

GPU Speed of Light Throughput

All

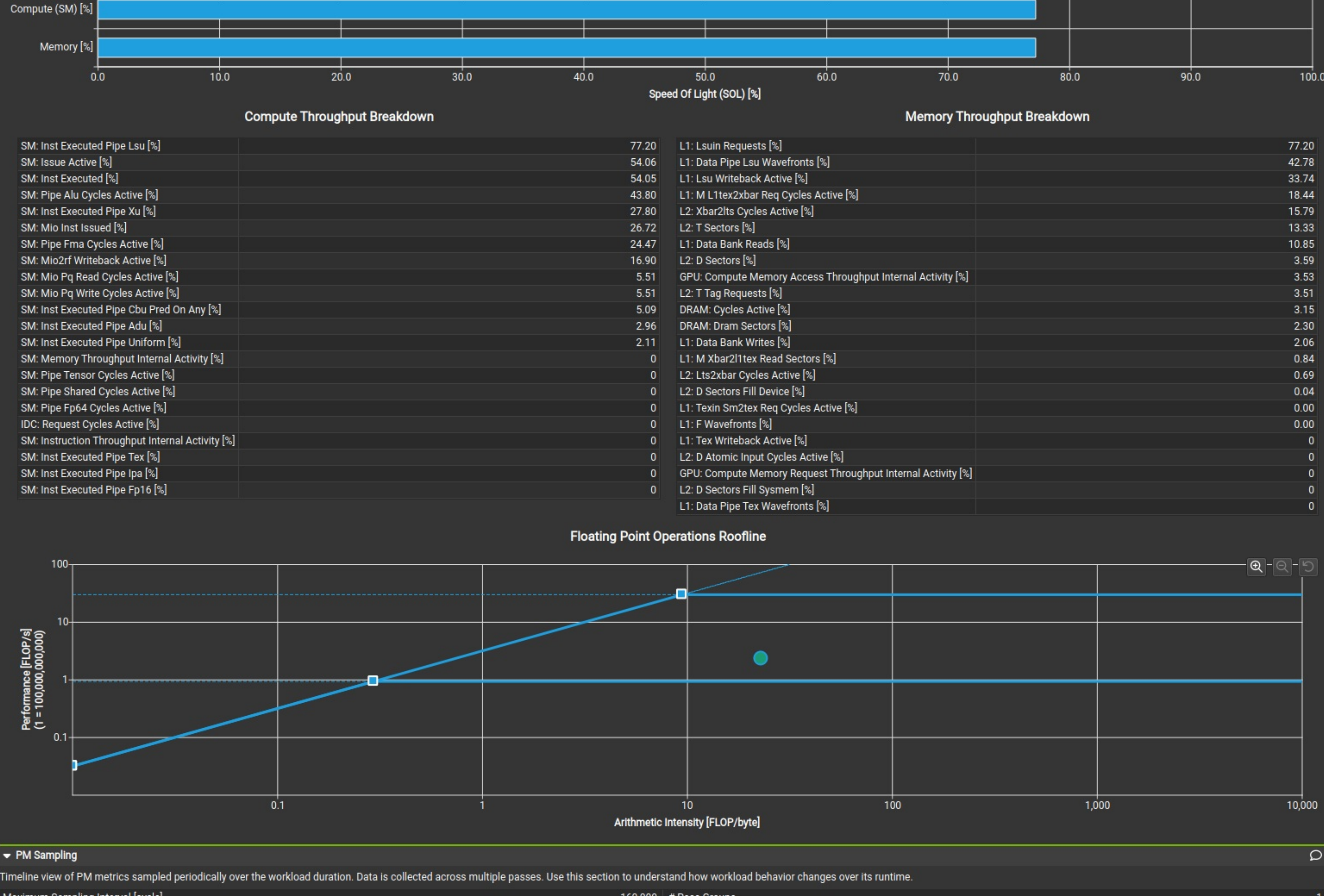
High-level overview of the throughput for compute and memory resources of the GPU. For each unit, the throughput reports the achieved percentage of utilization with respect to the theoretical maximum. Breakdowns show the throughput for each individual sub-metric of Compute and Memory to clearly identify the highest contributor. High-level overview of the utilization for compute and memory resources of the GPU presented as a routine chart.

