

Current	529 - applyMultiBandGainKernel	(6374.1, 1) x (1024, 1, 1)	2.08 ms	1,214,252	0 - Tesla T4	564.81 Mhz	[4558] OutdEQ3	
Summary	Details	Source	Context	Comments	Raw	Session		

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 rooftop chart.

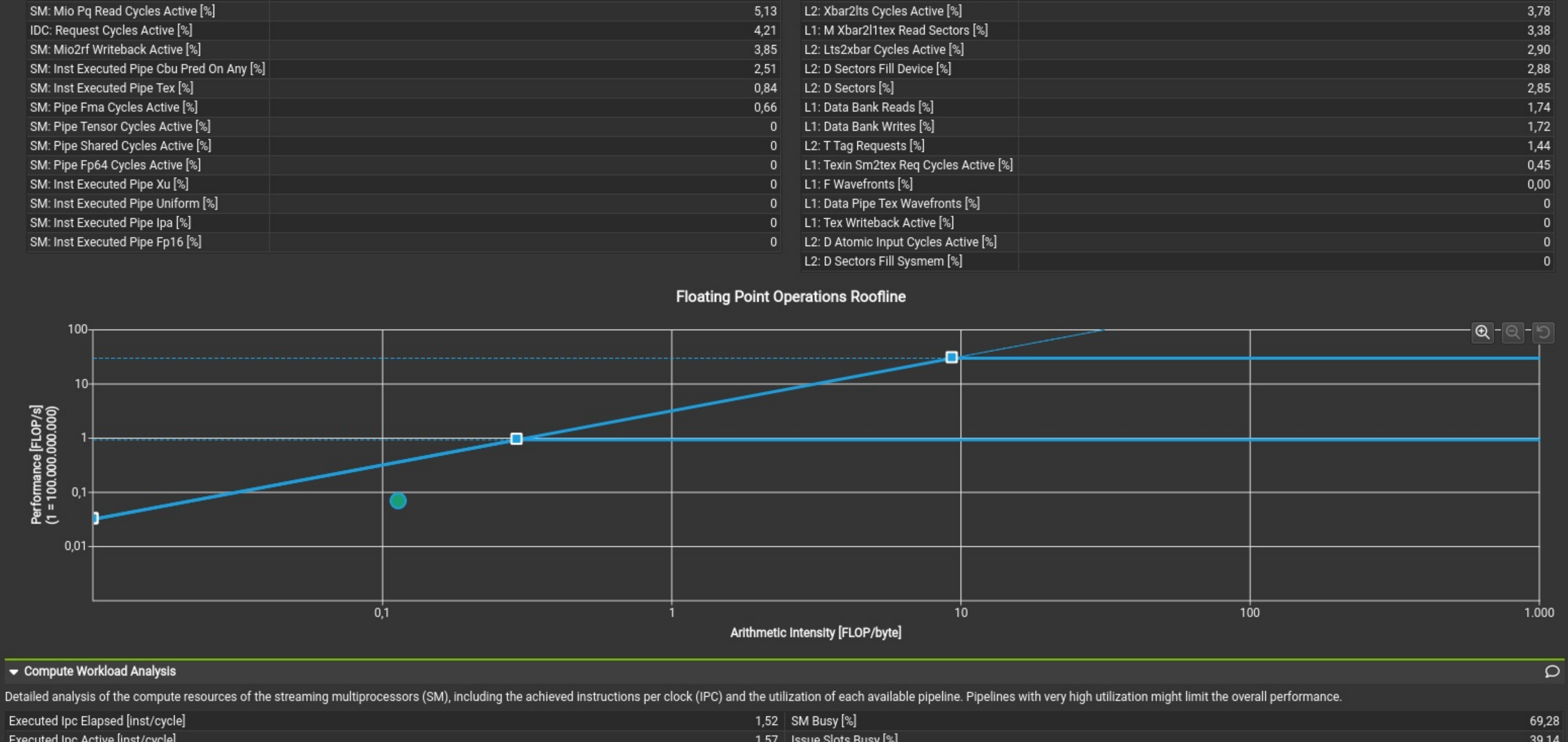
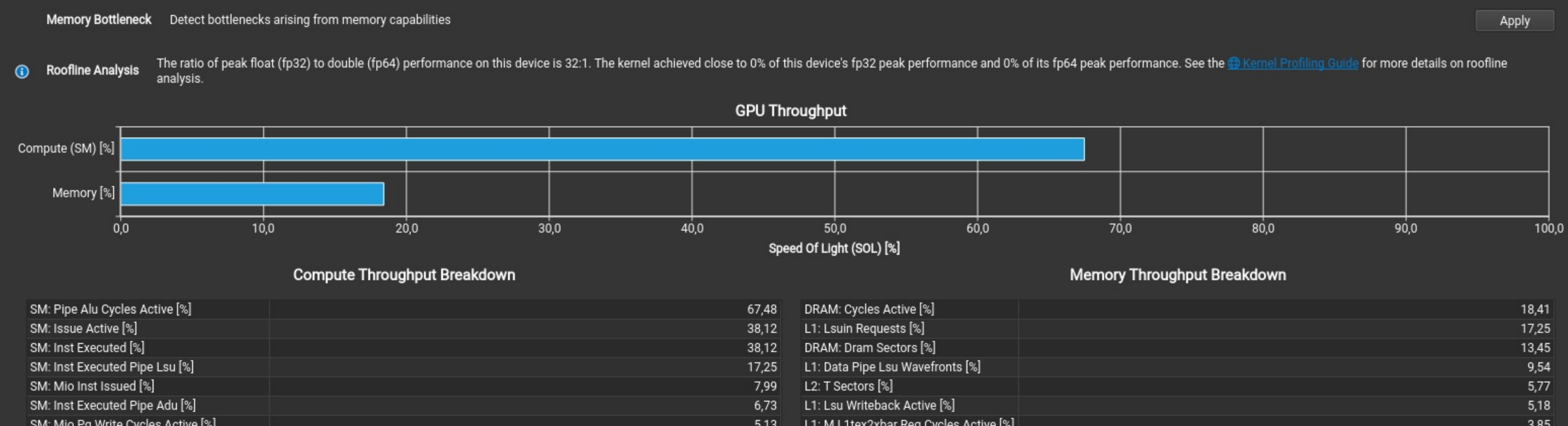
Compute (SM) Throughput [%]	67.48	Duration [ms]	2.08
Memory Throughput [%]	18.41	Elapsed Cycles [cycle]	1,214,252
L1/TEX Cache Throughput [%]	19.09	SM Active Cycles [cycle]	1,182,668.92
L2 Cache Throughput [%]	5.77	SM Frequency [MHz]	564.81
DRAM Throughput [%]	18.41	DRAM Frequency [GHz]	4.99

