

GPU Speed of Light Throughput

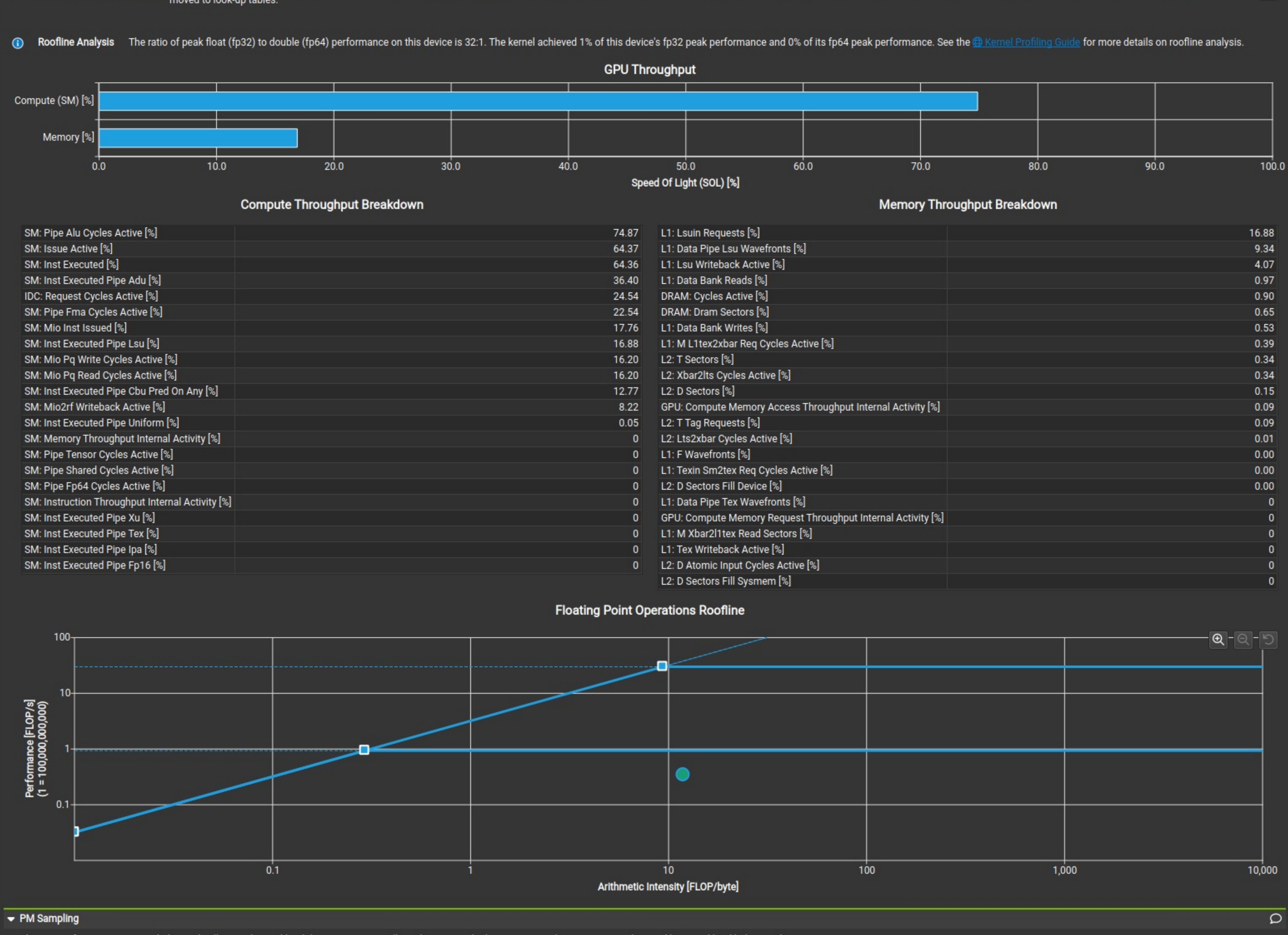
All

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 [%]	74.87	Duration [ms]	1.43
Memory Throughput [%]	16.88	Elapsed Cycles [cycle]	834,377
L1/TEX Cache Throughput [%]	18.68	SM Active Cycles [cycle]	810,912.55
L2 Cache Throughput [%]	0.34	SM Frequency [Mhz]	584.97
DRAM Throughput [%]	0.90	DRAM Frequency [Ghz]	5.00

**High Compute Throughput** Compute is more heavily utilized than Memory. Look at the [Compute Workload Analysis](#) section to see what the compute pipelines are spending their time doing. Also, consider whether any computation is redundant and could be reduced or moved to look-up tables.