Compute (SM) Throughput [%]	77.20	Duration [ms]	2.07
Memory Throughput [%]	77.20	Elapsed Cycles [cycle]	1,211,008
L1/TEX Cache Throughput [%]	85.88	SM Active Cycles [cycle]	1,181,089.43
L2 Cache Throughput [%]	15.79	SM Frequency [Mhz]	584.98
DRAM Throughput [%]	3.15	DRAM Frequency [Ghz]	4.99

- Balanced Throughput

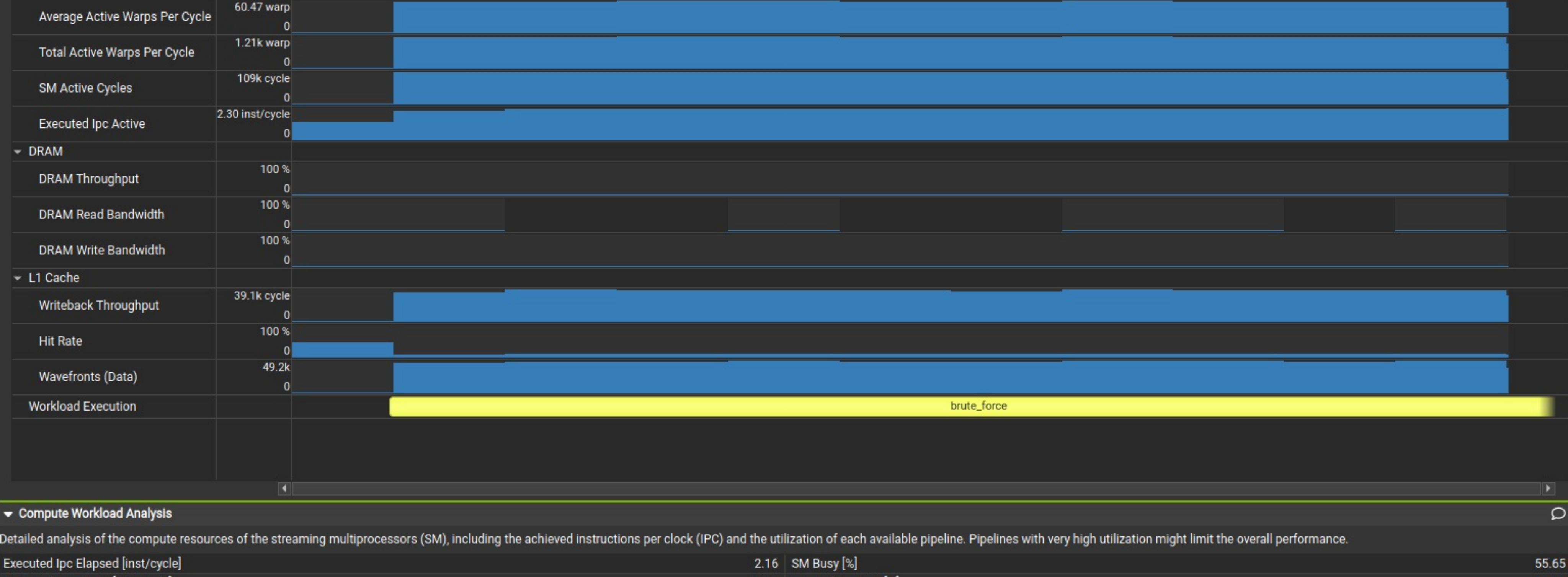
Compute and Memory are well-balanced. To reduce runtime, both computation and memory traffic must be reduced. Check both the [Compute Workload Analysis](#) and [Memory Workload Analysis](#) sections.
- Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 8% 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.