High Compute Throughput Compute is more heavily utilized than Memory. Look at the [F-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.

Compute Bottleneck Detect bottlenecks arising from compute capabilities

Memory Bottleneck Detect bottlenecks arising from memory capabilities

Roofline Analysis The ratio of peak float (fp32) to double (fp64) performance on this device is 32.1. The kernel achieved close to 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.



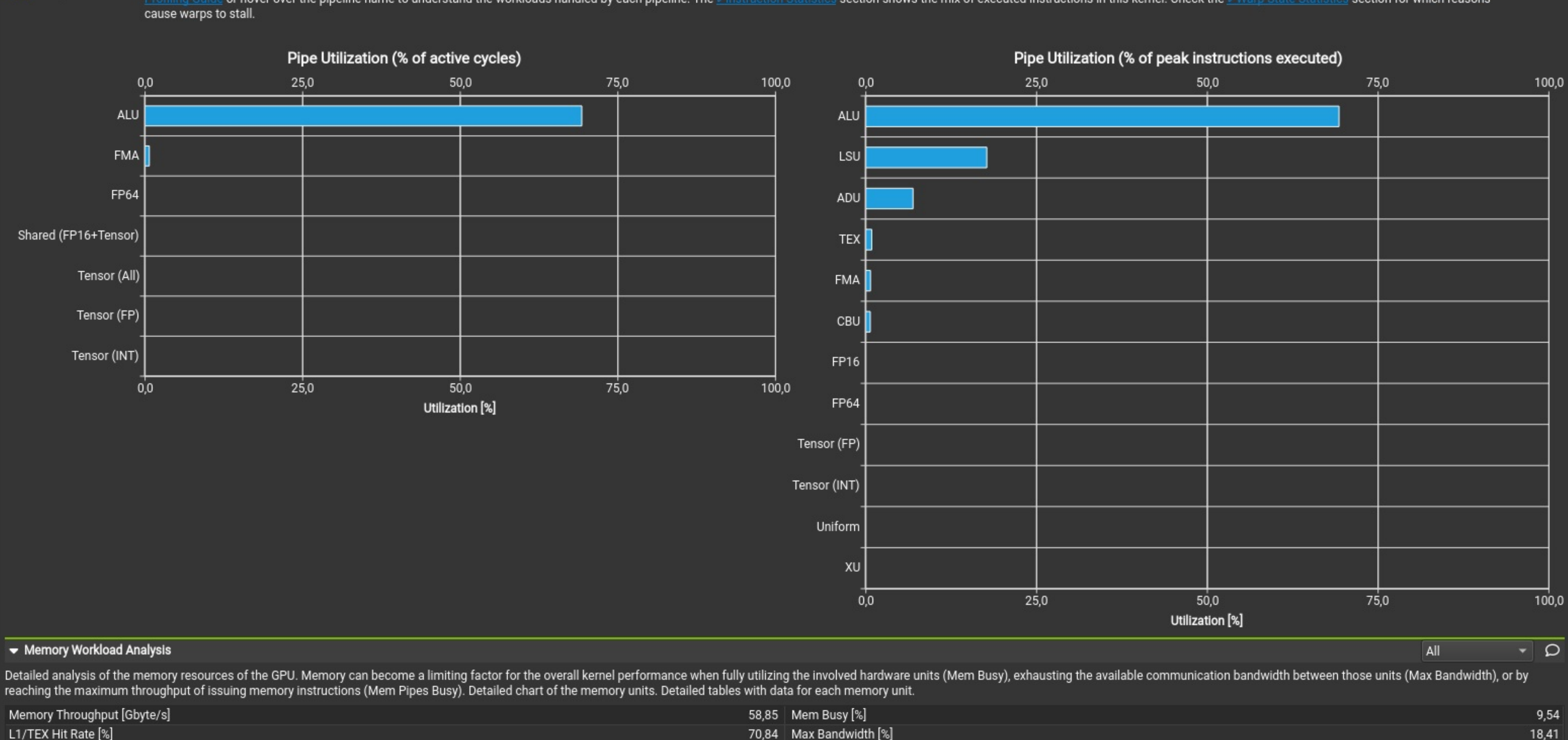
Compute Workload Analysis

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]	1.52	SM Busy [%]	69.28
Executed Ipc Active [Inst/cycle]	1.57	Issue Slots Busy [%]	39.14
Issued Ipc Active [Inst/cycle]	1.57		

High Utilization ALU is the highest-utilized pipeline (69.3%) 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 (69.3%) 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](#) section for more details on pipeline utilization. The [Instruction Sampling](#) section shows the mix of executed instructions in this kernel. Check the [Warp State Sampling](#) section for which reasons cause warps to stall.

