

GPU Speed of Light Throughput

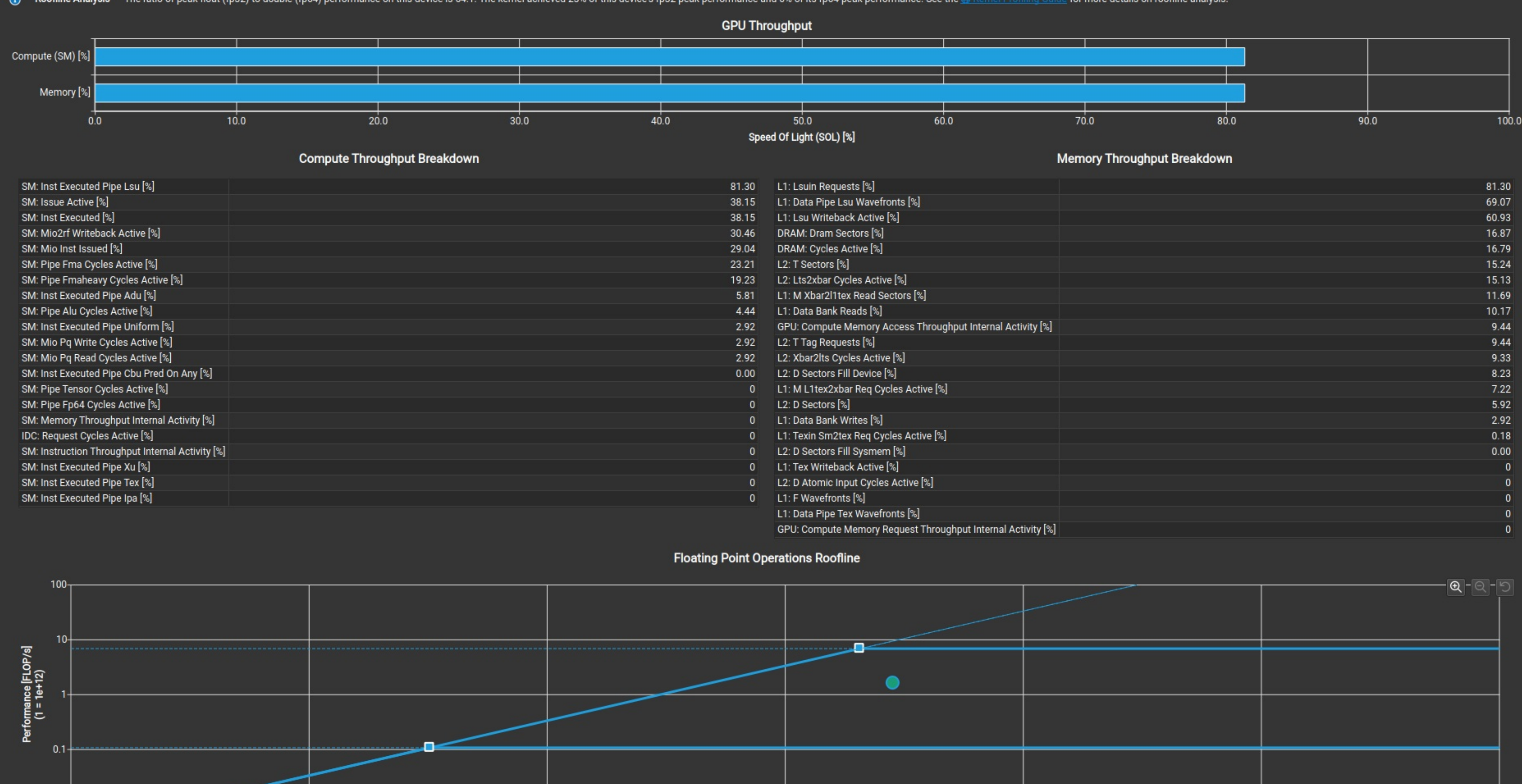
All

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 [%]	81.30	Duration [ms]	85.72
Memory Throughput [%]	81.30	Elapsed Cycles [cycle]	77151538
L1/TEX Cache Throughput [%]	81.52	SM Active Cycles [cycle]	76922855.97
L2 Cache Throughput [%]	15.24	SM Frequency [MHz]	900.00
DRAM Throughput [%]	16.67	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, it will likely need to be shifted from the most utilized to another unit. Start by analyzing workloads in the [Compute Workload Analysis](#) section.

