

GPU Speed of Light Throughput

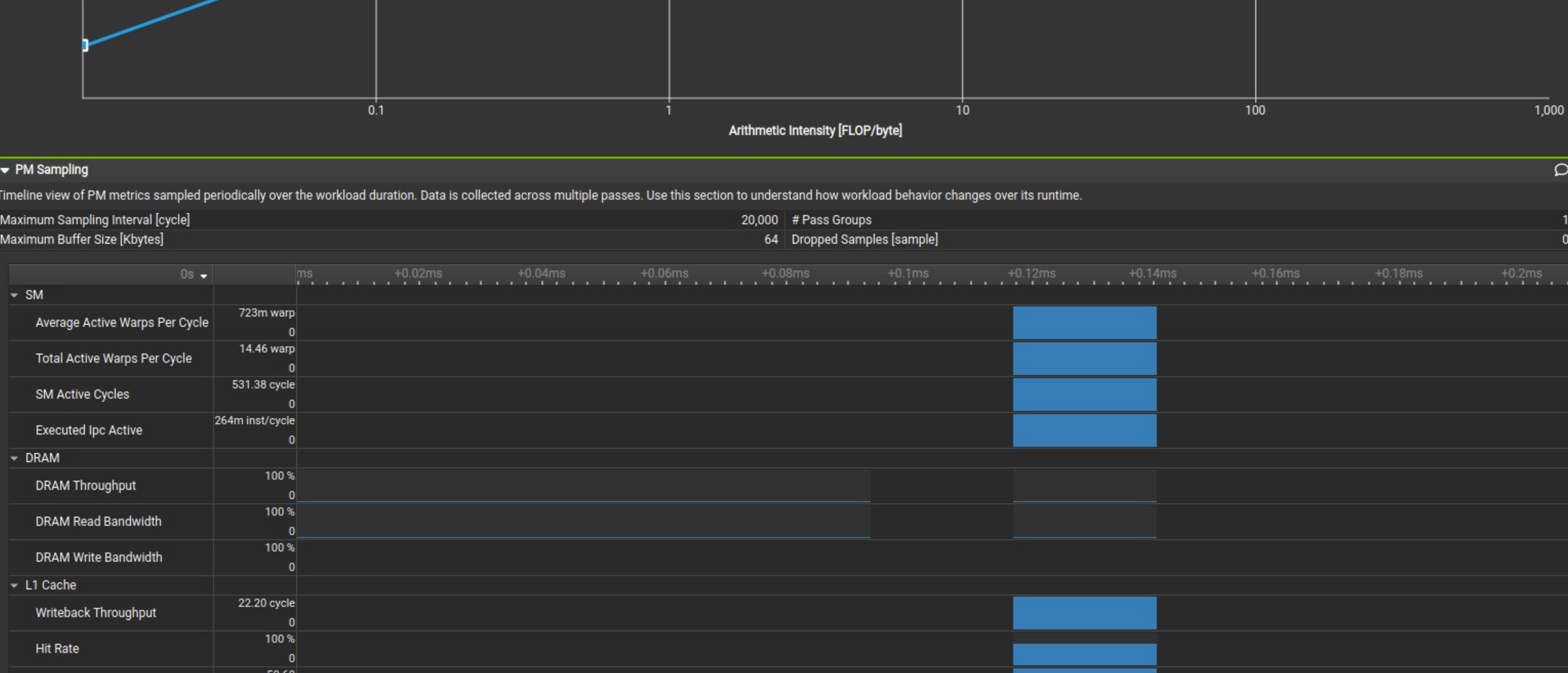