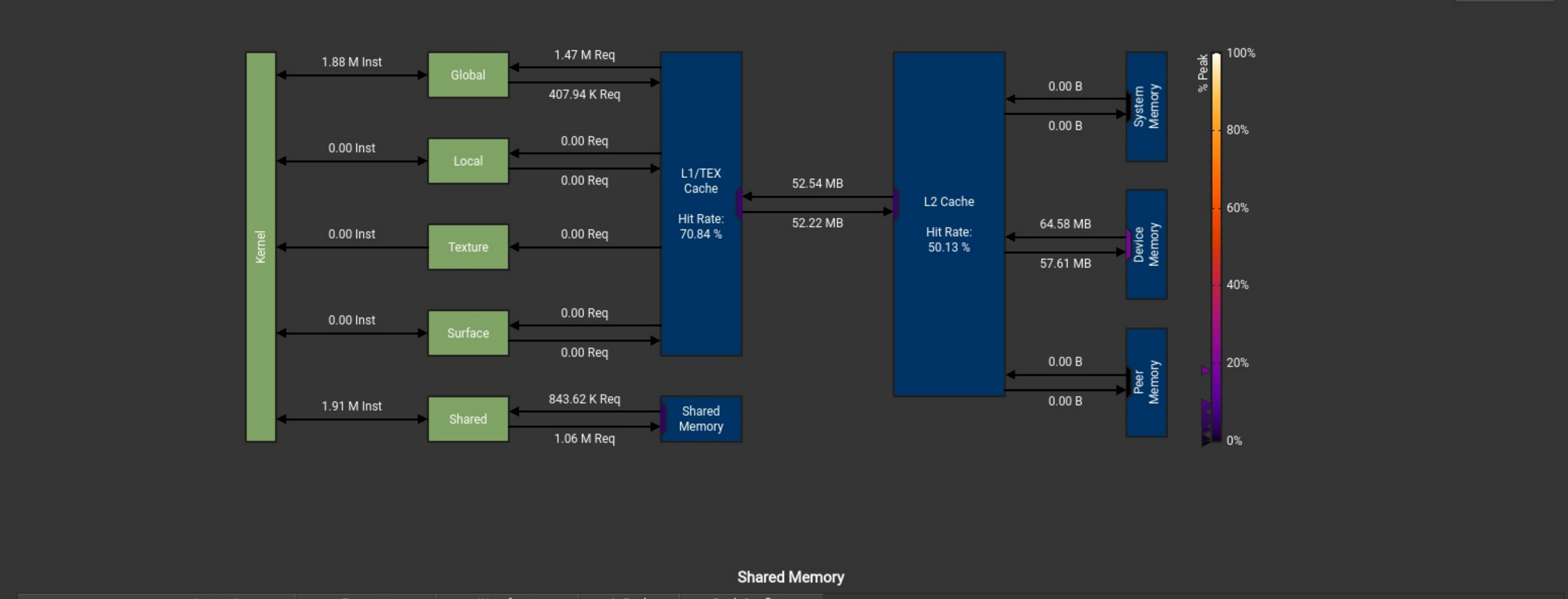


Memory Workload Analysis

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 Bus), 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]	58.85	Mem Busy [%]	9.54
L1/TEX Hit Rate [%]	70.84	Max Bandwidth [%]	18.41
L2 Hit Rate [%]	50.13	Mem Pipes Busy [%]	17.25



	Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	843,624	843,624	843,624	1.74	0
Shared Load Matrix	0	0	0	0	0
Shared Atomic	1,061,468	1,061,468	1,061,468	2.19	0
Other	-	-	420,684	0.87	0
Total	1,905,092	1,905,092	2,325,776	4.79	0

	Instructions	Requests	Wavefronts	% Peak	Sectors	Sectors/Req	Hit Rate	Bytes	Sector Misses to L2	% Peak to L2	Returns to SM
Global Load	1,468,017	1,468,017	1,468,017	3.02	3,998,881	2.72	58.94	127,964,192	1,631,748	3.38	1,468,017
Surface Load	0	0	0	0	0	0	0	0	0	0	0
Texture Load	0	0	0	0	0	0	0	0	0	0	0
Global Store	407,936	407,936	407,936	0.84	1,631,744	4	100	52,215,808	1,631,744	3.36	-
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 Reduction	0	0	0	0	0	0	0	0	0	0	-
DSMEM Reduction	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	see above
Surface Atomic ALU	0	0	0	0	0	0	0	0	0	0	see above
Surface Atomic CAS	0	0	0	0	0	0	0	0	0	0	-
Loads	1,468,017	1,468,017	1,468,017	3.02	3,998,881	2.72	58.94	127,964,192	1,631,748	3.38	1,468,017
Stores	407,936	407,936	407,936	0.84	1,631,744	4	100	52,215,808	1,631,744	3.36	-
Arithmetics & Reductions	0	0	0	0	0	0	0	0	0	0	-
Total	1,875,953	1,875,953	1,875,953	3.86	5,630,625	3.00	70.84	180,180,000	3,273,552	6.74	1,468,017

	Requests	Sectors	Sectors/Req	% Peak	Hit Rate	Bytes	Throughput	Sector Misses to Device	Sector Misses to System	Sector Misses to Peer
L1/TEX Load	410,512	1,641,808	4.00	2.89	0.61	52,537,856	25,303,352,084.46	1,631,748	0	0
L1/TEX Store	407,936	1,631,744	4	2.87	100	52,215,808	25,148,246,898.36	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	818,449	3,273,552	4.00	5.76	50.15	104,753,664	50,451,598,982.82	1,631,748	0	0
EOC Total	-	8	-	0.00	-	256	123,295.06	8	-	-
GPU Total	819,225	3,277,142	4.00	5.77	50.17	104,868,544	50,506,927,641.21	1,631,759	0	0