Roofline Analysis The ratio of peak float (fp32) to double (fp64) performance on this device is 64:1. The kernel achieved 23% 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.

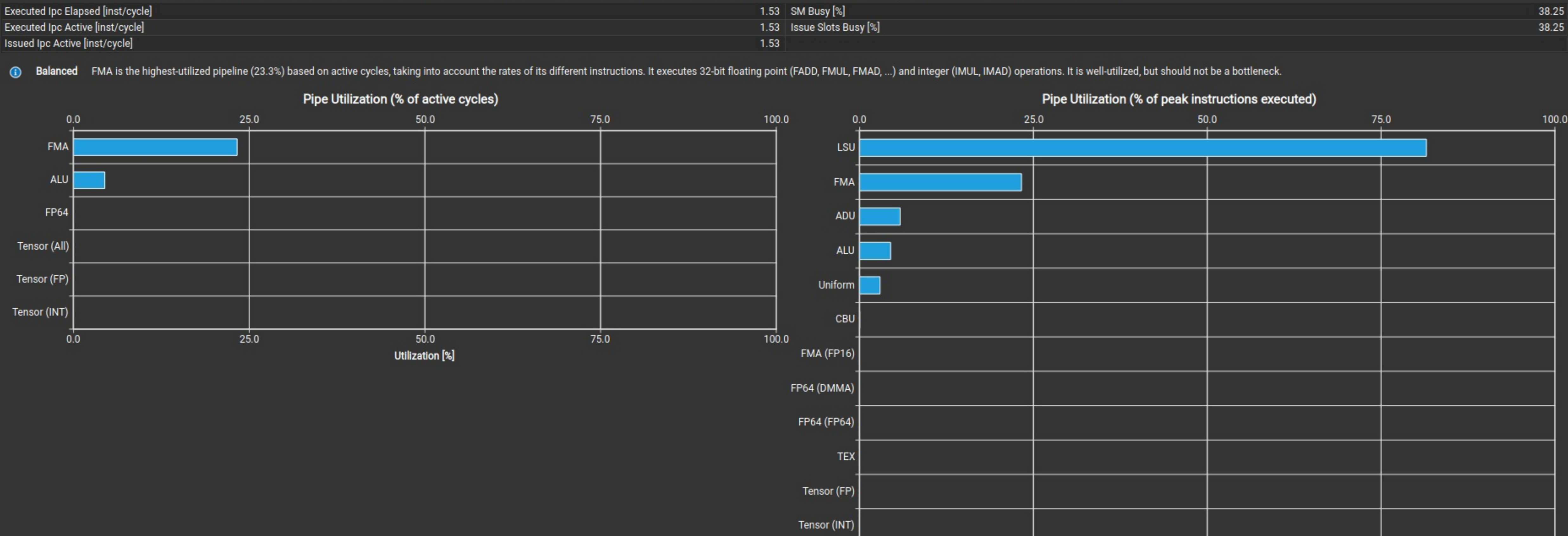
Maximum Sampling Interval [ms]	1.02	# Pass Groups	2
Maximum Buffer Size [Mbytes]	19.19	Dropped Samples [sample]	0

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]	1.53	SM Busy [%]	38.25
Executed ipc Active [inst/cycle]	1.53	Issue Slots Busy [%]	38.25
Issued ipc Active [inst/cycle]	1.53		

Balanced FMA is the highest-utilized pipeline (23.3%) based on active cycles, taking into account the rates of its different instructions. It executes 32-bit floating point (FADD, FMUL, FMAD, ...) and integer (MUL, IMAD) operations. It is well-utilized, but should not be a bottleneck.