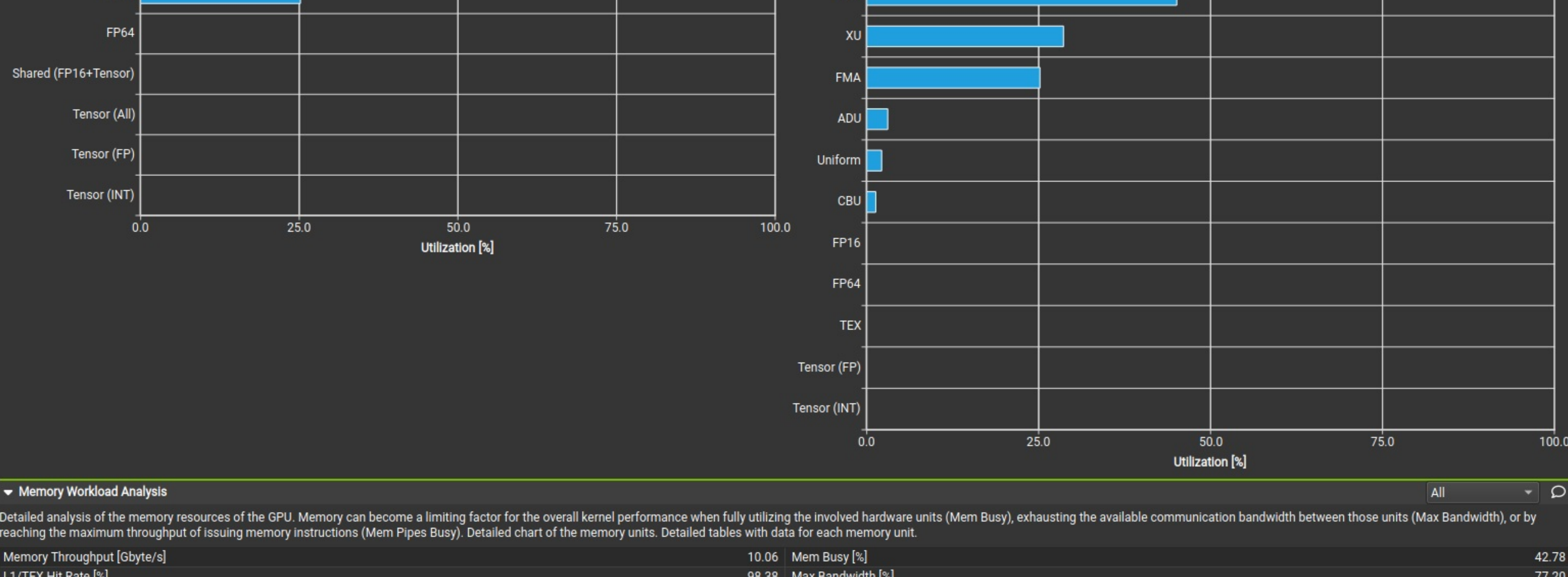
Compute Workload Analysis

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]	2.16	SM Busy [%]	55.65
Executed Ipc Active [inst/cycle]	2.23	Issue Slots Busy [%]	55.65
Issued Ipc Active [inst/cycle]	2.23		

- Balanced

ALU is the highest-utilized pipeline (45.1%) based on active cycles, taking into account the rates of its different instructions. It executes integer and logic operations. It is well-utilized, but should not be a bottleneck.



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.

Memory Throughput [Gbytes/s]	10.06	Mem Busy [%]	42.78
L1/TEX Hit Rate [%]	98.38	Max Bandwidth [%]	77.20
L2 Hit Rate [%]	100.20	Mem Pipes Busy [%]	77.20

- Memory L2 Compression

The optional metric lts\_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.

- L1TEX Global Load Access Pattern

Est. Speedup: 64.27%

The memory access pattern for global loads from L1TEX might not be optimal. On average, only 8.0 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 loads.

