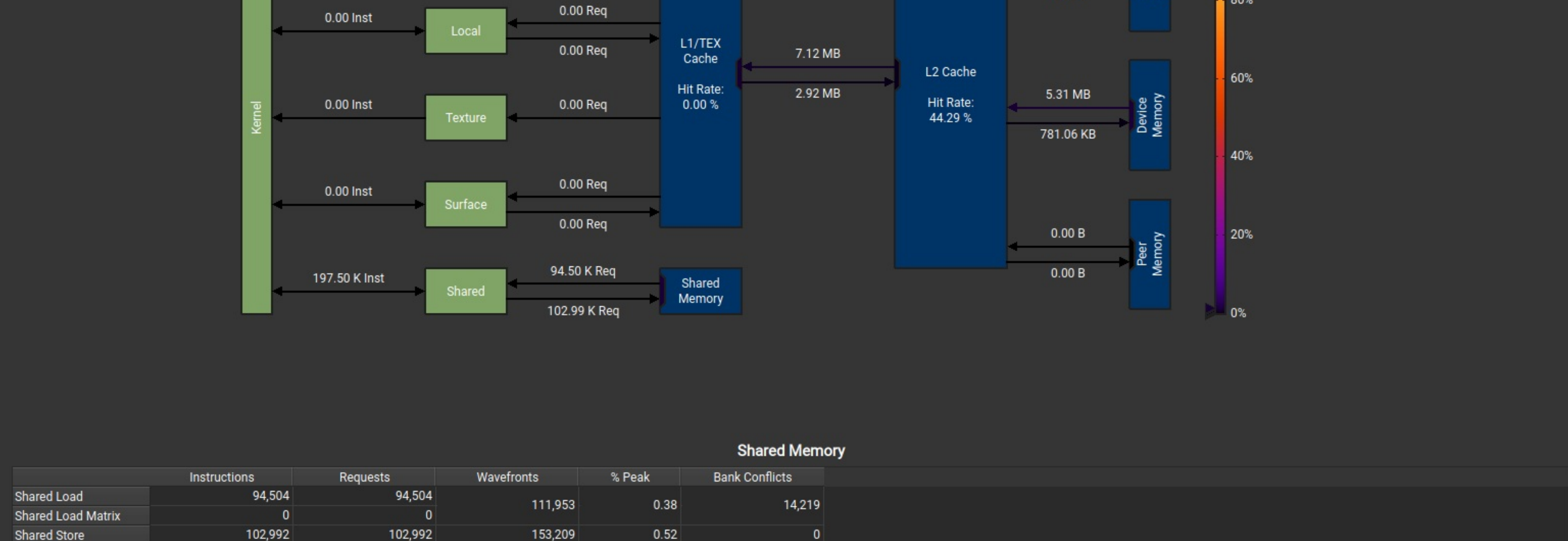
L1/TEX Global Store Access Pattern The memory access pattern for global stores to L1/TEX might not be optimal. On average, only 4.0 of the 32 bytes transmitted per sector are utilized by each thread. This could possibly be caused by a stride between threads. Check the [Source Counters](#) section for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

Shared Load Bank Conflicts 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 Counters](#) section for uncoalesced shared loads.

Est. Speedup: 0.56%

Memory Chart

Show As: Transfer Size



	Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	94,504	94,504	111,953	0.38	14,219
Shared Load Matrix	0	0	0	0	0
Shared Store	102,992	102,992	153,209	0.52	0
Shared Atomic	0	0	0	0	0
Other	-	-	139,393	0.48	0
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
Local Load	0	0	0	0	0	0	0	0	132,095	0.45	123,827
Global Load	33,792	33,792	33,792	0.12	132,096	3.91	0.00	4,227,072	0	0	0
Surface Load	0	0	0	0	0	0	0	0	0	0	0
Texture Load	0	0	0	0	0	0	0	0	0	0	0
Global Store	1,035	1,029	1,029	0.00	1,029	1	0.19	32,928	1,031	0.00	-
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	-
DSMEM Reduction	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	-
Global Atomic CAS	90,192	90,192	90,191	0.31	90,178	1.00	0	2,885,696	90,178	0.02	see above
Surface Atomic ALU	0	0	0	0	0	0	0	0	0	0	see above
Surface Atomic CAS	0	0	0	0	0	0	0	0	0	0	-
Loads	33,792	33,792	33,792	0.12	132,096	3.91	0.00	4,227,072	132,095	0.45	123,827
Stores	1,035	1,029	1,029	0.00	1,029	1	0.19	32,928	1,031	0.00	-
Atoms & Reductions	90,192	90,192	90,191	0.31	90,178	1.00	0	2,885,696	90,178	0.02	-

