

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.

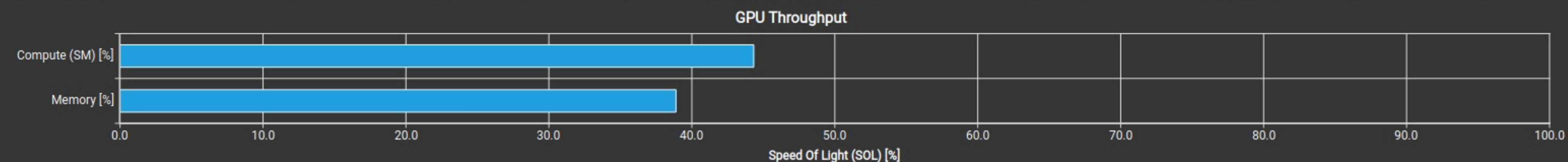
Compute (SM) Throughput [%]	44.33	Duration [ms]	4.22
Memory Throughput [%]	38.89	Elapsed Cycles [cycle]	2,468,445
L1/TEX Cache Throughput [%]	50.78	SM Active Cycles [cycle]	2,407,309.05
L2 Cache Throughput [%]	18.37	SM Frequency [Mhz]	584.99
DRAM Throughput [%]	30.51	DRAM Frequency [Ghz]	5.00