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

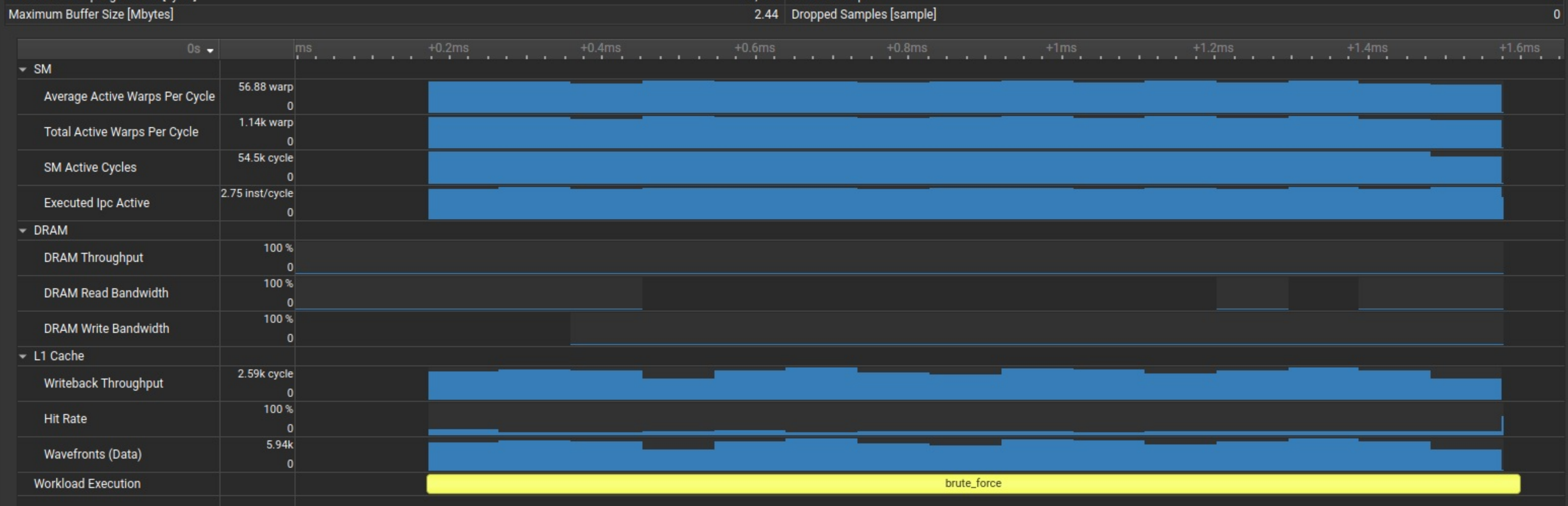


PM Sampling

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]80,000 # Pass Groups1

Maximum Buffer Size [Bytes]2.44 Dropped Samples [sample]0



Compute Workload Analysis

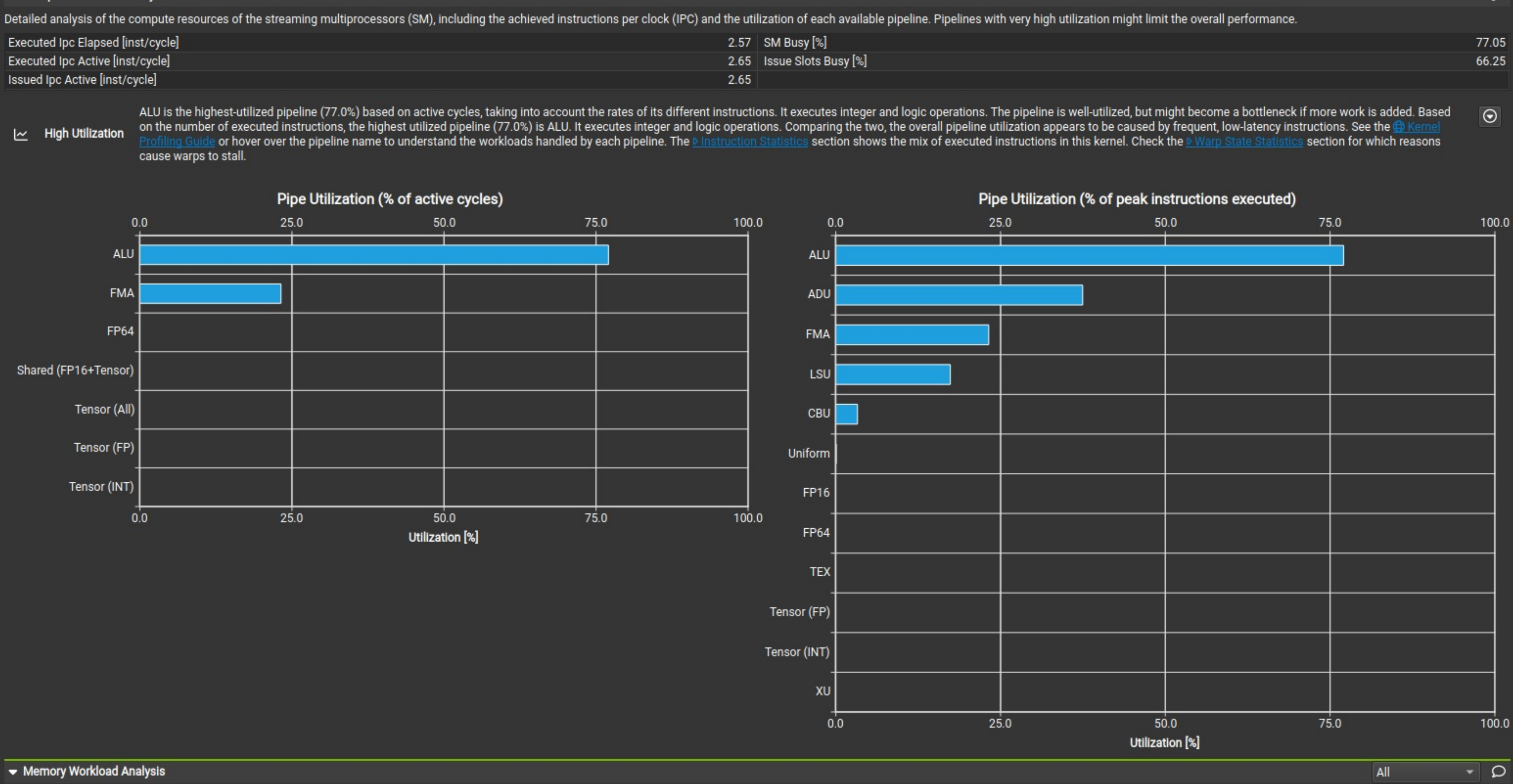
Detailed analysis of the compute resources of the streaming multiprocessors (SM), including the achieved instructions per clock (IPC) and the utilization of each available pipeline. Pipelines with very high utilization might limit the overall performance.

Executed Ipc Elapsed [inst/cycle]2.57 SM Busy [%]77.05

Executed Ipc Active [inst/cycle]2.65 Issue Slots Busy [%]66.25

Issued Ipc Active [inst/cycle]2.65

**High Utilization** ALU is the highest-utilized pipeline (77.0%) based on active cycles, taking into account the rates of its different instructions. It executes integer and logic operations. The pipeline is well-utilized, but might become a bottleneck if more work is added. Based on the number of executed instructions, the highest utilized pipeline (77.0%) is ALU. It executes integer and logic operations. Comparing the two, the overall pipeline utilization appears to be caused by frequent, low-latency instructions. See the [Kernel Profiling Guide](#) for more details on roofline analysis. If you have over the pipeline name to understand the workloads handled by each pipeline. The [Instruction Summary](#) section shows the mix of executed instructions in this kernel. Check the [High-Level Summary](#) section for which reasons cause warps to stall.



