

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.

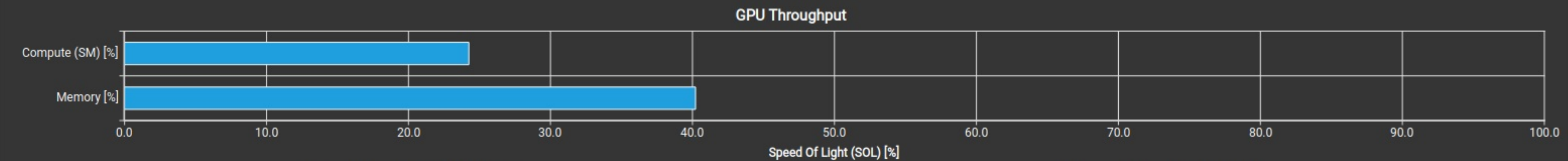
Compute (SM) Throughput [%]	24.25	Duration [us]	125.18
Memory Throughput [%]	40.22	Elapsed Cycles [cycle]	73,219
L1/TEX Cache Throughput [%]	33.29	SM Active Cycles [cycle]	66,013.80
L2 Cache Throughput [%]	11.61	SM Frequency [Mhz]	584.87
DRAM Throughput [%]	40.22	DRAM Frequency [Ghz]	4.97

