

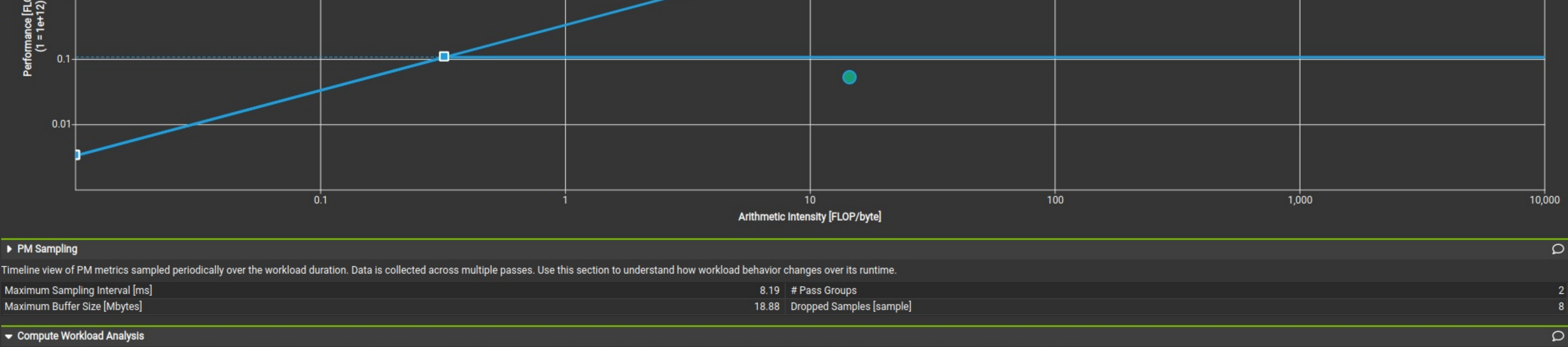
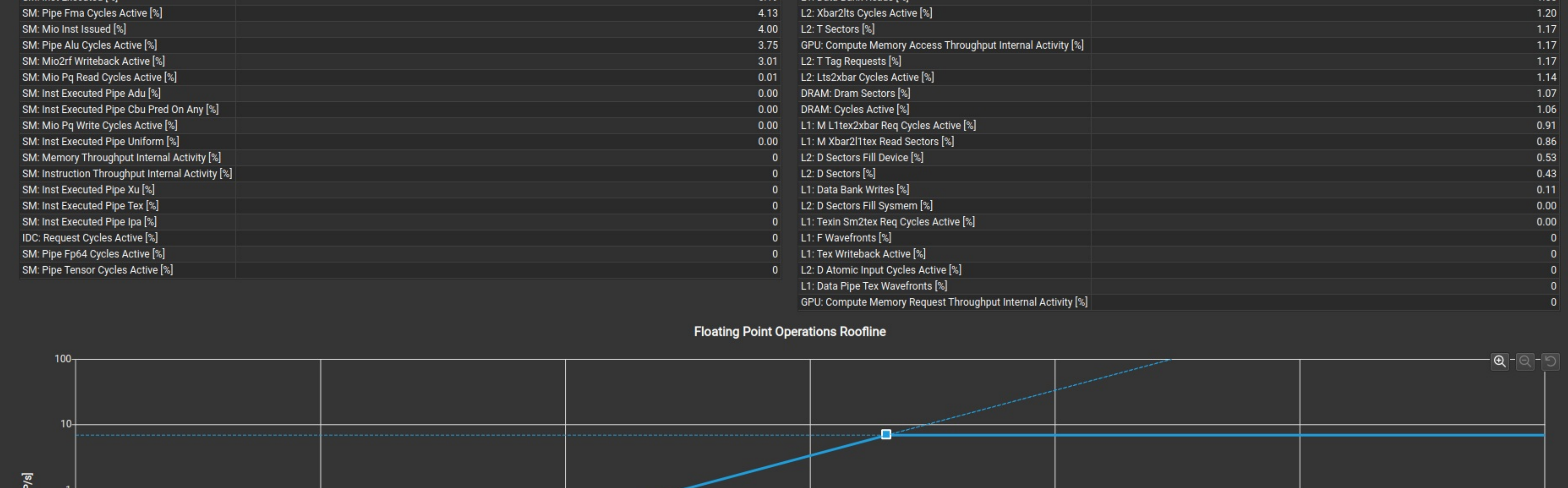
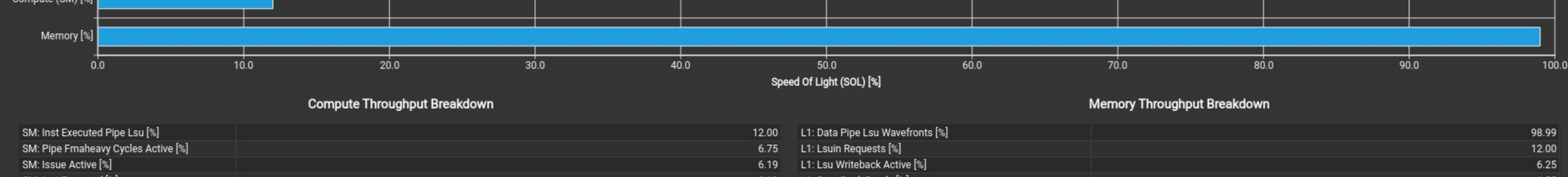
GPU Speed Of Light Throughput