Memory Workload Analysis

All

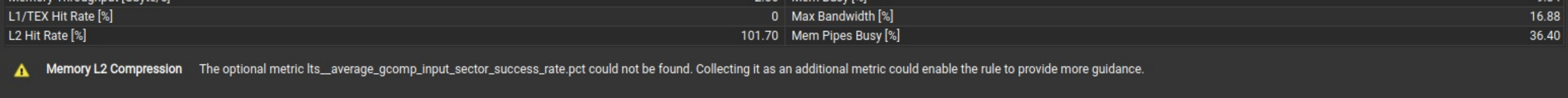
Detailed analysis of the memory resources of the GPU. Memory can become a limiting factor for the overall kernel performance when fully utilizing the involved hardware units (Mem Busy), exhausting the available communication bandwidth between those units (Max Bandwidth), or by reaching the maximum throughput of issuing memory instructions (Mem Pipes Busy). Detailed chart of the memory units.

Memory Throughput [bytes/s]2.86 Mem Busy [%]9.34

L1/TEX Hit Rate [%]0 Max Bandwidth [%]16.88

L2 Hit Rate [%]101.70 Mem Pipes Busy [%]36.40

**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.



Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	1,327,104	1,327,104	3.98	0
Shared Matrix	0	0	0	0
Shared Store	1,425,408	1,425,408	5.16	0
Shared Atomic	0	0	0	0
Other	-	-	32,937	0.10
Total	2,752,512	2,752,512	3,083,750	9.24

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	0	0
Global Load	0	0	0	0	0	0	0	0	0	0
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	0.10	131,072	4	0	4,194,304	131,072	0.39	-
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	0	0	0	0	0	0	0	0	0	see above
Surface Atomic ALU	0	0	0	0	0	0	0	0	0	see above
Surface Atomic CAS	0	0	0	0	0	0	0	0	0	see above
Loads	0	0	0	0	0	0	0	0	0	0
Stores	32,768	32,768	0.10	131,072	4	0	4,194,304	131,072	0.39	-
Atoms & Reductions	0	0	0	0	0	0	0	0	0	-
GPU Total	33,701	132,693	3.94	34	99.74	4,344,758	2,975,033,474.39	58	0	0

Requests	Sectors	Sectors/Req	% Peak	Hit Rate	Bytes	Throughput	Sector Misses to Device	Sector Misses to System	Sector Misses to Peer
L1/TEX Load	0	0	0	0	0	0	0	0	0
L1/TEX Store	32,768	131,072	4	0.34	100	4,194,304	2,940,614,273.21	0	0
L1/TEX Atomic ALU	0	0	0	0	0	0	0	0	0
L1/TEX Atomic CAS	0	0	0	0	0	0	0	0	0
L1/TEX Reduction	0	0	0	0	0	0	0	0	0
L1/TEX Total	32,768	131,072	4	0.34	100	4,194,304	2,940,614,273.21	0	0
EOC Total	24	-	0.00	-	768	338,442.55	24	-	-
GPU Total	33,701	132,693	3.94	0.34	99.74	4,344,758	2,975,033,474.39	58	0

Sectors	% Peak	Bytes	Throughput
Load	2,050	0.01	65,600
Store	125,564	0.88	4,018,048
Total	127,614	0.90	4,083,648

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 slot is skipped and no instruction is issued. Having many skipped issue slots indicates poor latency hiding.

Active Warps Per Scheduler [warp]7.60 No Eligible [%]26.60

Eligible Warps Per Scheduler [warp]3.13 One or More Eligible [%]73.40

Issued Warp Per Scheduler0.73



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.