	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	3.91	0.39	0.77	4,227,040	3,381,069,390.05	131,072	0	0
L1/TEX Store	1,031	1,031	1	0.00	99.81	32,992	26,389,208.84	2	0	0
L1/TEX Atomic ALU	140	229	1.64	0.00	100	7,328	5,861,424.66	0	0	0
L1/TEX Atomic CAS	89,986	89,963	1.00	0.26	100.00	2,878,816	2,302,669,635.77	2	0	0
L1/TEX Reduction	0	0	0	0	0	0	0	0	0	0
L1/TEX Total	124,850	223,087	1.79	0.65	41.26	7,138,784	5,710,077,043.18	131,076	0	0
EOC Total	-	42	-	0.00	-	344	1,075,021.12	42	-	-
GPU Total	135,763	235,064	1.73	0.69	44.11	7,520,128	6,015,101,487.11	131,306	0	0

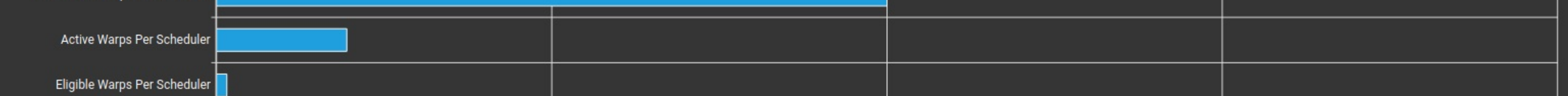
	Sectors	% Peak	Bytes	Throughput
Load	168,862	1.33	5,307,854	4,245,360,771.97
Store	24,468	0.20	791,055	625,740,842.12
Total	190,270	1.53	6,088,540	4,870,101,613.09

▼ 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 pipelines 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 the warp is skipped and no instruction is issued. Having many skipped issue slots indicates poor warp hiding.

Active Warps Per Scheduler [warp]	1.55	No Eligible Warps [%]	92.09
Eligible Warps Per Scheduler [warp]	0.12	One or More Eligible [%]	7.91
Issued Warp Per Scheduler	0.08		

Issue Slot Utilization 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 1.55 active warps per scheduler, but only an average of 0.12 warps are 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 Counters](#) sections.



▼ 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 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 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]	19.65	Avg. Active Threads Per Warp	18.55
Warp Cycles Per Executed Instruction [cycle]	22.03	Avg. Not Predicated Off Threads Per Warp	15.66

Long Scoreboard Stall 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 kernel configuration. Consider moving frequently used data to shared memory. This stall type represents about 50.4% of the total average of 19.7 cycles between issuing two instructions.

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 reduce 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 15.7 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.

► Instruction Statistics

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 instructions 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]	2,253,628	Avg. Executed Instructions Per Scheduler [inst]	14,085.17
Issued Instructions [inst]	2,526,420	Avg. Issued Instructions Per Scheduler [inst]	15,791.38

► NVLink Topology

NVLink Topology diagram shows logical NVLink connections with transmit/receive throughput.

► NVLink Tables

Detailed tables with properties for each NVLink.

► NUMA Affinity

Non-uniform memory access (NUMA) affinities based on compute and memory distances for all GPUs.

▼ Launch Statistics

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.