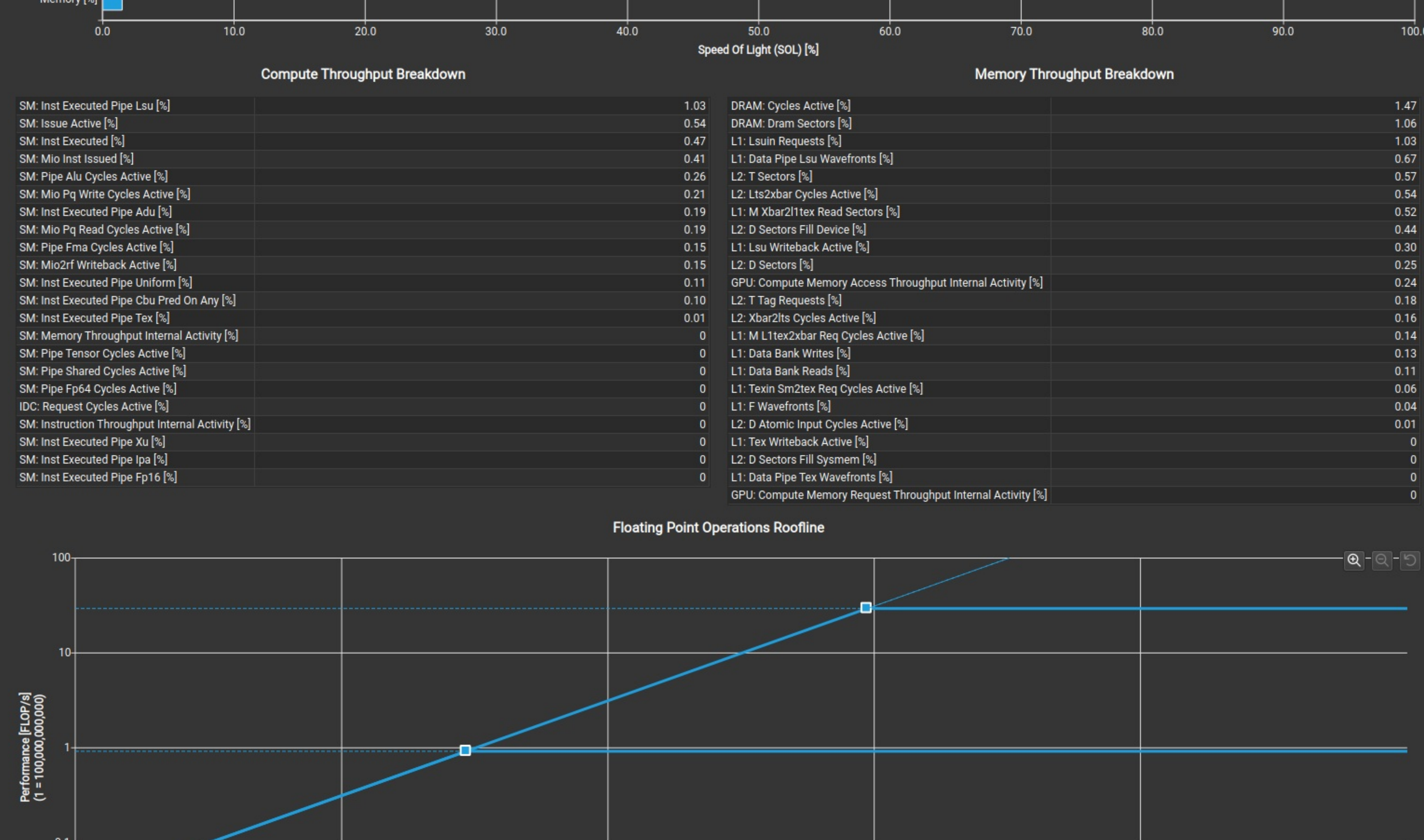
All

