

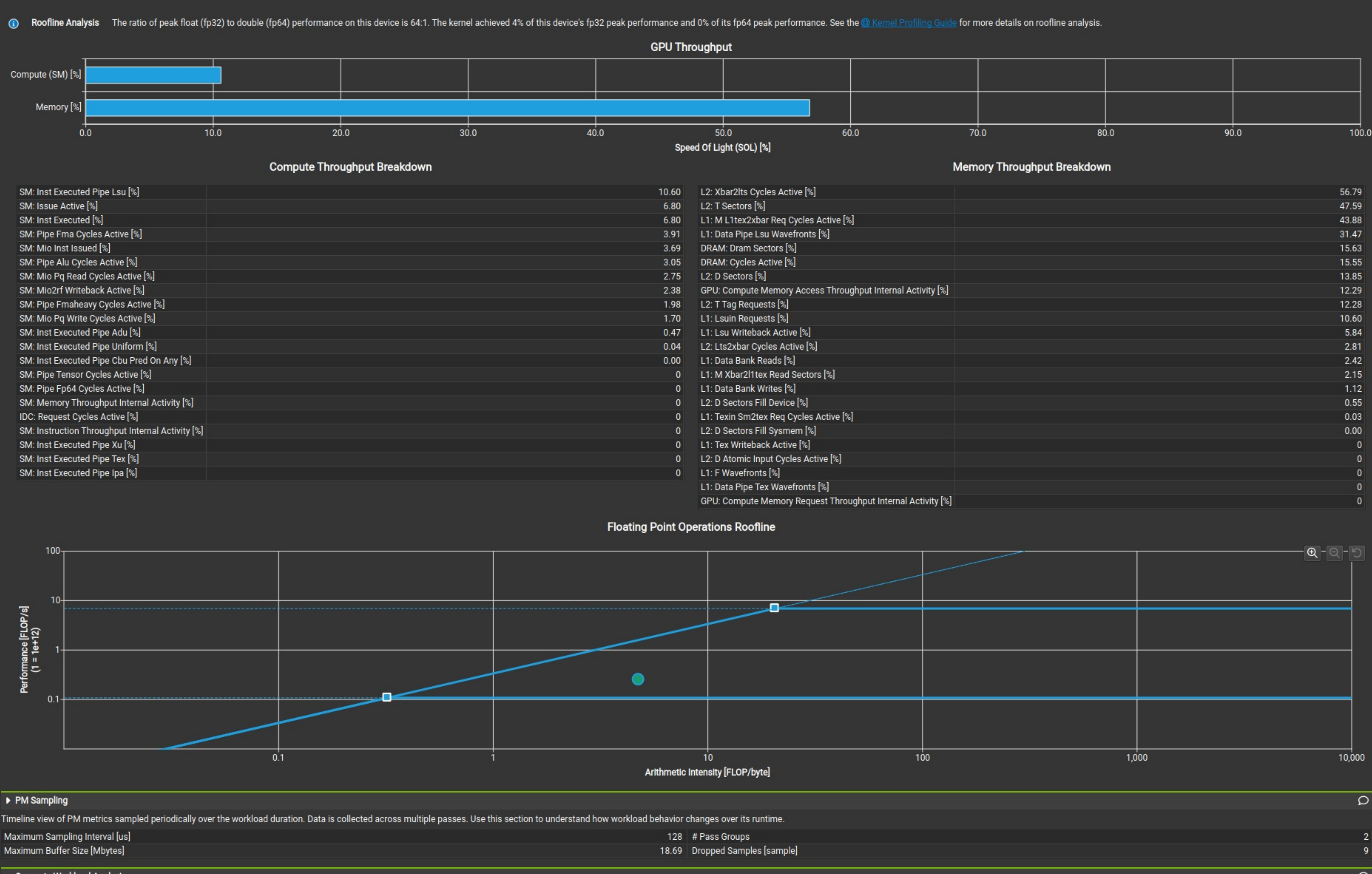
GPU Speed of Light Throughput

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 rooftop chart.