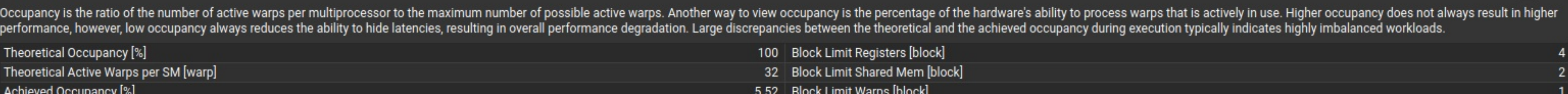
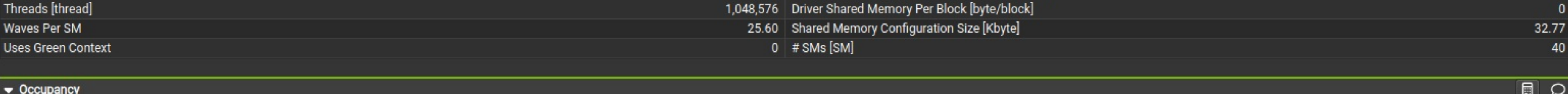
Grid Size	1,024	Function Cache Configuration	CachePreferNone
Registers Per Thread [register/thread]	16	Static Shared Memory Per Block [byte/block]	12,292
Block Size	1,024	Dynamic Shared Memory Per Block [byte/block]	0
Threads [thread]	1,048,576	Driver Shared Memory Per Block [byte/block]	0
Waves Per SM	25.60	Shared Memory Configuration Size [kbyte]	32.77
Uses Green Context	0	# SMs [SM]	40

▼ 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]	4
Theoretical Active Warps per SM [warp]	32	Block Limit Shared Mem [block]	2
Achieved Occupancy [%]	5.52	Block Limit Wmps [block]	1
Achieved Active Warps Per SM [warp]	1.77	Block Limit SM [block]	16

Achieved Occupancy The difference between calculated theoretical (100.0%) and measured achieved occupancy (5.5%) 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. See the [CUDA Best Practices Guide](#) for more details on optimizing occupancy.



▼ 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]	712,002.62	Average L1 Active Cycles [cycle]	712,002.62
Average L2 Active Cycles [cycle]	112,571	Average SMSP Active Cycles [cycle]	199,678.09
Average DRAM Active Cycles [cycle]	95,135	Total SM Elapsed Cycles [cycle]	29,323,856
Total L1 Elapsed Cycles [cycle]	112,571	Total L2 Elapsed Cycles [cycle]	34,203,456
Total SMSP Elapsed Cycles [cycle]	95,135	Total DRAM Elapsed Cycles [cycle]	49,689,280

SMSPs Workload Imbalance One or more SMSPs have a much lower number of active cycles than the average number of active cycles. Maximum instance value is 55.15% above the average, while the minimum instance value is 74.30% below the average.

L2 Slices Workload Imbalance One or more L2 Slices have a much higher number of active cycles than the average number of active cycles. Maximum instance value is 89.03% above the average, while the minimum instance value is 30.52% below the average.

Workload Distribution

	Average	Min	Max	Sum
SM Active Cycles	712,002.62	697,496	726,899	28,480,105
SMP Active Cycles	199,678.09	51,319	445,224	31,948,495
L1 Active Cycles	712,002.62	697,496	726,899	28,480,105
L2 Active Cycles	112,571	78,209	1,026,468	3,602,272
DRAM Active Cycles	95,135	94,240	96,024	701,080

► 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]	414,393	Branch Efficiency [%]	98.54
Branch Instructions Ratio [%]	0.18	Avg. Divergent Branches	32

Uncoalesced Shared Accesses This kernel has uncoalesced shared accesses resulting in a total of 30622 excessive wavefronts (12% of the total 265162 wavefronts). Check the L1 Wavefronts Shared Excessive table for the primary source locations. See the [CUDA Best Practices Guide](#) has an example on optimizing shared memory accesses.

L1 Wavefronts Shared Excessive

Location	Value	Value (%)
kernels.cu:207 (0x78ab3f038240 in reduce_argmin)	15,311	59
kernels.cu:207 (0x78ab3f038240 in reduce_argmin)	15,311	59
kernels.cu:215 (0x78ab3f038240 in reduce_argmin)	0	0
kernels.cu:215 (0x78ab3f038240 in reduce_argmin)	0	0
kernels.cu:206 (0x78ab3f038240 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 [unwanted sections](#) to focus on selected performance aspects and make profiling faster.