	Sectors	% Peak	Bytes	Throughput
Load	2,017,969	9.73	64,575,008	31,100,701,240.66
Store	1,800,256	8.68	57,608,256	27,745,364,873.32
Total	3,818,227	18.41	122,183,264	58,846,066,113.98

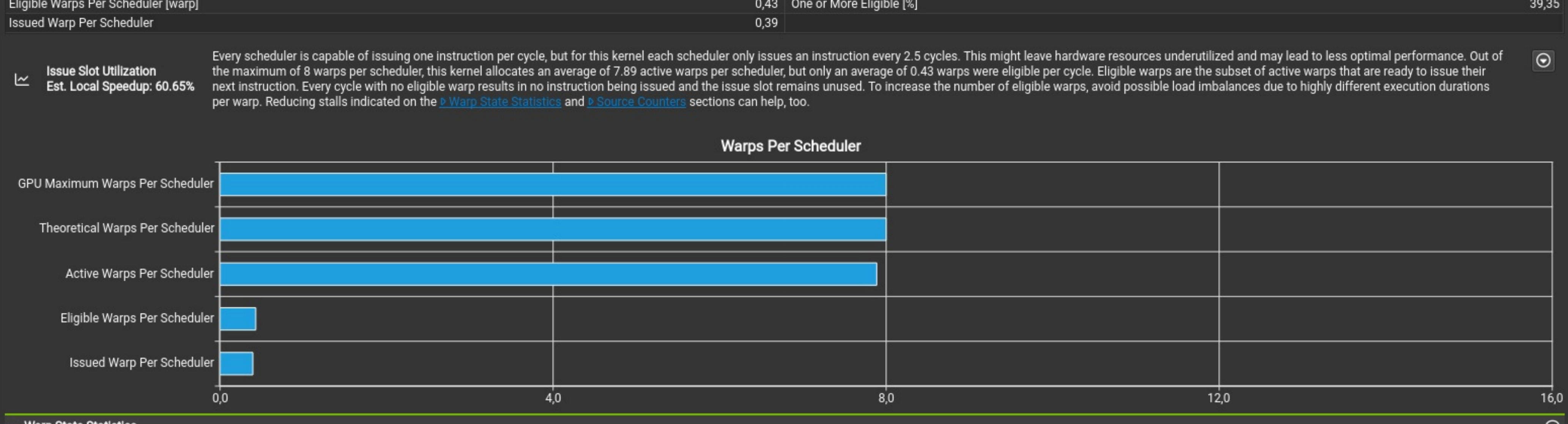
Scheduler Statistics

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 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]	7.89	No Eligible [%]	60.65
Eligible Warps Per Scheduler [warp]	0.43	One or More Eligible [%]	39.35
Issued Warp Per Scheduler	0.39		

Issue Slot Utilization Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 2.5 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 7.89 active warps per scheduler, but only an average of 0.43 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 Sampling](#) and [Source Doubles](#) sections can help too.