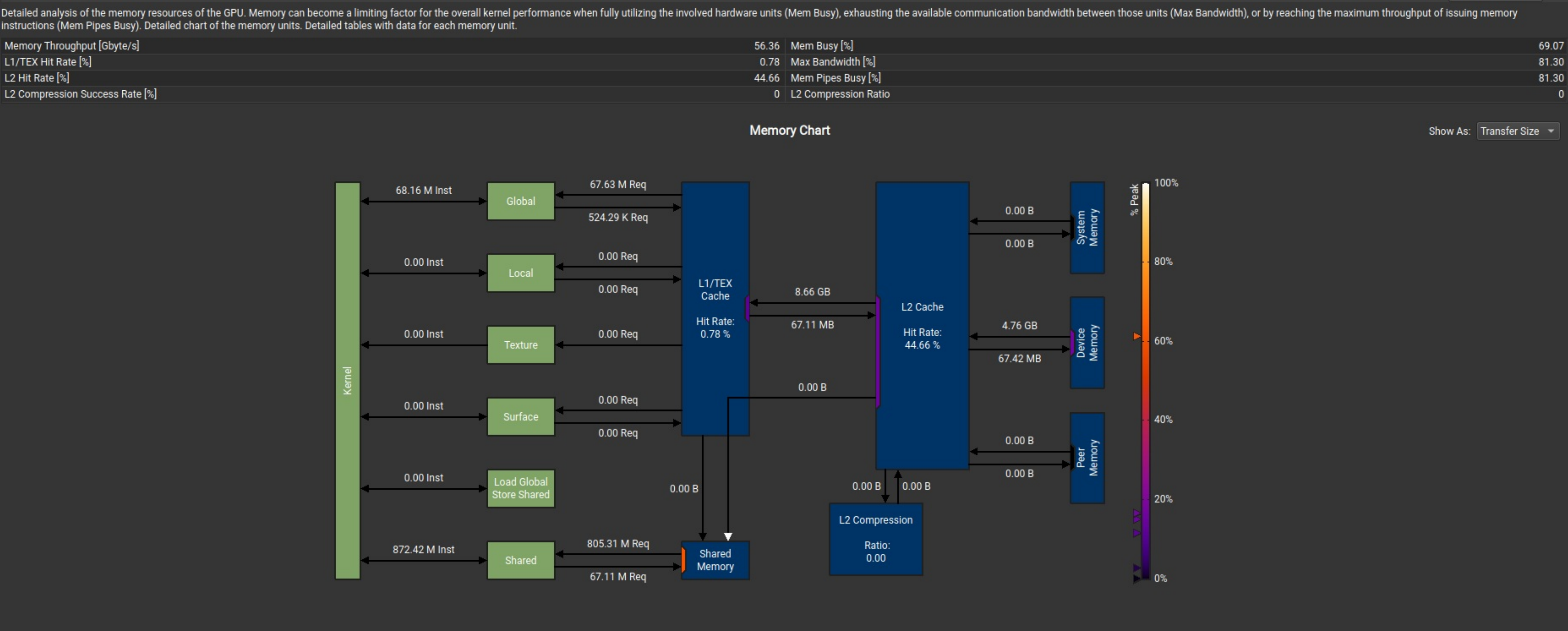


Memory Workload Analysis

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 [GB/s]	56.36	Mem Busy [%]	69.07
L1/TEX Hit Rate [%]	0.78	Max Bandwidth [%]	81.30
L2 Hit Rate [%]	44.66	Mem Pipes Busy [%]	81.30
L2 Compression Success Rate [%]	0	L2 Compression Ratio	0



	Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	805306368	805306368	1342301236	88.00	79789
Global Load To Shared Store (access)	0	0	0	0	0
Global Load To Shared Store (bypass)	67108864	67108864	67108864	0.72	0
Shared Store From Global Load	0	0	0	0	0
Shared Atomic	0	0	0	0	0
Other	-	11630599	2.68	7296262	
Total	872415232	872415232	1421040599	61.41	7376041