All

Highest 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 [%]	12.00	Duration [s]	2.65
Memory Throughput [%]	98.99	Elapsed Cycles [cycle]	2387515271
L1/L2 Cache Throughput [%]	99.11	SM Active Cycles [cycle]	2384490716.17
L2 Cache Throughput [%]	1.20	SM Frequency [MHz]	900.00
DRAM Throughput [%]	1.07	DRAM Frequency [GHz]	6.99

High Throughput The kernel is utilizing greater than 80.0% of the available compute or memory performance of the device. To further improve performance, work will likely need to be shifted from the most utilized to another unit. Start by analyzing L1 in the [Memory Archival Analysis](#) section.



PM Sampling

All

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.

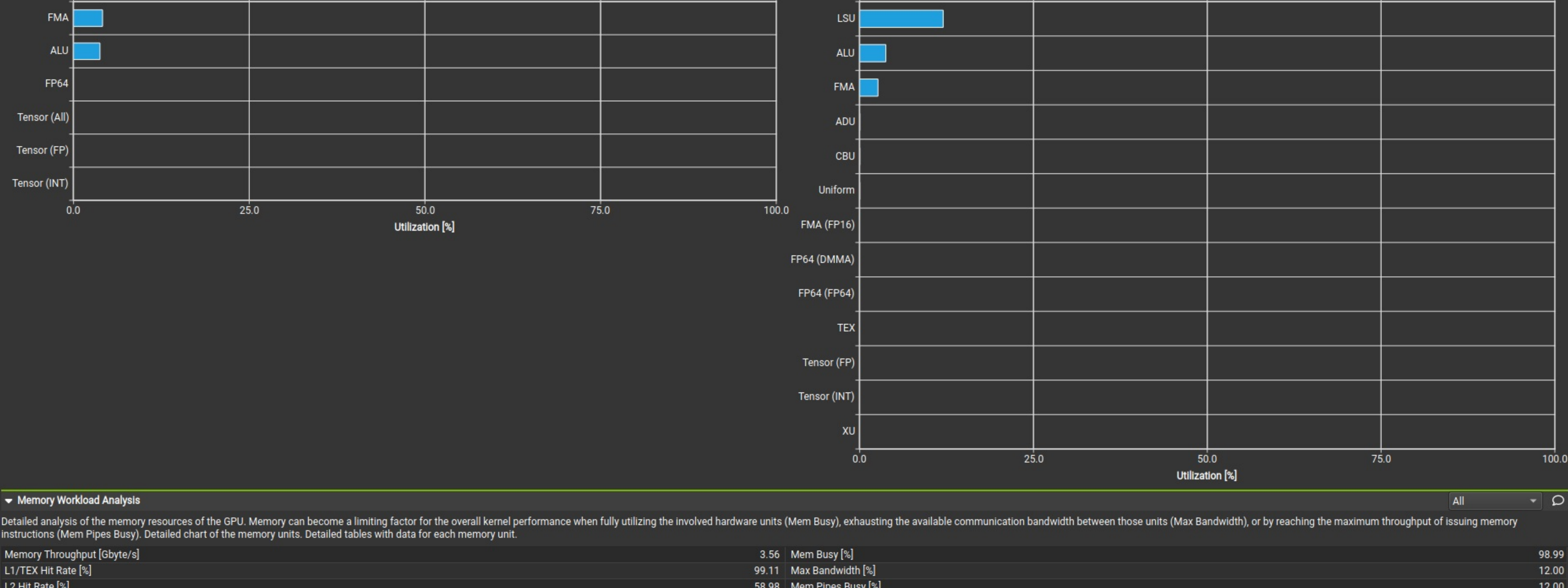
Maximum Sampling Interval [ms]	8.19	# Pass Groups	2
Maximum Buffer Size [Mbytes]	18.88	Dropped Samples [sample]	8