Warp State Statistics

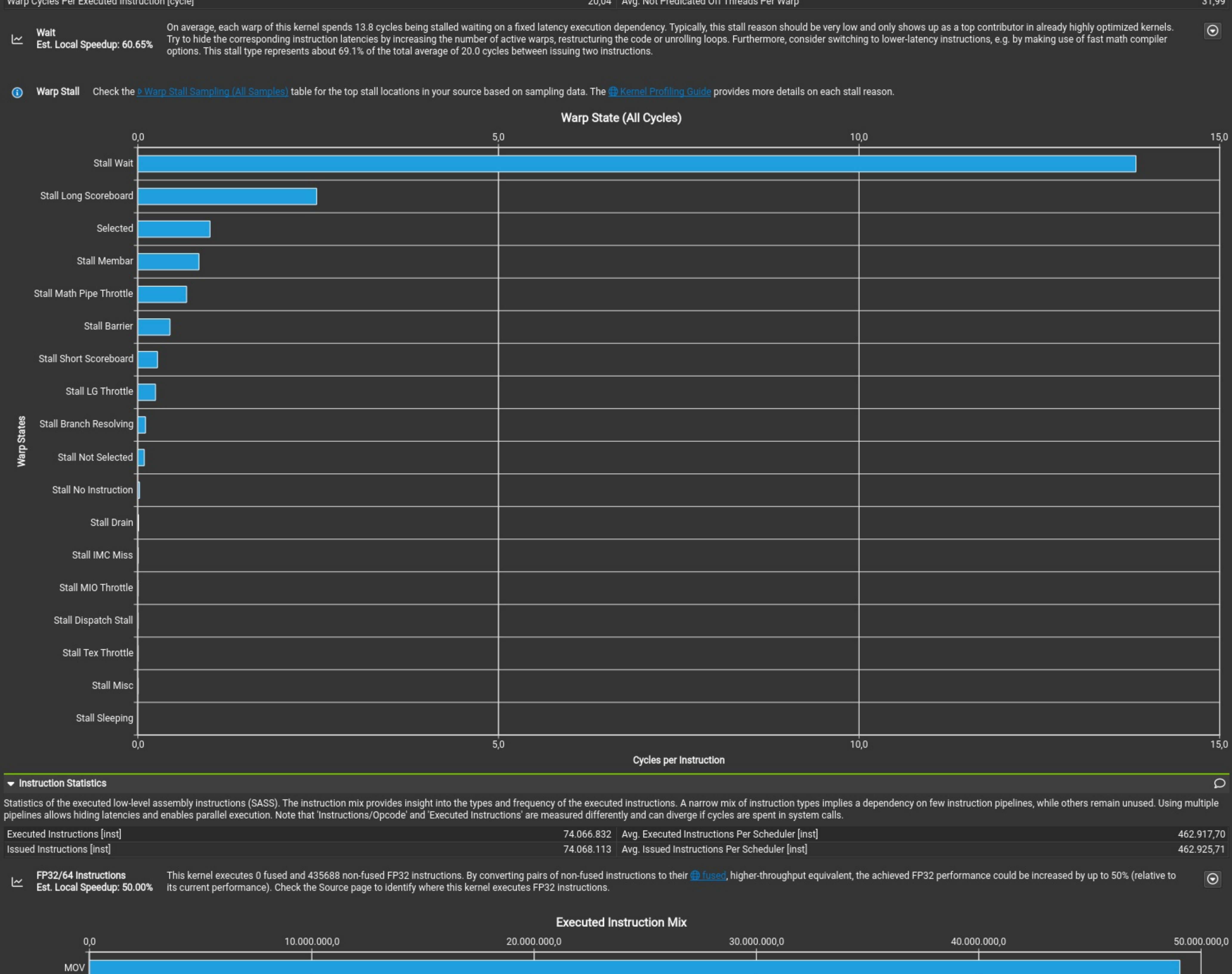
All

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 Issued Instruction [cycle]	20.04	Avg. Stalls Per Warp	32.00
Warp Cycles Per Executed Instruction [cycle]	20.04	Avg. Not Predicted Off Threads Per Warp	31.99

Wait On average, each warp of this kernel spends 19.8 cycles being stalled waiting on a fixed library execution dependency. Typically, this reason should be very low and only shows up as a top contributor in already highly optimized kernels. Try to hide the corresponding instruction latencies by increasing the number of active warps, restructuring the code or unrolling loops. Furthermore, consider switching to lower-latency instructions, e.g. by making use of fast math compiler options. This stall type represents about 69.1% of the total average of 20.0 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.