	Instructions	Requests	Wavefronts	% Peak	Sectors	Sectors/Req	Hit Rate	Bytes	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	67721428	2.93
Global Load	67633152	67633152	0	0	270462061	4.00	0.01	8654789592	270506016	11.69	67721428	-
Global Load To Shared Store (access)	0	0	67634962	2.92	0	0	0	0	0	0	-	-
Global Load To Shared Store (bypass)	0	0	0	0	0	0	0	0	0	0	-	-
Surface Load	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
Global Store	524288	524288	524288	0.02	2097152	4	100.00	67108864	2097152	0.09	-	-
Local Store	0	0	0	0	0	0	0	0	0	0	-	-
Surface Store	0	0	0	0	0	0	0	0	0	0	-	-
Global Reducers	0	0	0	0	0	0	0	0	0	0	-	-
DSMEM Reduction	0	0	0	0	0	0	0	0	0	0	-	-
Surface Reduction	0	0	0	0	0	0	0	0	0	0	-	-
Global Atomic CAS	0	0	0	0	0	0	0	0	0	0	see above	see above
Surface 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	see above	see above
Surface Atomic CAS	0	0	0	0	0	0	0	0	0	0	-	-
Loads	67633152	67633152	67634962	2.92	270462061	4.00	0.01	8654789592	270506016	11.69	67721428	2.93
Stores	524288	524288	524288	0.02	2097152	4	100.00	67108864	2097152	0.09	-	-
Atomics & Reductions	0	0	0	0	0	0	0	0	0	0	-	-
Total	68157440	68157440	68159250	2.95	272559213	4.00	0.78	8721894816	272603168	11.78	67721428	2.93

	Requests	Sectors	Sectors/Req	% Peak	Hit Rate	Bytes	Throughput	Sector Misses to Device	Sector Misses to System	Sector Misses to Peer
L1/TEX Load	16839815	270506168	1.61	15.11	44.64	8656197376	100977623267.14	150177716	0	0
L1/TEX Store	524288	2097152	4	0.12	100	67108864	762848023.99	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	168814170	272603064	1.61	15.23	45.06	8723298048	101766277728.52	151057640	0	0
ECO Total	-	0	-	0	-	0	10185298503.57	151038592	15072	0
GPU Total	169031714	272851371	1.61	15.24	44.65	8731243872	10185298503.57	151038592	15072	0

	First	Hit Rate	Last	Hit Rate	Normal	Hit Rate	Normal Demote	Hit Rate
L1/TEX Load	0	0	0	0	270458500	44.78	0	0
L1/TEX Store	0	0	0	0	2097152	100	0	0
L1/TEX Atomic	0	0	0	0	0	0	-	-
L1/TEX Total	0	0	0	0	272781844	45.30	0	0
GPU Total	251956	81.51	0	0	273772273	44.84	0	0