Compute (SM) Throughput [%]	10.60	Duration [ms]	656.22
Memory Throughput [%]	56.79	Elapsed Cycles [cycle]	590598959
L1/TEX Cache Throughput [%]	87.75	SM Active Cycles [cycle]	57725317.87
L2 Cache Throughput [%]	56.79	SM Frequency [MHz]	900.00
DRAM Throughput [%]	15.63	DRAM Frequency [GHz]	6.99

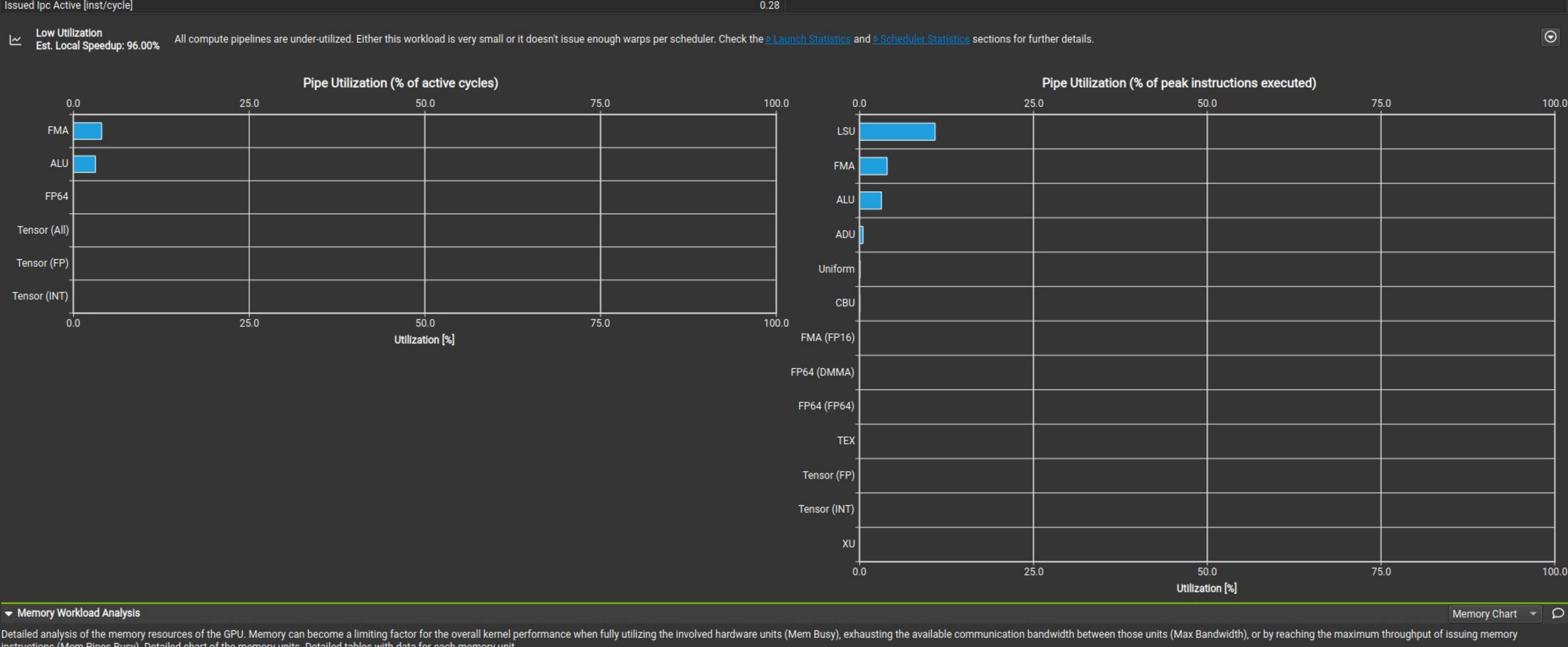
This kernel exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of this device. Achieved compute throughput and/or memory bandwidth below 60.0% of peak typically indicate latency issues. Look at the [Scheduler Statistics](#) and [Warp State Statistics](#) for potential reasons.

Roofline Analysis The ratio of peak float (fp32) to double (fp64) performance on this device is 64:1. The kernel achieved 4% 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.

Metric	Value
Maximum Sampling Interval [us]	128
Maximum Buffer Size [Mbytes]	18.69
Dropped Samples [sample]	9



Memory Workload Analysis

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.