- L1TEX Local Load Access Pattern

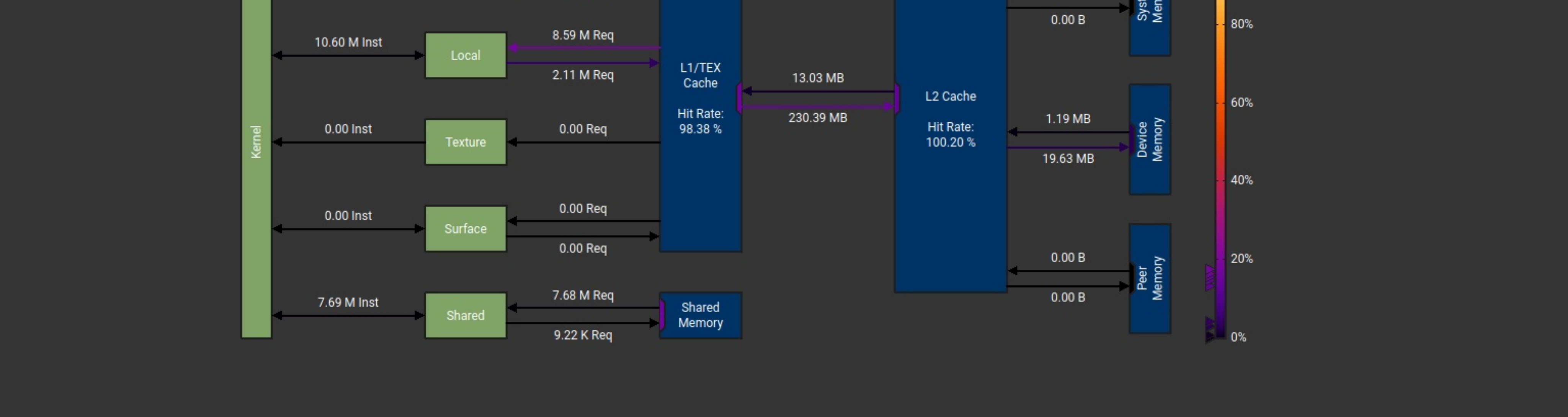
Est. Speedup: 84.88%

The memory access pattern for local loads from L1TEX might not be optimal. On average, only 0.3 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 local loads.

- L1TEX Local Store Access Pattern

Est. Speedup: 84.75%

The memory access pattern for local stores to L1TEX might not be optimal. On average, only 0.3 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 local stores.



	Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	7,680,000	7,680,000	17.66	15.79	0
Shared Load Matrix	0	0	7,680,000	0	0
Shared Store	9,216	9,216	9,216	0.02	0
Shared Atomic	0	0	0	0	0
Other	-	-	340,992	0.70	0
Total	7,689,216	7,689,216	8,030,208	16.51	0

	Instructions	Requests	Wavefronts	% Peak	Sectors/Req	Hit Rate	Bytes	Sector Misses to L2	% Peak to L2	Returns to SM
Local Load	8,587,264	8,587,264	8,587,264	17.66	32,951,360	3.84	98.82	1,054,443,520	32	0
Global Load	107,520	107,520	107,520	0.22	128,000	1.19	93.94	4,096,000	390,774	0.80
Surface Load	0	0	0	0	0	0	0	0	0	0
Texture Load	0	0	0	0	0	0	0	0	0	0
Global Store	32,768	32,768	32,768	0.07	131,072	4	0	4,194,304	7,199,556	14.80
Local Store	2,015,232	2,113,536	2,113,536	4.35	8,060,928	3.81	98.19	257,949,696	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
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	0	0	0	0	0	0	0	0	0	0
Surface Atomic ALU	0	0	0	0	0	0	0	0	0	0
Surface Atomic CAS	0	0	0	0	0	0	0	0	0	0
Loads	8,694,784	8,694,784	8,694,784	17.88	33,079,360	3.80	98.80	1,058,539,520	390,774	0.80
Stores	2,048,000	2,146,304	2,146,304	4.41	8,192,000	3.82	96.62	262,144,000	7,199,556	14.80
Atoms & Reductions	0	0	0	0	0	0	0	0	0	0
Total	10,635,784	10,841,088	10,841,088	19.16	41,270,360	7.66	98.81	1,320,683,520	407,928	0.80