Latency Issue

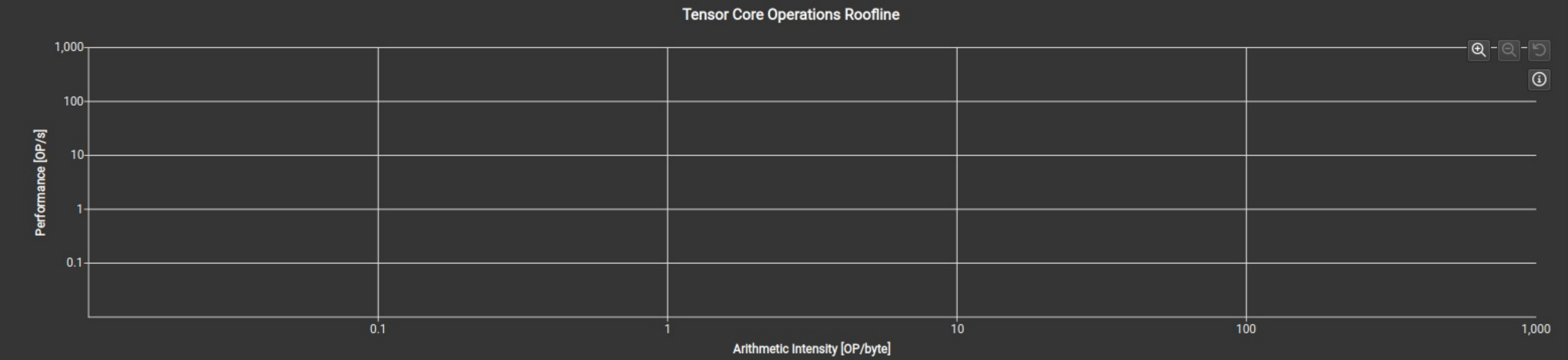
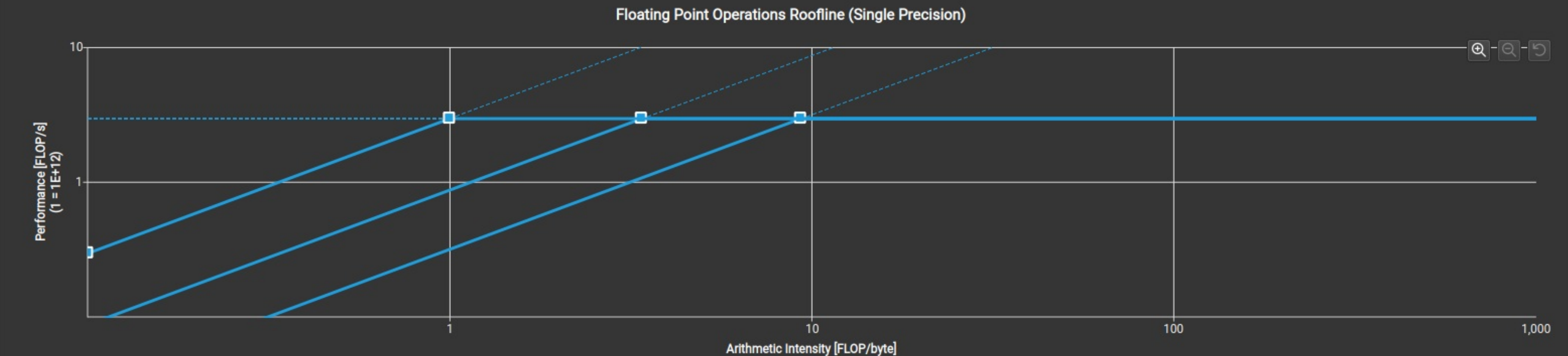
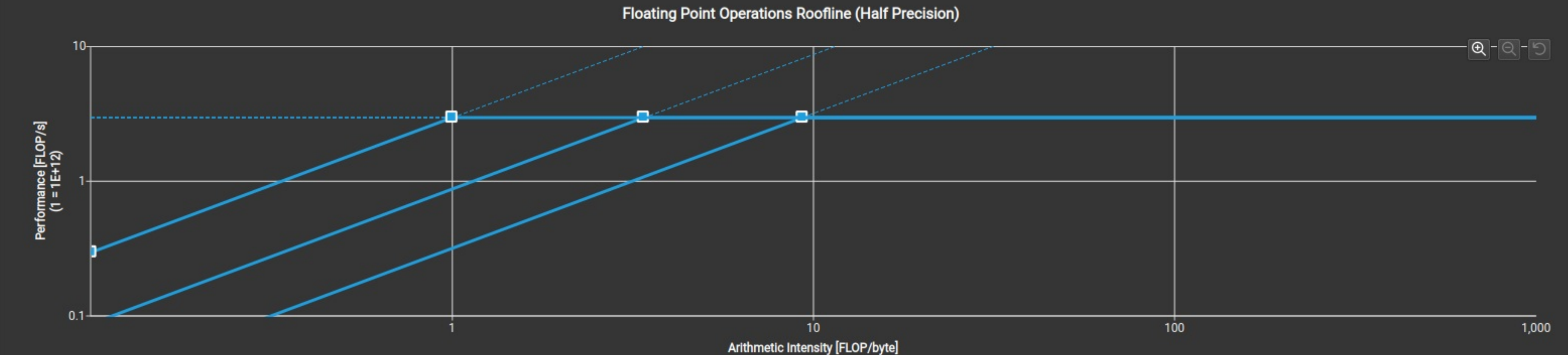
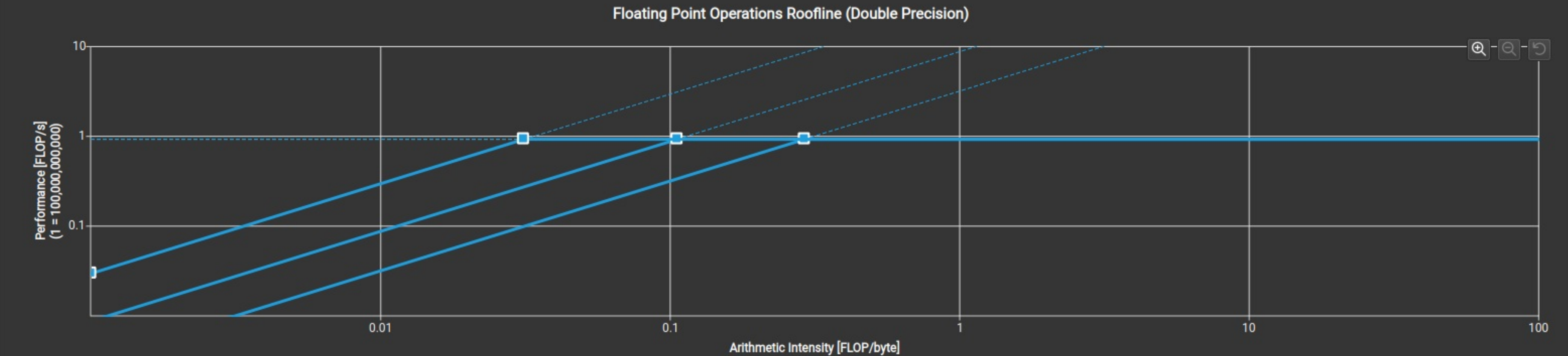
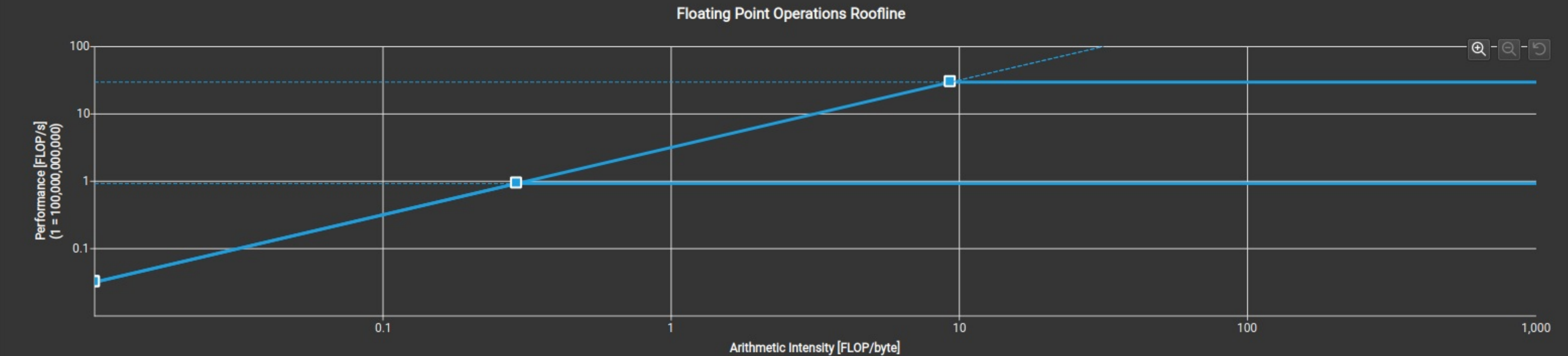
This kernel exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of this device. Achieved compute throughput and/or memory bandwidth below 60.0% of peak typically indicate latency issues. Look at [Scheduler Statistics](#) and [a Warp State Statistics](#) for potential reasons.