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]	0.25	SM Busy [%]	6.76
Executed ipc Active [inst/cycle]	0.25	Issue Slots Busy [%]	6.20
Issued ipc Active [inst/cycle]	0.25		

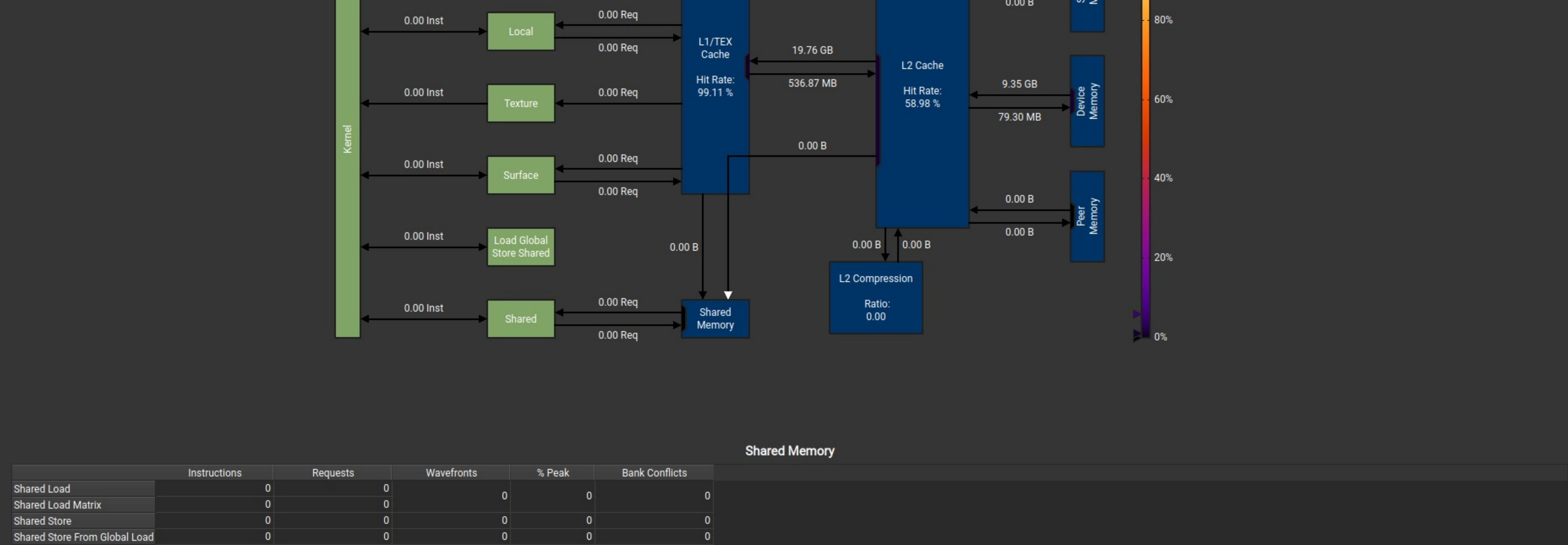


Memory Workload Analysis

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]	3.56	Mem Busy [%]	98.99
L1/TEX Hit Rate [%]	99.11	Max Bandwidth [%]	12.00
L2 Hit Rate [%]	58.98	Mem Pipes Busy [%]	12.00
L2 Compression Success Rate [%]	0	L2 Compression Ratio	0



Shared Memory

	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 Store From Global Load	0	0	0	0	0
Shared Atomic	0	-	1064960	0.00	0
Other	-	-	-	-	-
Total	-	-	1064960	0.00	0