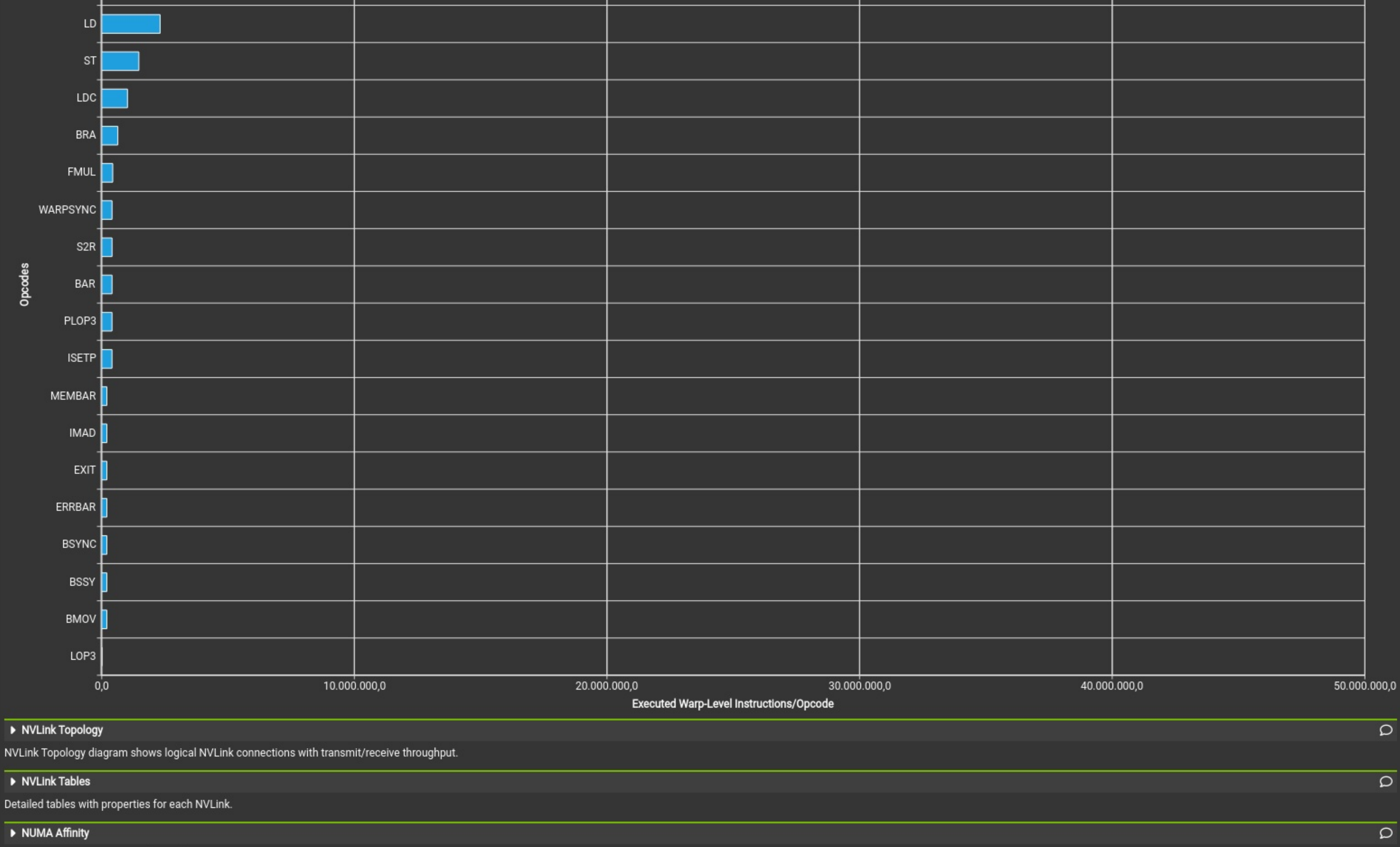
Instruction Statistics

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]	74,066,832	Avg. Executed Instructions Per Scheduler [Inst]	462,917.70
Issued Instructions [Inst]	74,068,113	Avg. Issued Instructions Per Scheduler [Inst]	462,925.71

FP32/64 Instructions This kernel executes 0 fused and 43568 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

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	6,374	Function Cache Configuration	CachePreferNone
Registers Per Thread [register/thread]	25	Static Shared Memory Per Block [Kbyte/block]	3,200
Block Size	1,024	Dynamic Shared Memory Per Block [Kbyte/block]	0
Threads [thread]	6,526,976	Driver Shared Memory Per Block [Kbyte/block]	0
Waves Per SM	159,35	Shared Memory Configuration Size [Kbyte]	32,77

Occupancy

All

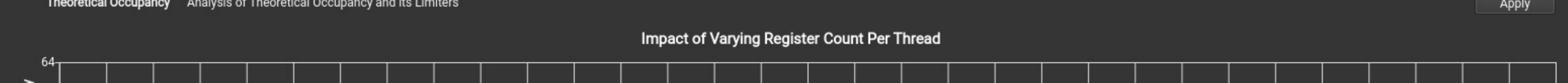
Occupancy is the ratio of the number of active warps per multiprocessor to the maximum number of possible active warps. Another way to view occupancy is the percentage of the hardware's ability to process warps that is actively in use. Higher occupancy does not always result in higher performance, however, 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 [%]	100	Block Limit Registers [block]	2
Theoretical Active Warps per SM [warp]	32	Block Limit Shared Mem [block]	3
Achieved Occupancy [%]	98.50	Block Limit Warps [block]	1
Achieved Active Warps Per SM [warp]	31.52	Block Limit SM [block]	16

Occupancy Limiters This kernel's theoretical occupancy is not impacted by any block limit.