Warp Cycles Per Issued Instruction [cycle]10.36 Avg. Active Threads Per Warp22.22

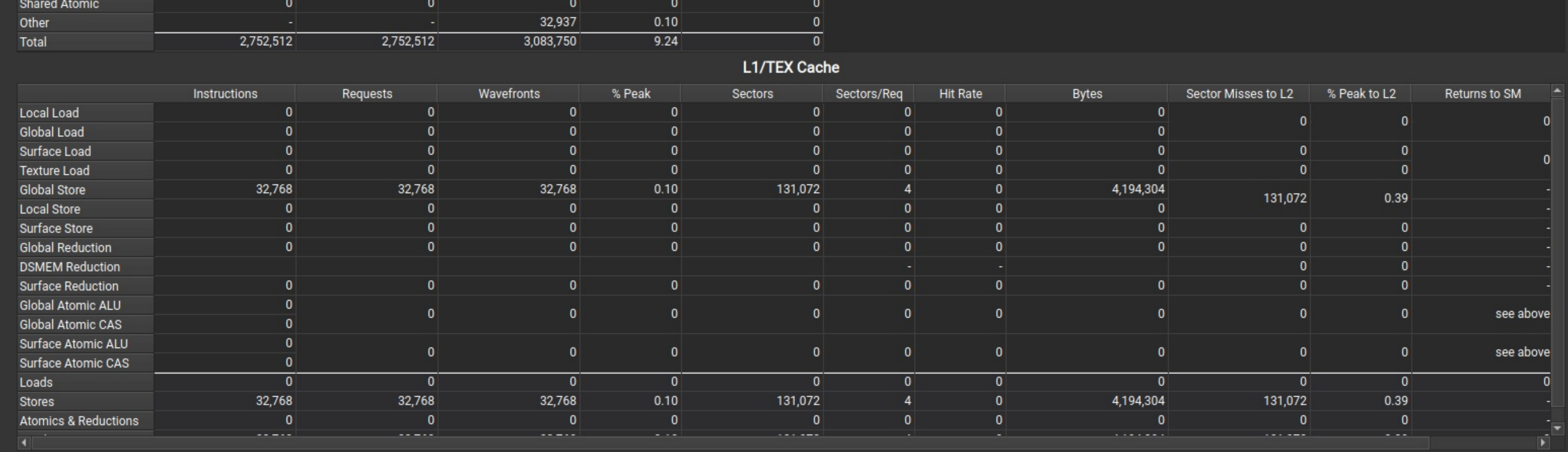
Warp Cycles Per Executed Instruction [cycle]10.36 Avg. Not Predicated Off Threads Per Warp20.24

**Not Selected Stalls** On average, each warp of this kernel spends 3.3 cycles being stalled waiting for the micro scheduler to select the warp to issue. Not selected warps are eligible warps that were not picked by the scheduler to issue that cycle as another warp was selected. A high number of not selected warps typically means you have sufficient warps to cover warp latencies and you may consider reducing the number of active warps to possibly increase cache coherence and data locality. This stall type represents about 31.5% of the total average of 10.4 cycles between issuing two instructions.

**Ext. Local Speedup: 31.55%**

**Launch Statistics** Check the [Warp Stall Sampling \(All Samples\)](#) table for the top stall locations in your source based on sampling data. The [Kernel Profiling Guide](#) provides more details on each stall reason.

**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 lower the number of active threads in a warp per cycle. This kernel achieves an average of 22.2 threads being active per cycle. This is further reduced to 20.2 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.



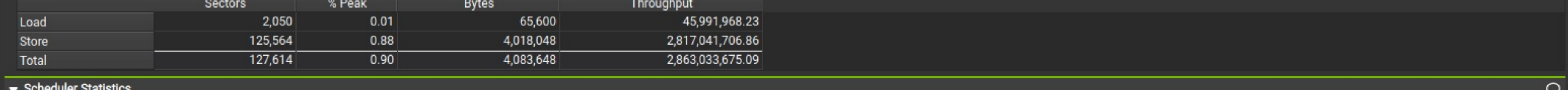
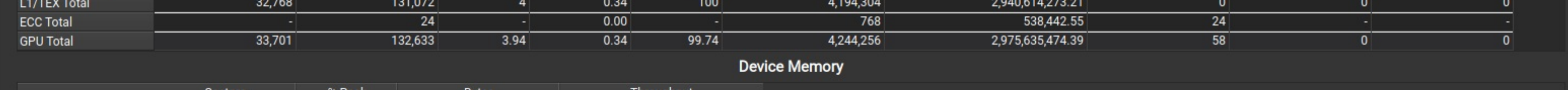
Instruction Statistics