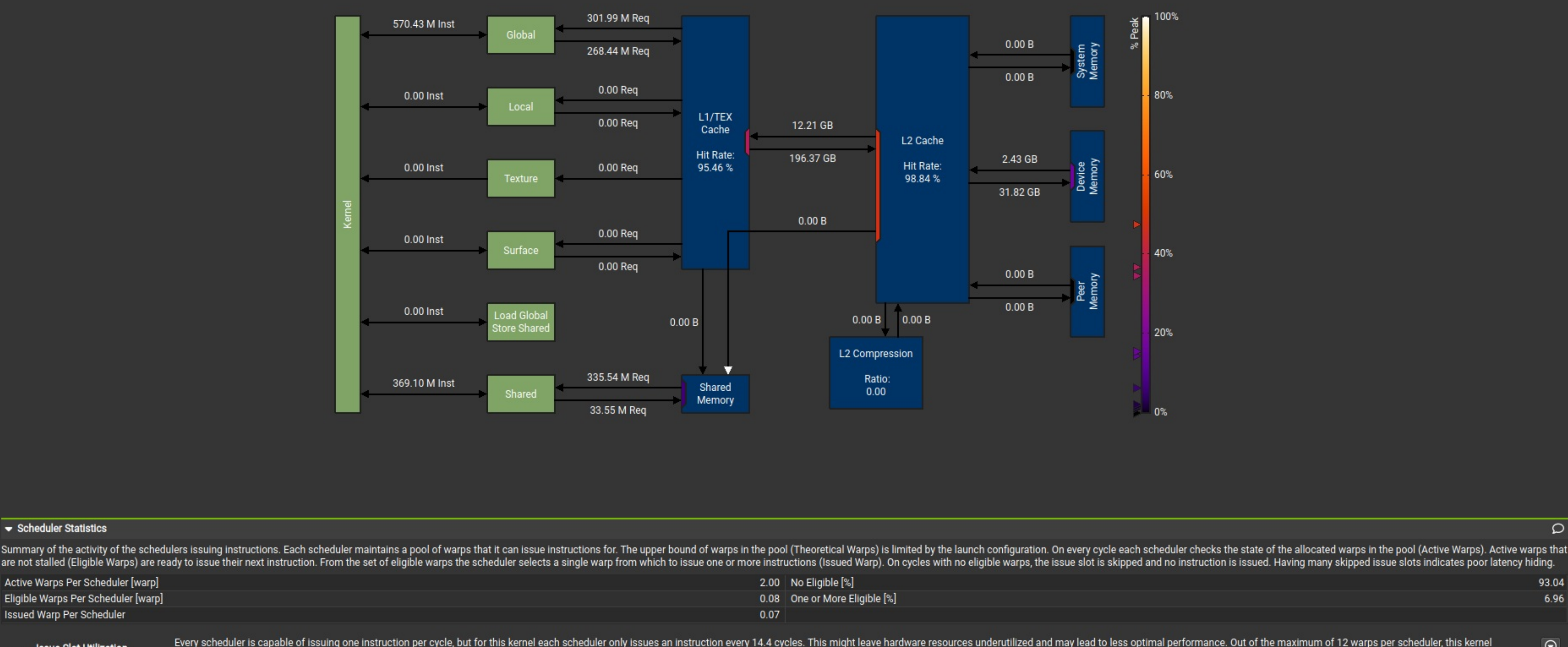
Metric	Value
Memory Throughput [Gbyte/s]	52.20
L1/TEX Hit Rate [%]	95.46
L2 Hit Rate [%]	98.84
L2 Compression Success Rate [%]	0

L1/TEX Global Load Access Pattern Est. Speedup: 75.60% The memory access pattern for global loads from L1/TEX might not be optimal. On average, only 4.4 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 Est. Speedup: 76.76% 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 Est. Speedup: 43.88% The memory access pattern for shared loads might not be optimal and causes on average a 3.2 - way bank conflict across all 335544320 shared load requests. This results in 536870912 bank conflicts, which represent 50.00% of the overall 1073741824 wavefronts for shared loads. Check the [Source Counters](#) section for uncoalesced shared loads.

Shared Store Bank Conflicts Est. Speedup: 12.46% The memory access pattern for shared stores might not be optimal and causes on average a 1.2 - way bank conflict across all 33554432 shared store requests. This results in 5544275 bank conflicts, which represent 14.20% of the overall 39041207 wavefronts for shared stores. Check the [Source Counters](#) section for uncoalesced shared stores.