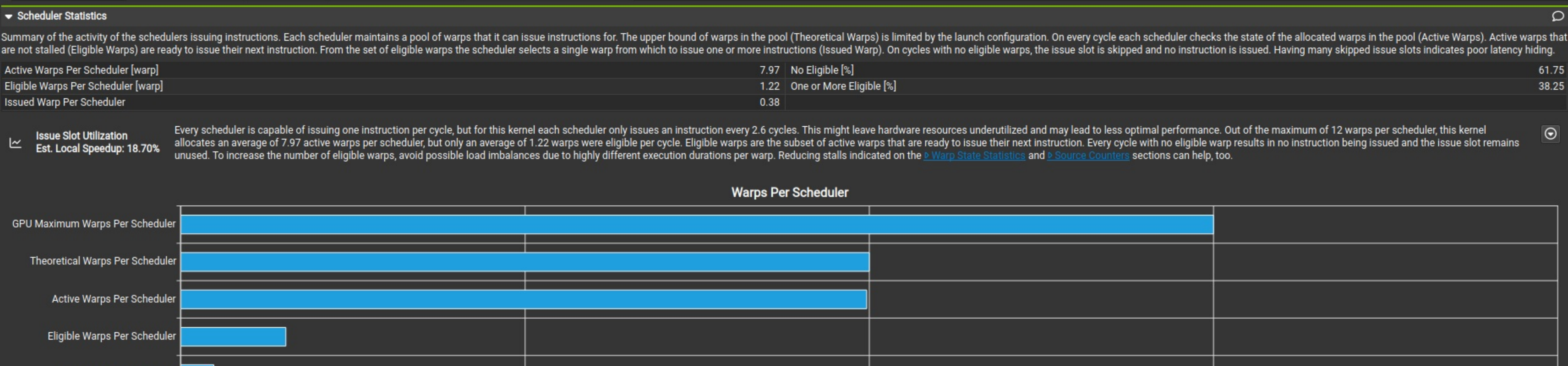
	Sectors	% Peak	Bytes	Throughput
Load	148864684	16.55	4763669888	55569850776.91
Store	2109228	0.23	67421696	786497316.97
Total	150971012	16.73	4831091684	5635536093.88

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.

Active Warps Per Scheduler [warp]	7.97	No Eligible [%]	61.75
Eligible Warps Per Scheduler [warp]	1.22	One or More Eligible [%]	38.25
Issued Warp Per Scheduler	0.38		

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 2.6 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 7.97 active warps per scheduler, but only an average of 1.22 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 [Warp State Statistics](#) sections can help, too.