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 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]85,936,672 Avg. Executed Instructions Per Scheduler [inst]537,104.20

Issued Instructions [inst]85,950,832 Avg. Issued Instructions Per Scheduler [inst]537,192.70

**FP32 Non-Fused Instructions** This kernel executes 0 fused and 464,588 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 50% (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 Configuration

Summary of the configurations 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 Size1,024 Function Cache ConfigurationCachePreferNone

Registers Per Thread [register/thread]21 Static Shared Memory Per Block [byte/block]0

Block Size1,024 Dynamic Shared Memory Per Block [byte/block]3,037

Threads [thread]1,048,576 Shared Shared Memory Per Block [byte/block]0

Waves Per SM25.60 Driver Memory Configuration Size [kbyte]32.77

Uses Green Context0 # 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

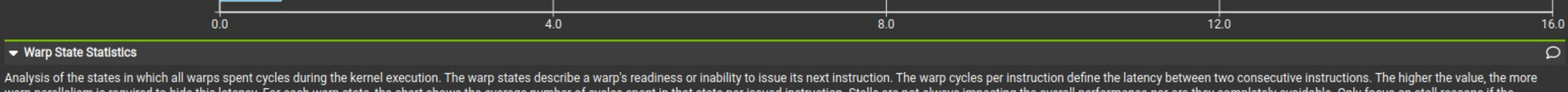
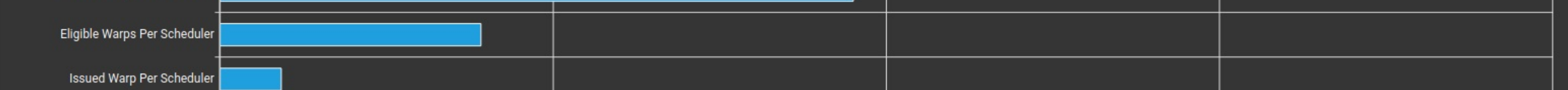
Theoretical Active Warps per SM [warp]32 Block Limit Shared Mem [block]10

Achieved Occupancy [%]85.84 Block Limit Warps [block]1

Achieved Active Warps per SM [warp]27.47 Block Limit SM [block]16

**Achieved Occupancy: 85.84%** The difference between calculated theoretical (100.0%) and measured achieved occupancy (85.84%) can be the result of warp scheduling overheads or workload imbalances during the kernel execution. Look the [Occupancy Guide](#) for more details on optimizing occupancy.

**Ext. Local Speedup: 14.16%**



GPU and Memory Workload Distribution

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

Average P0 Active Cycles [cycle]810,912.55 Average L1 Active Cycles [cycle]810,912.55

Average L2 Active Cycles [cycle]44,885.06 Average SMSP Active Cycles [cycle]731,912.74

Average DRAM Active Cycles [cycle]63,807 Total SM Elapsed Cycles [cycle]33,380,240

Total L1 Elapsed Cycles [cycle]33,380,240 Total L2 Elapsed Cycles [cycle]39,022,432

Total SMSP Elapsed Cycles [cycle]133,520,960 Total DRAM Elapsed Cycles [cycle]57,003,008

Workload Distribution

	Average	Min	Max	Sum
SM Active Cycles	810,912.55	793,239	829,584	32,436,502
SMSP Active Cycles	731,912.74	690,806	768,979	117,106,039
L1 Active Cycles	810,912.55	793,239	829,584	32,436,502
L2 Active Cycles	44,885.06	42,271	45,488	1,436,322
DRAM Active Cycles	63,807	62,164	64,728	510,456

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 [%]9,380,332 Branch Efficiency [%]0.11 Avg. Divergent Branches2,004.40

Branch Instructions Ratio [%]0.11

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 individual sections to focus on selected performance aspects and make profiling faster.