L1/TEX Cache

	Instructions	Requests	Wavefronts	% Peak	Sectors	Sectors/Req	Hit Rate	Bytes	0	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	0
Global Load	4295491584	4295491584	0	0	70883373390	16.50	99.13	2268267948480	617100240	0.86	4473239771	6.25	
Global Load To Shared Store (access)	0	0	19331629032	26.99	0	0	0	0	0	0	0	0	0
Global Load To Shared Store (bypass)	0	0	0	0	0	0	0	0	0	0	0	0	0
Surface Load	0	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	0
Global Store	524288	524288	4243116	0.01	16777216	32	1.08	536870912	16777216	0.02			
Local Store	0	0	0	0	0	0	0	0	0	0	0	0	0
Surface Store	0	0	0	0	0	0	0	0	0	0	0	0	0
Global Reduction	0	0	0	0	0	0	0	0	0	0	0	0	0
DSMCM Reduction	0	0	0	0	0	0	0	0	0	0	0	0	0
Surface Reduction	0	0	0	0	0	0	0	0	0	0	0	0	0
Global Atomic ALU	0	0	0	0	0	0	0	0	0	0	0	0	0
Global Atomic CAS	0	0	0	0	0	0	0	0	0	0	0	0	0
Surface Atomic ALU	0	0	0	0	0	0	0	0	0	0	0	0	0
Surface Atomic CAS	0	0	0	0	0	0	0	0	0	0	0	0	0
Loads	4295491584	4295491584	19331629032	26.99	70883373390	16.50	99.13	2268267948480	617100240	0.86	4473239771	6.25	
Stores	524288	524288	4243116	0.01	16777216	32	1.08	536870912	16777216	0.02			
Atomic & Reductions	0	0	0	0	0	0	0	0	0	0	0	0	0
Total	4296015872	4296015872	19335872148	27.00	70900150606	16.50	99.11	2268804819392	633877456	0.89	4473239771	6.25	

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	617325381	617342272	1.00	1.12	52.85	19754952704	7446845336.42	290818375	0	0
L1/TEX Store	24102091	16777216	0.70	0.03	100	536870912	202379358.09	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	633322810	633974587	1.00	1.15	54.05	20283968784	7648270515.03	290885939	0	0
ECG Total	0	0	-	-	-	0	0	0	0	0
GPU Total	642200911	643764739	1.00	1.17	56.51	20600471648	7765572943.58	292156557	40927	0

L2 Cache Eviction Policies

	First	Hit Rate	Last	Hit Rate	Normal	Hit Rate	Normal Demote	Hit Rate
L1/TEX Load	0	0	0	0	617230996	52.80	0	0
L1/TEX Store	0	0	0	0	16777216	100	0	0
L1/TEX Atomic	0	0	0	0	0	0	0	0
L1/TEX Total	0	0	0	0	633798132	54.07	0	0
GPU Total	9960216	88.46	0	0	633729702	54.08	0	0

Device Memory

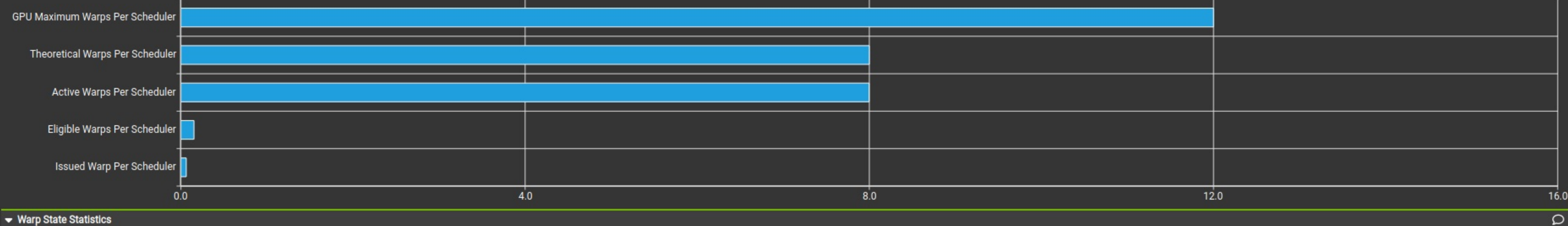
	Sectors	% Peak	Bytes	Throughput
Load	292232640	1.05	9301540480	3525165392.96
Store	2478252	0.01	79304064	29894533.69
Total	294713892	1.06	9430844544	355509926.66