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.



Compute Throughput Breakdown		Memory Throughput Breakdown	
SM: Inst Executed Pipe Lsu [%]	24.25	DRAM: Cycles Active [%]	40.22
SM: Issue Active [%]	10.52	DRAM: Dram Sectors [%]	29.21
SM: Inst Executed [%]	10.46	L1: Lsuin Requests [%]	24.25
SM: Mio Inst Issued [%]	9.68	L1: Data Pipe Lsu Wavefronts [%]	16.64
SM: Pipe Alu Cycles Active [%]	6.18	L1: M Xbar2l1tex Read Sectors [%]	13.58
SM: Mio Pq Write Cycles Active [%]	5.38	L2: T Sectors [%]	11.61
SM: Mio Pq Read Cycles Active [%]	5.32	L2: Lts2xbar Cycles Active [%]	11.54
SM: Inst Executed Pipe Adu [%]	4.80	L2: D Sectors Fill Device [%]	11.50
SM: Mio2rf Writeback Active [%]	4.52	L1: Lsu Writeback Active [%]	7.90
SM: Pipe Fma Cycles Active [%]	3.85	L2: D Sectors [%]	6.02
SM: Inst Executed Pipe Cbu Pred On Any [%]	2.58	L1: Data Bank Writes [%]	3.76
SM: Inst Executed Pipe Uniform [%]	1.13	L1: Data Bank Reads [%]	3.07
SM: Memory Throughput Internal Activity [%]	0	GPU: Compute Memory Access Throughput Internal Activity [%]	2.98
SM: Pipe Tensor Cycles Active [%]	0	L2: T Tag Requests [%]	2.96
SM: Pipe Shared Cycles Active [%]	0	L1: M L1tex2xbar Req Cycles Active [%]	1.92
SM: Pipe Fp64 Cycles Active [%]	0	L2: Xbar2lts Cycles Active [%]	1.86
IDC: Request Cycles Active [%]	0	L1: Texin Sm2tex Req Cycles Active [%]	0.36
SM: Instruction Throughput Internal Activity [%]	0	L1: F Wavefronts [%]	0.00
SM: Inst Executed Pipe Xu [%]	0	L2: D Atomic Input Cycles Active [%]	0
SM: Inst Executed Pipe Tex [%]	0	L1: Tex Writeback Active [%]	0
SM: Inst Executed Pipe Ipa [%]	0	L2: D Sectors Fill Sysmem [%]	0
SM: Inst Executed Pipe Fp16 [%]	0	L1: Data Pipe Tex Wavefronts [%]	0
		GPU: Compute Memory Request Throughput Internal Activity [%]	0



	# Operations	# Operations / Cycle	# Operations / s	Peak %	Peak Operations / Cycle	Peak Operations / s
Src:fp16,bf16,tf32 Dst:fp32	0	0	0	0	40,960	23,742.23
Src:fp16 Dst:fp16	0	0	0	0	40,960	23,742.23
Src:int1	0	0	0	0	81,920	47,484.47
Src:int4	0	0	0	0	81,920	47,484.47
Src:int8	0	0	0	0	81,920	47,484.47

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]	66,013.80	Average L1 Active Cycles [cycle]	66,013.80
Average L2 Active Cycles [cycle]	76,641.19	Average SMSP Active Cycles [cycle]	44,283.61
Average DRAM Active Cycles [cycle]	250,079.50	Total SM Elapsed Cycles [cycle]	2,902,488
Total L1 Elapsed Cycles [cycle]	2,902,488	Total L2 Elapsed Cycles [cycle]	3,424,032
Total SMSP Elapsed Cycles [cycle]	11,609,952	Total DRAM Elapsed Cycles [cycle]	4,974,592

L2 Slices Workload Imbalance

Est. Speedup: 5.44%

One or more L2 Slices have a much higher number of active cycles than the average number of active cycles. Maximum instance value is 7.59% above the average, while the minimum instance value is 3.46% below the average.

Workload Distribution

	Average	Min	Max	Sum
SM Active Cycles	66,013.80	64,872	67,406	2,640,552
SMSP Active Cycles	44,283.61	42,615	46,415	7,085,378
L1 Active Cycles	66,013.80	64,872	67,406	2,640,552
L2 Active Cycles	76,641.19	73,993	82,939	2,452,518
DRAM Active Cycles	250,079.50	249,132	251,116	2,000,636