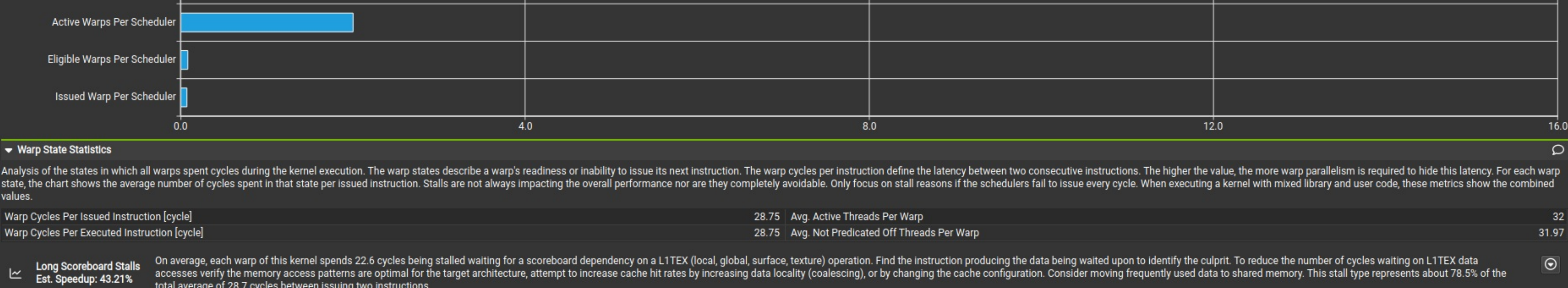


Scheduler Statistics

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 Warps). 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.

Metric	Value
Active Warps Per Scheduler [warp]	2.00
Eligible Warps Per Scheduler [warp]	0.08
Issued Warp Per Scheduler	0.07

Issue Slot Utilization Est. Local Speedup: 43.21% Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 14.4 cycles. This might leave hardware resources underutilized and may lead to less optimal performance. Out of the maximum of 12 warps per scheduler, this kernel allocates an average of 2.00 active warps per scheduler, but only an average of 0.08 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 [Source Counters](#) sections can help, too.