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 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]	6.00	No Eligible [%]	92.80
Eligible Warps Per Scheduler [warp]	0.15	One or More Eligible [%]	6.20
Issued Warp Per Scheduler	0.06		

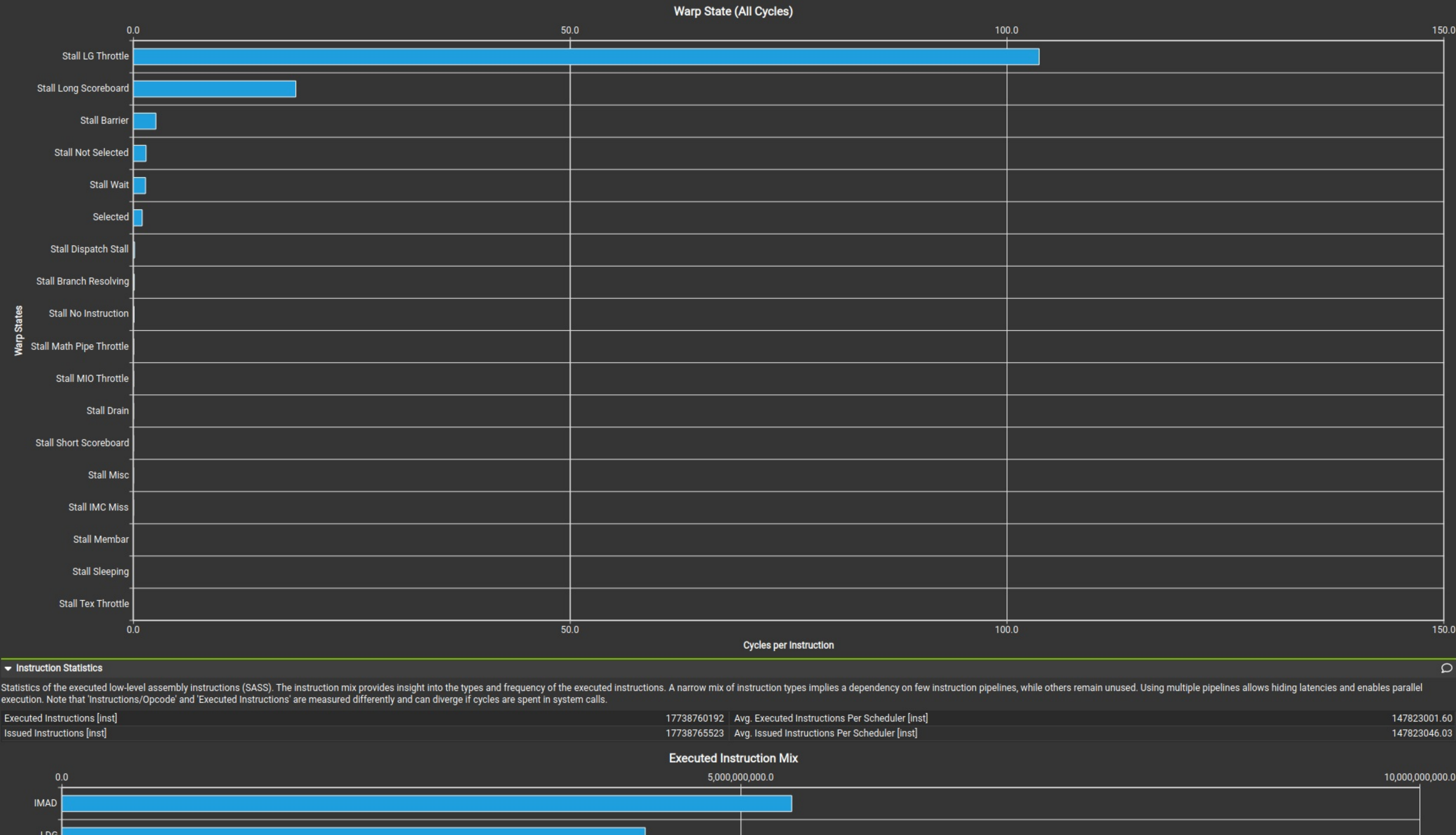


Warp State Statistics

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 Executed Instruction [cycle]	129.01	Avg. Active Threads Per Warp	32
Warp Cycles Per Issued Instruction [cycle]	129.01	Avg. Not Predicated Off Threads Per Warp	32.00