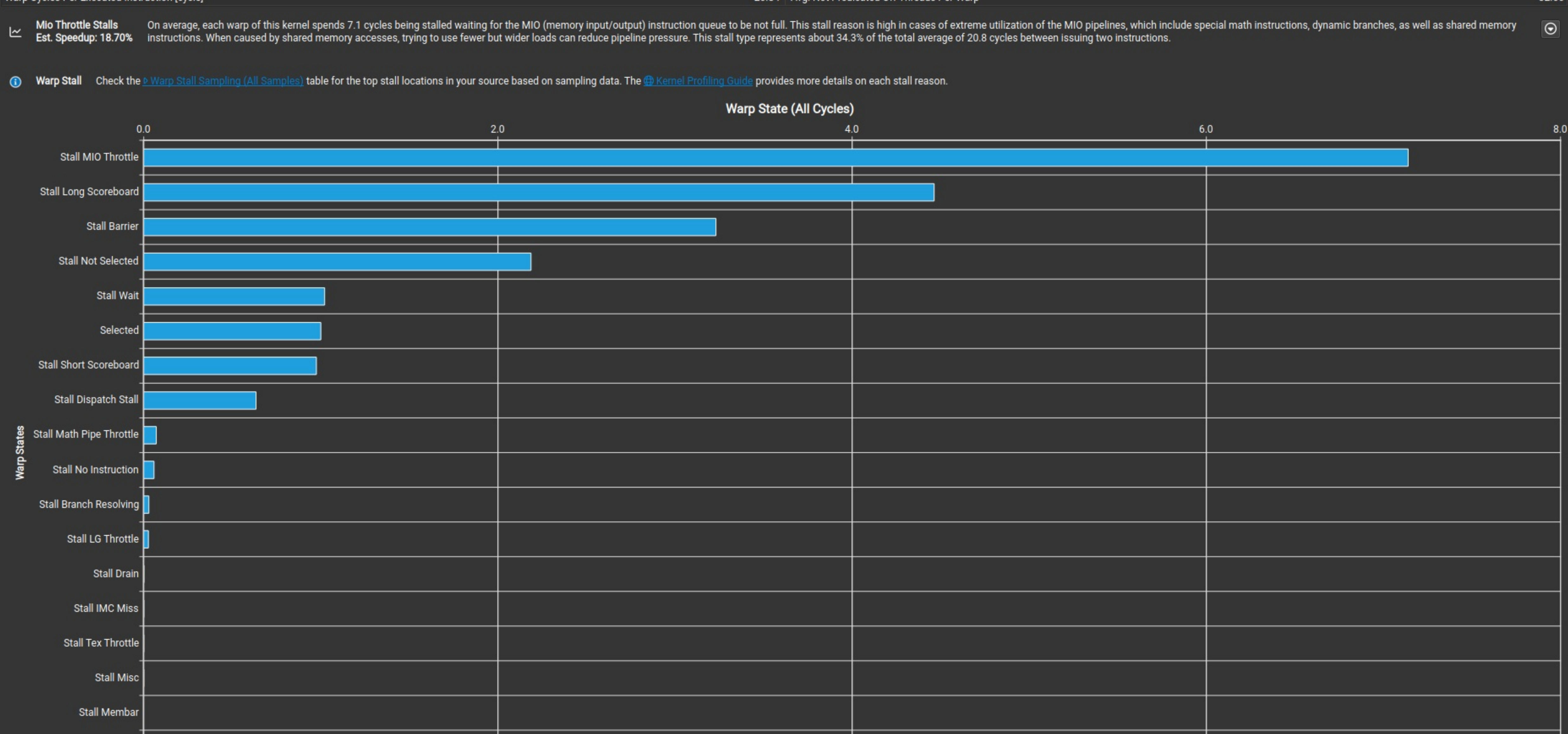
Warp State Statistics

Analysis of the states in which all warps spent cycles during the kernel execution. The warp state describes 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]	20.84	Avg. Active Threads Per Warp	32
Warp Cycles Per Executed Instruction [cycle]	20.84	Avg. Not Predicted Off Threads Per Warp	32.00

Issue Slot Utilization: 18.70% On average, each warp of this kernel spends 2.1 cycles being stalled waiting for the MIO (memory input/output) instruction queue to be not full. This stall reason is high in cases of extreme utilization of the MIO pipelines, which include special math instructions, dynamic branches, as well as shared memory instructions. When caused by shared memory accesses, trying to use fewer but wider loads can reduce pipeline pressure. This stall type represents about 34.3% of the total average of 20.8 cycles between issuing two instructions.