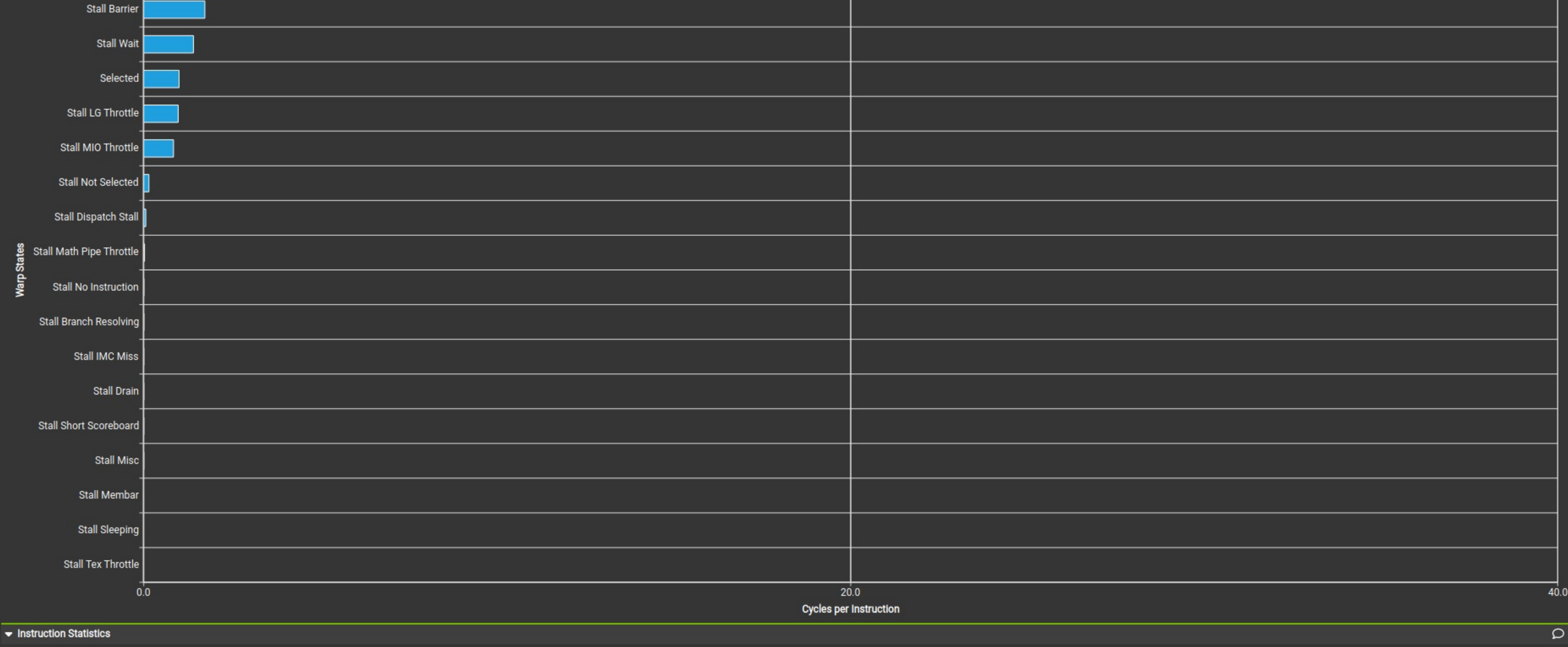


Warp State Statistics

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 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.

Metric	Value
Warp Cycles Per Issued Instruction [cycle]	28.75
Warp Cycles Per Executed Instruction [cycle]	28.75

Long Scoreboard Stall Est. Speedup: 43.21% On average, each warp of this kernel spends 22.6 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 78.5% of the total average of 28.7 cycles between issuing two instructions.



Instruction Statistics

Statistics of the executed low-level assembly instructions (SAS). 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.

Metric	Value
Executed Instructions [Inst]	4819732200
Issued Instructions [Inst]	4819772239

Executed Instruction Mix

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.

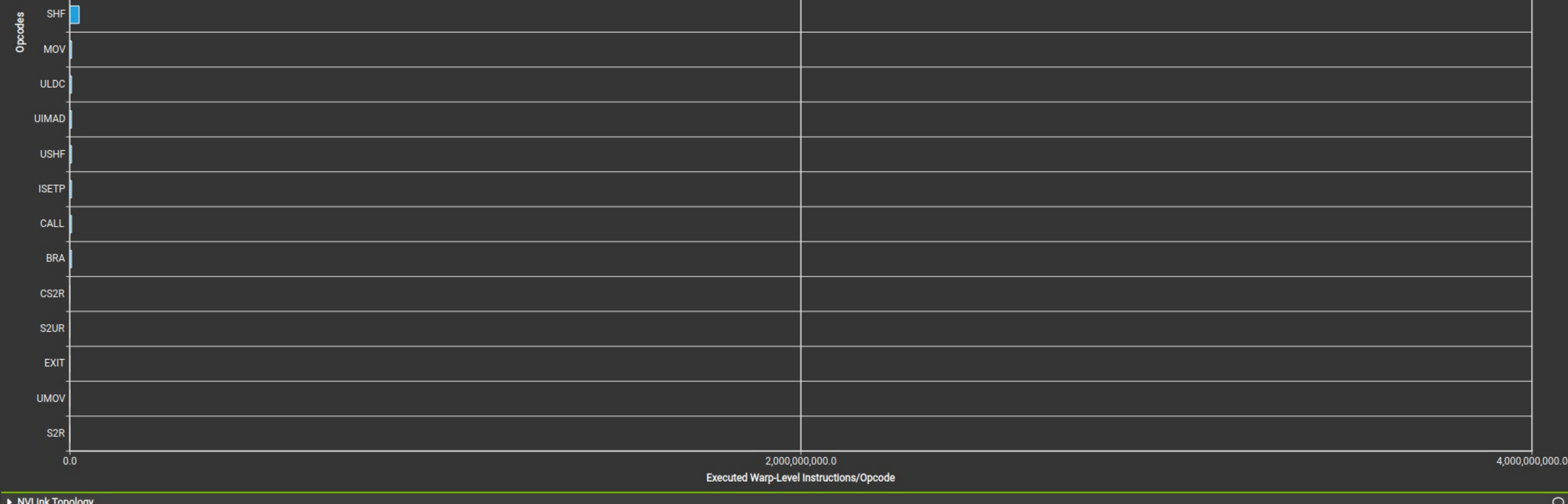
Metric	Value
Grid Size	1024
Registers Per Thread [register/thread]	218
Block Size	256
Threads Thread	262144
Waves Per SM	34.13
Uses Green Context	0

Occupancy

Occupancy is the ratio of the number of active warps per multiprocessor to the maximum number of possible 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.

Metric	Value
Theoretical Occupancy [%]	16.67
Average SM Active Cycles [cycle]	9
Achieved Occupancy [%]	16.66
Achieved Active Warps Per SM [warp]	8.00

Theoretical Occupancy Est. Speedup: 43.21% The 2.00 theoretical warps per scheduler this kernel can issue according to its occupancy are below the hardware maximum of 12. This kernel's theoretical occupancy (16.7%) is limited by the number of required registers. This kernel's theoretical occupancy (16.7%) is limited by the required amount of shared memory.