	Requests	Sectors	Sectors/Req	% Peak	Hit Rate	Bytes	Throughput	Sector Misses to Device	Sector Misses to System	Sector Misses to Peer
L1/TEX Load	102,080	391,009	3.84	0.69	99.99	12,541,088	6,055,075,187.04	32	0	0
L1/TEX Store	1,896,583	7,192,962	3.79	12.70	99.67	230,174,784	111,187,813,021.70	23,572	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	1,997,188	7,594,479	3.80	13.41	99.69	243,023,328	117,394,407,345.58	23,868	0	0
ECC Total	-	23,844	-	0.01	-	763,008	368,577,258.39	23,844	-	-
GPU Total	1,988,589	7,549,973	3.80	13.33	99.69	241,599,136	116,706,439,745.25	23,971	0	0

	Sectors	% Peak	Bytes	Throughput
Load	37,183	0.18	1,189,856	574,769,677.86
Store	613,384	2.97	19,628,288	9,481,605,144.38
Total	650,567	3.15	20,818,144	10,056,374,822.23

Scheduler Statistics

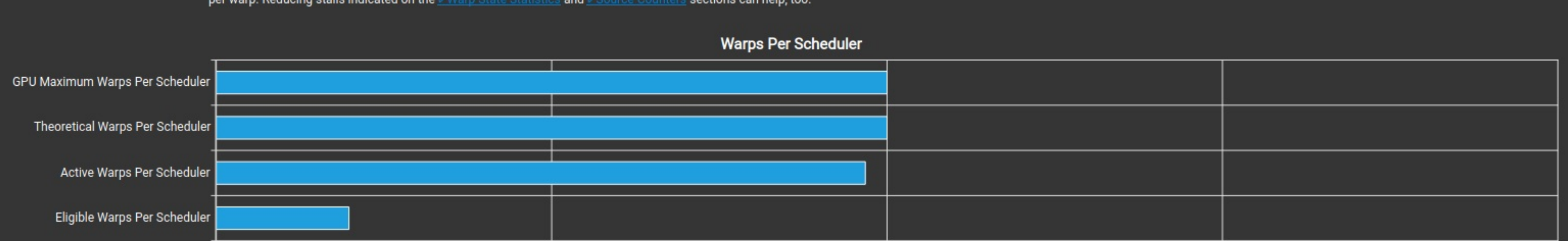
Summary of the activity of the schedulers issuing instructions. Each scheduler maintains a pool of warps that it can issue instructions for. The upper bound of warps in the pool (Theoretical Warps) is limited by the launch configuration. On every cycle each scheduler checks the state of the allocated 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 slots is skipped and no instruction is issued. Having many skipped issue slots indicates poor latency hiding.

Active Warps Per Scheduler [warp]	7.74	No Eligible [%]	41.87
Eligible Warps Per Scheduler [warp]	1.58	One or More Eligible [%]	58.13
Issued Warp Per Scheduler	0.58		

- Issue Slot Utilization

Est. Local Speedup: 22.80%

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 1.7 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 7.74 active warps per scheduler, but only an average of 1.58 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

Any scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 1.7 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 7.74 active warps per scheduler, but only an average of 1.58 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 Cycles Per Issued Instruction [cycle]	13.32	Avg. Active Threads Per Warp	26.06
Warp Cycles Per Executed Instruction [cycle]	13.32	Avg. Not Predicted Off Threads Per Warp	24.61

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 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]	105,146,208	Avg. Executed Instructions Per Scheduler [inst]	657,163.80
Issued Instructions [inst]	105,155,175	Avg. Issued Instructions Per Scheduler [inst]	657,219.84

- FP32 Non-Fused Instructions

Est. Speedup: 6.30%

This kernel executes 6604800 fused and 6604800 non-fused FP32 instructions. By converting pairs of non-fused instructions to their [fused](#), higher-throughput equivalent, the achieved FP32 performance could be increased by up to 25% (relative to its current performance). Check the Source page to identify where this kernel executes FP32 instructions.