Achieved Occupancy Analysis of the Achieved Occupancy

Theoretical Occupancy Analysis of Theoretical Occupancy and its Limiters

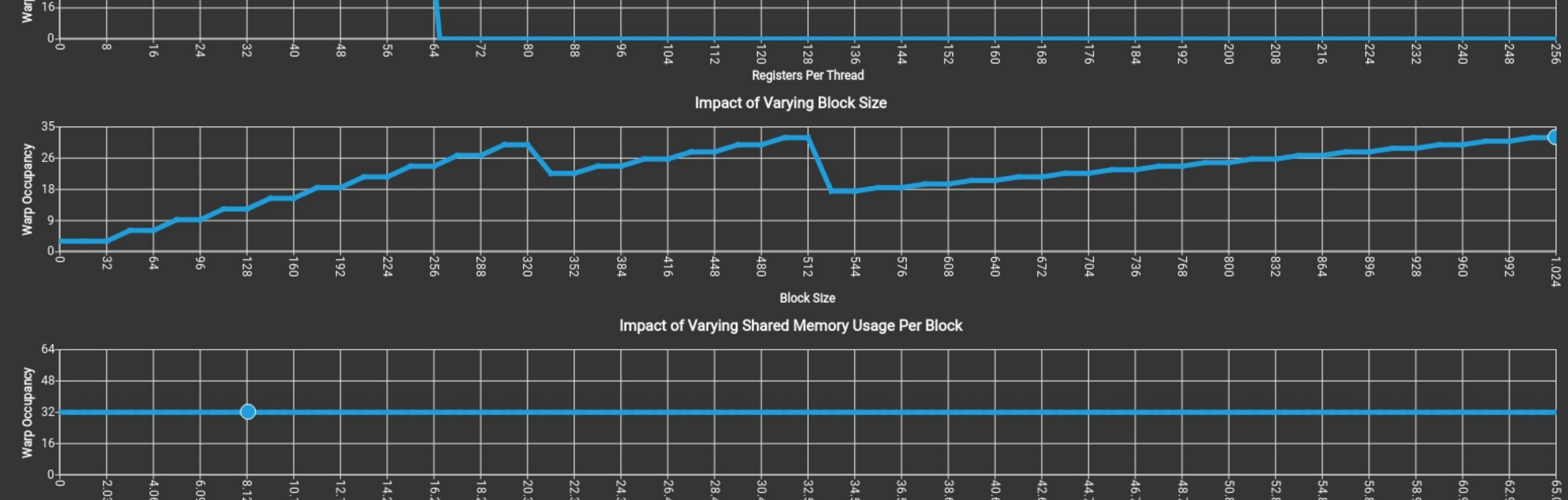


Source Counters

All

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]	1,451,369	Branch Efficiency [%]	100.00
Branch Instructions Ratio [%]	0.02	Arch. Divergent Branches	0.01



Location	Value	Value (%)	Location	Value	Value (%)
jnvwa0.cu41.0x7405cd033800 in applyMtx...	2,350	0.00	jnvwa0.cu43.0x7405cd032490 in applyMtx...	203,959	0.00
jnvwa0.cu40.0x7405cd035520 in applyMtx...	2,143	0.00	jnvwa0.cu40.0x7405cd035520 in applyMtx...	203,968	0.00
jnvwa0.cu42.0x7405cd033460 in applyMtx...	2,059	0.00	jnvwa0.cu40.0x7405cd033520 in applyMtx...	203,968	0.00
jnvwa0.cu40.0x7405cd033460 in applyMtx...	720	0.00	jnvwa0.cu40.0x7405cd033520 in applyMtx...	203,968	0.00
jnvwa0.cu40.0x7405cd033520 in applyMtx...	456	0.00	jnvwa0.cu40.0x7405cd033520 in applyMtx...	203,968	0.00

To augment your report even further, you might want to learn about [system metrics](#) and [setting your own rules](#). You might also want to consider [providing individual metrics](#).