GPU and Memory Workload Distribution

Analysis of workload distribution in active cycles of SM, SMP, SMSP, L1 & L2 caches, and DRAM

Metric	Value
Average SM Active Cycles [cycle]	57725317.87
Average L2 Active Cycles [cycle]	54959980.17
Average DRAM Active Cycles [cycle]	713652565.33
Total L1 Elapsed Cycles [cycle]	17719378630
Total SMSP Elapsed Cycles [cycle]	70877514520

Workload Distribution

Metric	Average	Min	Max	Sum
SM Active Cycles	57725317.87	574445160	591126704	17317605536
SMSP Active Cycles	57718638.59	574492971	590864554	17328000999
L2 Active Cycles	57725317.87	574445160	591126704	17317605536
DRAM Active Cycles	713652565.33	548630405	552685089	13189199524

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.

Metric	Value
Branch Instructions [Inst]	8396800
Branch Instructions Ratio [%]	0.00

Uncoalesced Global Accesses

This kernel has uncoalesced global accesses resulting in a total of 1503238536 excessive sectors (87% of the total 1707296256 wavefronts). Check the L2 Theoretical Sectors Global Excessive table for the primary source locations. The [CUDA Programming Guide](#) has additional information on reducing uncoalesced device memory accesses.

Est. Speedup: 83.57%

L2 Theoretical Sectors Global Excessive

Location	Value	Value (%)
05_2d_block_tiling cu 31 (0x500d95070 in gemm2dBlockTiling) ▶	117,446,912	1
05_2d_block_tiling cu 31 (0x500d95060 in gemm2dBlockTiling) ▶	117,446,912	1
05_2d_block_tiling cu 31 (0x500d95050 in gemm2dBlockTiling) ▶	117,446,912	1
05_2d_block_tiling cu 31 (0x500d95040 in gemm2dBlockTiling) ▶	117,446,912	1
05_2d_block_tiling cu 31 (0x500d95030 in gemm2dBlockTiling) ▶	117,446,912	1

Uncoalesced Shared Accesses

This kernel has uncoalesced shared accesses resulting in a total of 536870912 excessive wavefronts (48% of the total 1107296256 wavefronts). Check the L1 Wavefronts Shared Excessive table for the primary source locations. The [CUDA Best Practices Guide](#) has an example on optimizing shared memory accesses.

Est. Speedup: 47.39%

L1 Wavefronts Shared Excessive

Location	Value	Value (%)
05_2d_block_tiling cu 72 (0x500d950d0 in gemm2dBlockTiling) ▶	16,777,216	0
05_2d_block_tiling cu 72 (0x500d950c0 in gemm2dBlockTiling) ▶	16,777,216	0
05_2d_block_tiling cu 72 (0x500d950b0 in gemm2dBlockTiling) ▶	16,777,216	0
05_2d_block_tiling cu 72 (0x500d950a0 in gemm2dBlockTiling) ▶	16,777,216	0
05_2d_block_tiling cu 72 (0x500d95090 in gemm2dBlockTiling) ▶	16,777,216	0

Warp Stall Sampling (All Samples)

Location	Value	Value (%)
05_2d_block_tiling cu 50 (0x500d95090 in gemm2d) ▶	137,458	0
05_2d_block_tiling cu 50 (0x500d95080 in gemm2d) ▶	61,800	0
05_2d_block_tiling cu 50 (0x500d95070 in gemm2d) ▶	6,728	0
05_2d_block_tiling cu 50 (0x500d95060 in gemm2d) ▶	60,580	0
05_2d_block_tiling cu 50 (0x500d95050 in gemm2d) ▶	60,458	0

Most Instructions Executed

Location	Value	Value (%)
05_2d_block_tiling cu 50 (0x500d95090 in gemm2d) ▶	4,194,308	0
05_2d_block_tiling cu 50 (0x500d95080 in gemm2d) ▶	4,194,308	0
05_2d_block_tiling cu 50 (0x500d95070 in gemm2d) ▶	4,194,308	0
05_2d_block_tiling cu 50 (0x500d95060 in gemm2d) ▶	4,194,308	0
05_2d_block_tiling cu 50 (0x500d95050 in gemm2d) ▶	4,194,308	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 enable [Source Counters](#) to focus on selected performance aspects and make profiling faster.