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

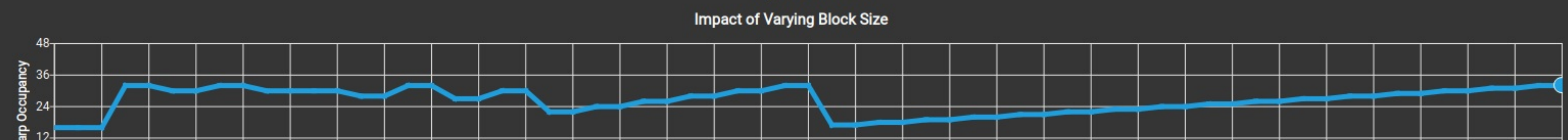
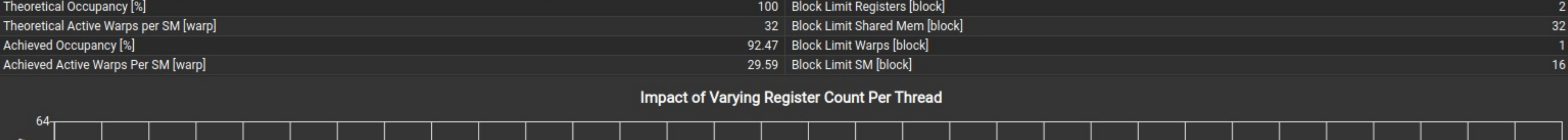
Summary of the configuration used to launch the kernel. The launch configuration defines the size of the kernel grid, the division of the grid into blocks, and the GPU resources needed to execute the kernel. Choosing an efficient launch configuration maximizes device utilization.

Grid Size	1,024	Function Cache Configuration	CachePreferNone
Registers Per Thread [register/thread]	21	Static Shared Memory Per Block [byte/block]	0
Block Size	1,024	Dynamic Shared Memory Per Block [byte/block]	900
Threads [thread]	1,048,576	Diver Shared Memory Per Block [byte/block]	0
Waves Per SM	23.60	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]	2
Theoretical Active Warps per SM [warp]	32	Block Limit Shared Mem [block]	32
Achieved Occupancy [%]	92.47	Block Limit Warps [block]	1
Achieved Active Warps Per SM [warp]	29.59	Block Limit SM [block]	16



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]	1,181,089.43	Average L1 Active Cycles [cycle]	1,181,089.43
Average L2 Active Cycles [cycle]	1,100,650.38	Average SMSP Active Cycles [cycle]	1,130,614.20
Average DRAM Active Cycles [cycle]	325,283.50	Total SM Elapsed Cycles [cycle]	48,629,392
Total L1 Elapsed Cycles [cycle]	48,629,392	Total L2 Elapsed Cycles [cycle]	56,637,152
Total SMSP Elapsed Cycles [cycle]	194,517,568	Total DRAM Elapsed Cycles [cycle]	82,569,216

- L2 Slices Workload Imbalance

Est. Speedup: 5.19%

One or more L2 Slices have a much lower number of active cycles than the average number of active cycles. Maximum instance value is 8.34% above the average, while the minimum instance value is 12.08% below the average.

	Average	Min	Max	Sum
SM Active Cycles	1,181,089.43	1,156,600	1,208,780	47,243,577
SMSP Active Cycles	1,130,614.20	1,077,068	1,178,307	160,998,272
L1 Active Cycles	1,181,089.43	1,156,600	1,208,780	47,243,577
L2 Active Cycles	1,100,650.38	967,641	1,200,793	35,220,812
DRAM Active Cycles	325,283.50	301,544	342,840	2,602,268

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]	8,943,456	Branch Efficiency [%]	97.73
Branch Instructions Rate [%]	0.09	Avg. Divergent Branches	1,131.80