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 (SM) Throughput [%]	1.03	Duration [%]	13.09
Memory Throughput [%]	1.47	Elapsed Cycles [cycle]	7,622
L1/TEX Cache Throughput [%]	14.98	SM Active Cycles [cycle]	516.67
L2 Cache Throughput [%]	0.57	SM Frequency [Mhz]	582.32
DRAM Throughput [%]	1.47	DRAM Frequency [Ghz]	4.89

☒ **Small Grid** This kernel grid is too small to fill the available resources on this device, resulting in only 0.1 full waves across all SMs. Look at [Launch Statistics](#) for more details.

☒ **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 [kbytes]	64	Dropped Samples [sample]	0

SM

DRAM

L1 Cache

Workload Execution

723m warp

14.46 warp

531.38 cycle

264m inst/cycle

100 %

100 %

100 %

22.20 cycle

100 %

50.60

0

0ms

+0.02ms

+0.04ms

+0.06ms

+0.08ms

+0.1ms

+0.12ms

+0.14ms

+0.16ms

+0.18ms

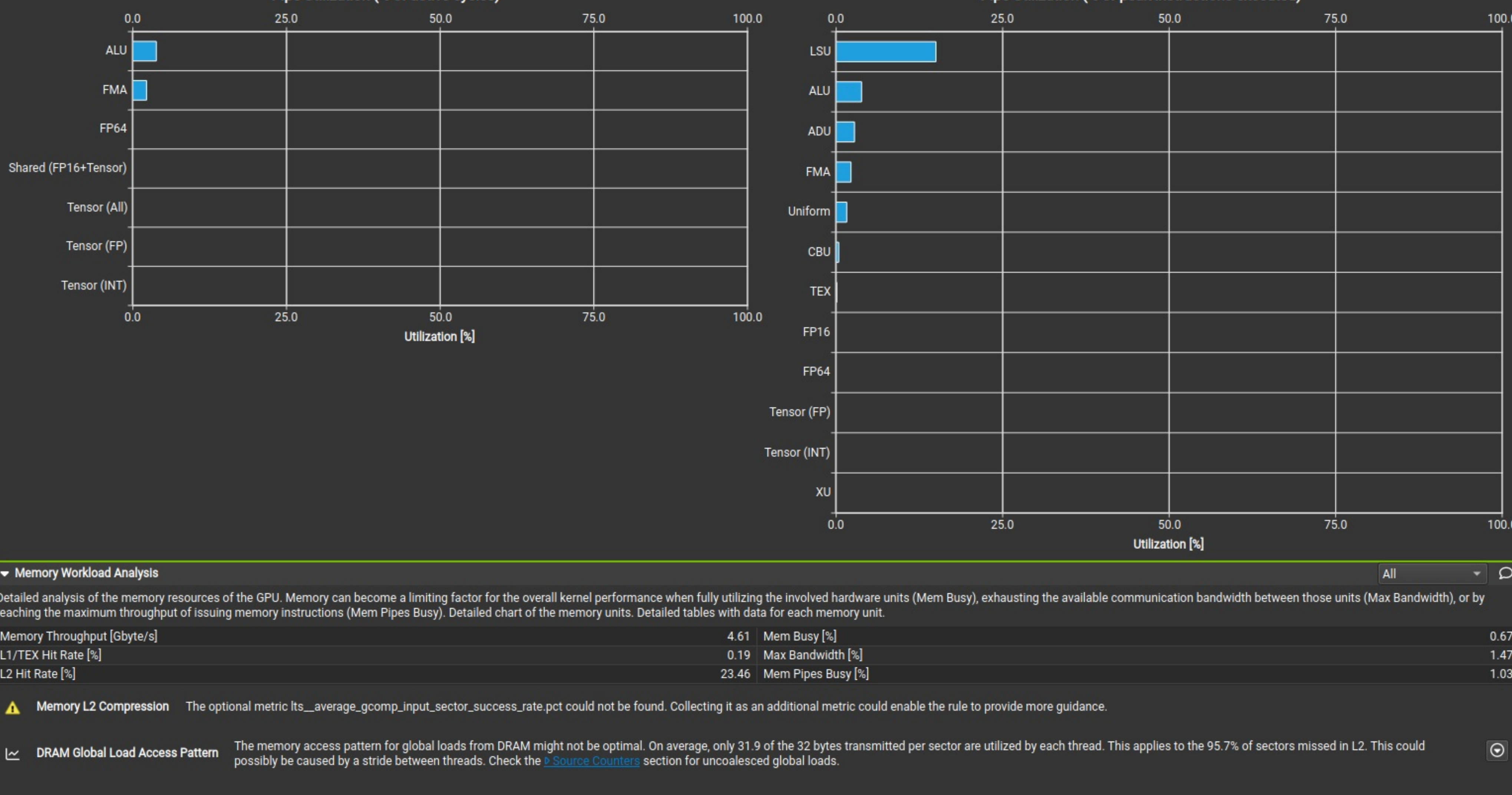
+0.2ms

red-box_region

☒ **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 [pc Elapsed] [inst/cycle]	0.02	SM Busy [%]	7.80
Executed [pc Active] [inst/cycle]	0.19	Issue Slots Busy [%]	7.80
Issued [pc Active] [inst/cycle]	0.31		

☒ **Low Utilization** **Est. Local Speedup: 96.19%** 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.



☒ **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 tables of the memory units. Detailed tables with data for each memory unit.

Memory Throughput [Gbyte/s]	4.61	Mem Busy [%]	0.57
L1/TEX Hit Rate [%]	0.19	Max Bandwidth [%]	1.47
L2 Hit Rate [%]	23.46	Mem Pipes Busy [%]	1.03

☒ **Memory L2 Compression** The optional metric `its_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.

☒ **DRAM Global Load Access Pattern** The memory access pattern for global loads from DRAM might not be optimal. On average, only 31.9 of the 32 bytes transmitted per sector are utilized by each thread. This applies to the 95.7% of sectors missed in L2. This could possibly be caused by a stride between threads. Check the [Schedul Statistics](#) 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 4.6 of the 32 bytes transmitted per sector are utilized by each thread. This could possibly be caused by a stride between threads. Check the [Launch Statistics](#) section for uncoalesced global stores.

Memory Chart

Show As: Transfer Size