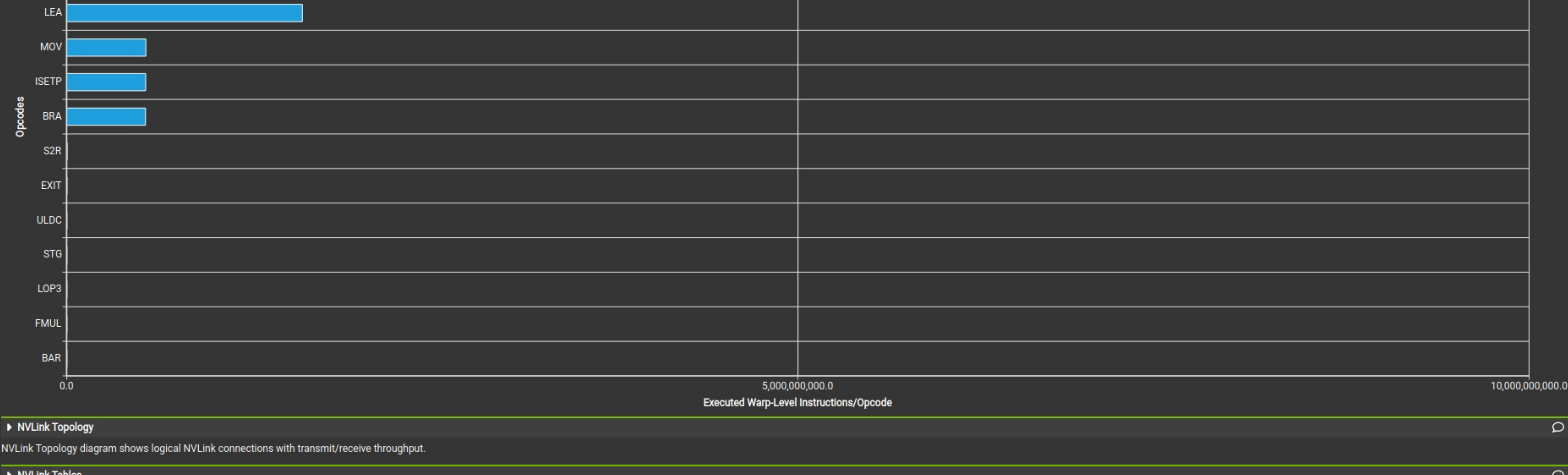


Instruction Statistics

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]	1773876192	Avg. Executed Instructions Per Scheduler [inst]	1478223901.66
Issued Instructions [inst]	1773876523	Avg. Issued Instructions Per Scheduler [inst]	1478223946.02



NVLink Topology

All

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

NVLink Tables

All

Detailed tables with properties for each NVLink.

NUMA Affinity

All

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

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	16384	Function Cache Configuration	CachePreference
Registers Per Thread [register/thread]	36	Static Shared Memory Per Block [byte/block]	None
Block Size	1024	Dynamic Shared Memory Per Block [byte/block]	2384498114.66
Threads Thread	16777216	Driver Shared Memory Per Block [kbyte/block]	1.02
Waves Per SM	546.13	Shared Memory Configuration Size [kbyte]	8.19
Uses Green Context	0	# SMs [SM]	30

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 [%]	66.67	Block Limit Registers [block]	1
Theoretical Active Warps per SM [warp]	1024	Block Limit Shared Memory [block]	8
Achieved Occupancy [%]	66.65	Block Limit Warps [block]	1
Achieved Active Warps Per SM [warp]	31.99	Block Limit SM [block]	16

The 8.00 theoretical warps per scheduler this kernel can issue according to its occupancy are below the hardware maximum of 12. This kernel's theoretical occupancy (66.7%) is limited by the number of required registers. This kernel's theoretical occupancy (66.7%) is limited by the number of warps within each block.



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 [Interactive Analysis](#) to focus on selected performance aspects and make profiling faster.