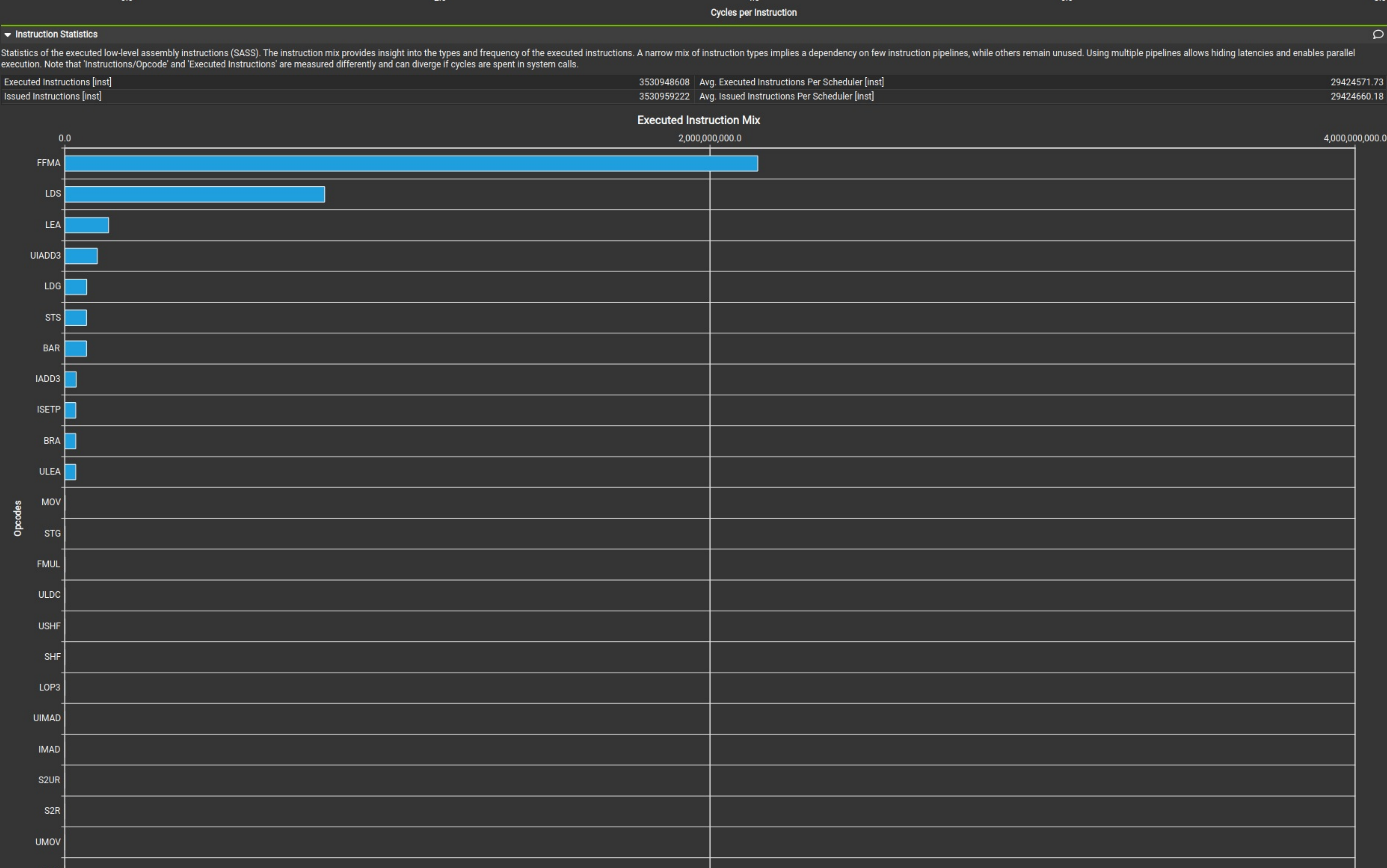
Warp Stall Check the [Warp Stall Sampling \(All Samples\)](#) table for the top stall locations in your source based on sampling data. The [Kernel Profiling Guide](#) provides more details on each stall reason.



Instructions Executed

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]	3530948608	Avg. Executed Instructions Per Scheduler [inst]	294245171.73
Issued Instructions [inst]	3530959222	Avg. Issued Instructions Per Scheduler [inst]	29424665.18



GPU and Memory Workload Distribution

Analysis of workload distribution by activity of SM, SMP, SMSP, L1 & L2 caches, and DRAM

SM Active Cycles [cycle]	76922855.97	Average L1 Active Cycles [cycle]	76922855.97
Average L2 Active Cycles [cycle]	74454824	Average SMSP Active Cycles [cycle]	76922855.97
Average DRAM Active Cycles [cycle]	72111749	Total SM Elapsed Cycles [cycle]	2314126540
Total L1 Elapsed Cycles [cycle]	74450511.17	Total L2 Elapsed Cycles [cycle]	1789915632
Total SMSP Elapsed Cycles [cycle]	100647741.33	Total DRAM Elapsed Cycles [cycle]	3597320192

Workload Distribution

	Average	Min	Max	Sum
SM Active Cycles	76922855.97	76644324	77111749	2307685679
SMSP Active Cycles	7693249.24	76619978	77153992	9231898909
L2 Active Cycles	76922855.97	76644324	77111749	2307685679
DRAM Active Cycles	74450511.17	74460041	74464225	1786812268
DRAM Active Cycles	100647741.33	100575616	100692592	603886448

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]	33685504	Branch Efficiency [%]	100
Branch Instructions Ratio [%]	0.01	Avg. Divergent Branches	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 disable [Warp Stall Sampling](#) to focus on selected performance aspects and make profiling faster.