Latency Issue

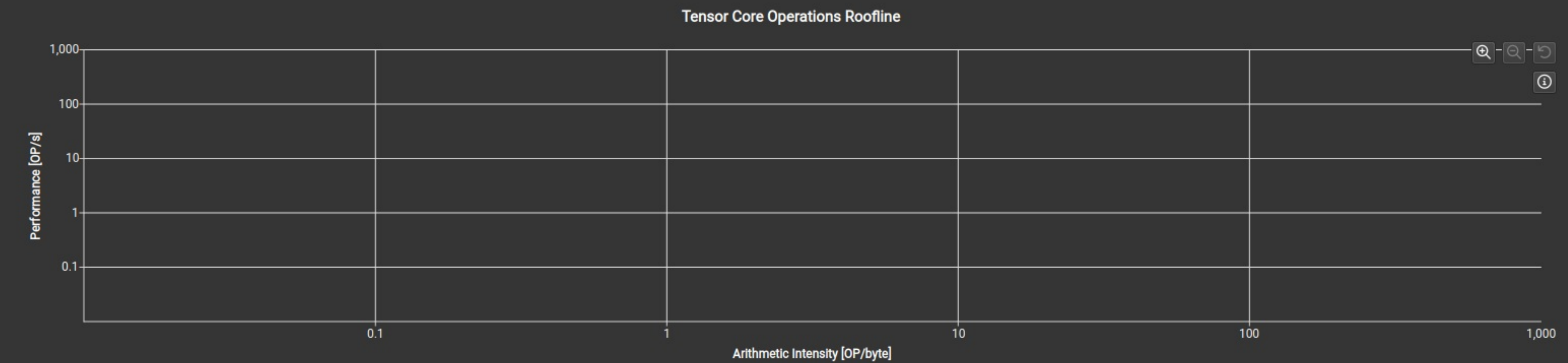
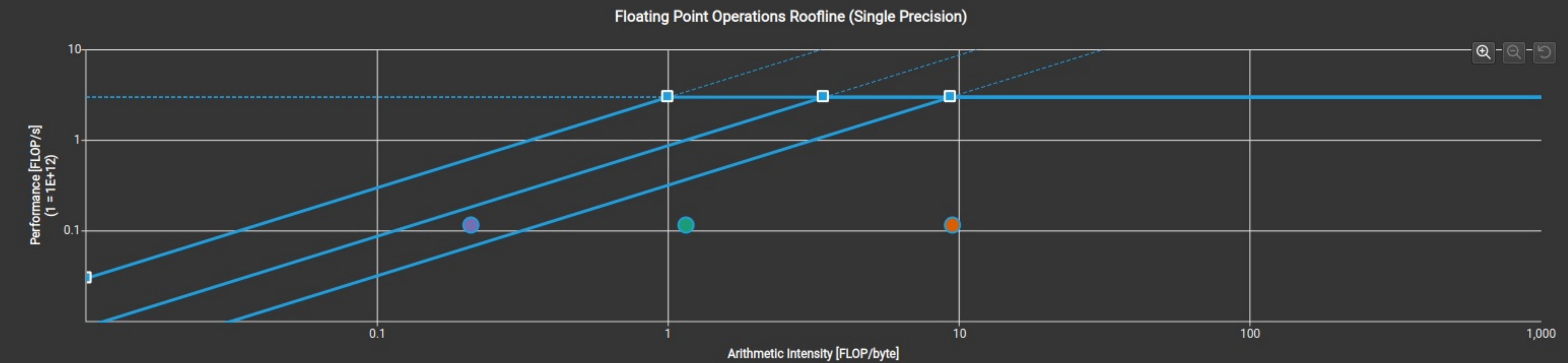
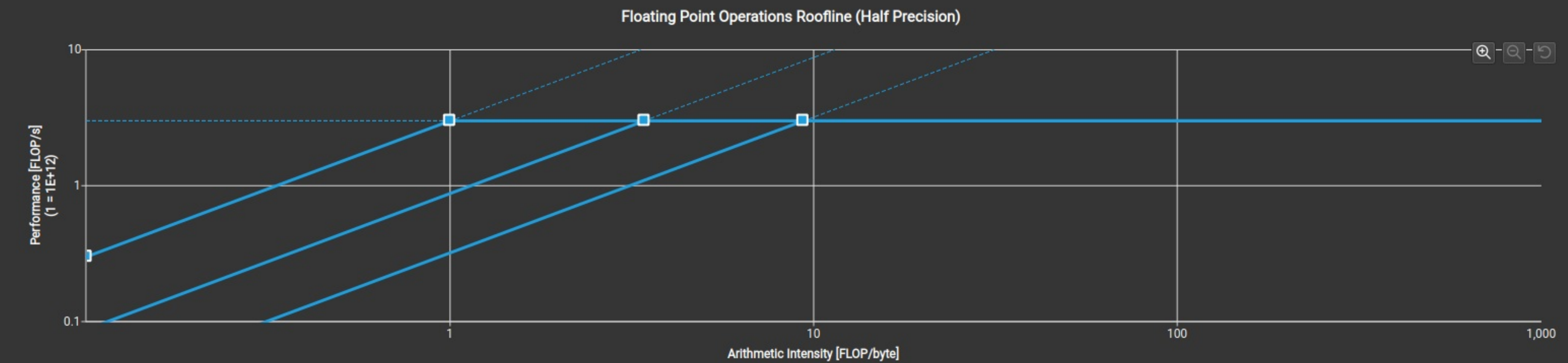
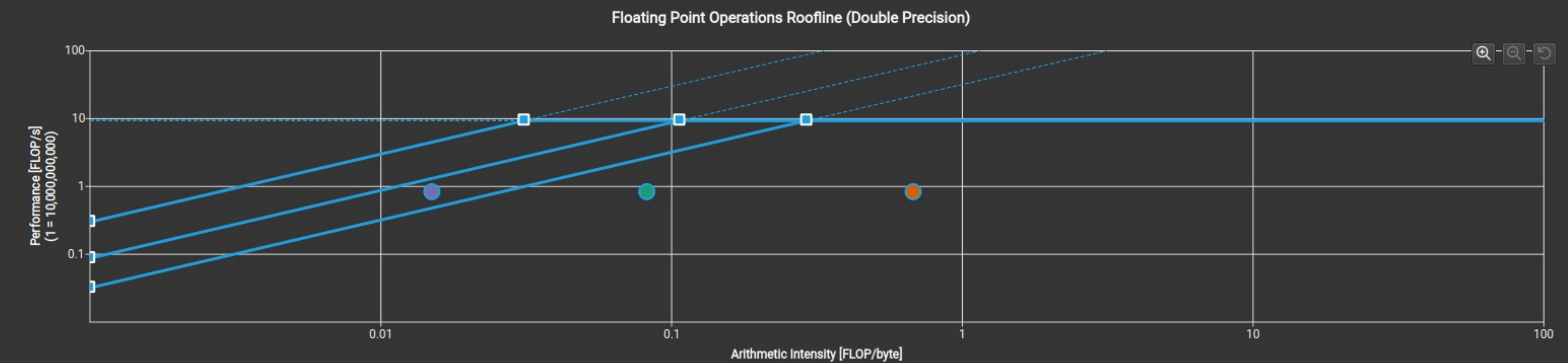
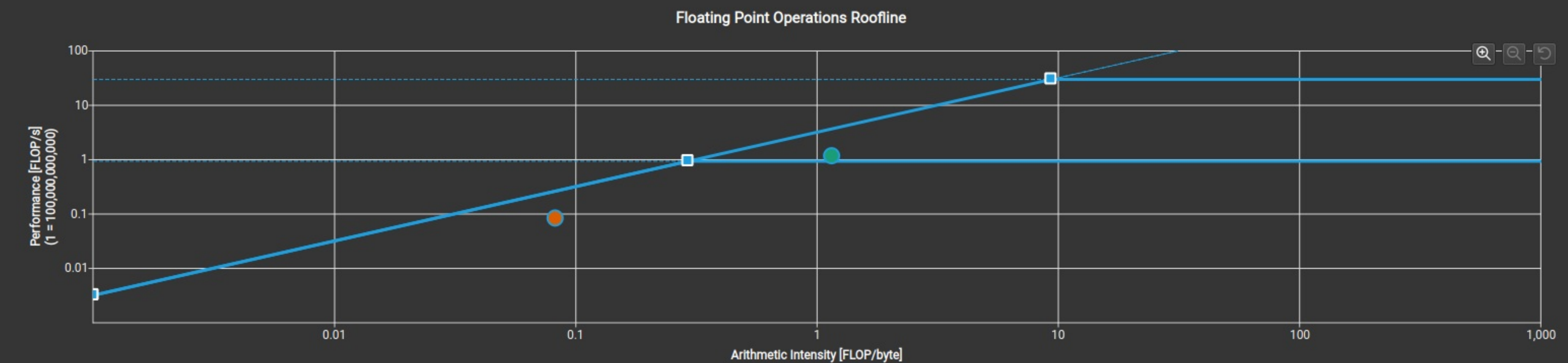
This kernel exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of this device. Achieved compute throughput and/or memory bandwidth below 60.0% of peak typically indicate latency issues. Look at [Scheduler Statistics](#) and [Warp State Statistics](#) for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

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 4% of this device's fp32 peak performance and 9% of its fp64 peak performance. See the [Kernel Profiling Guide](#) for more details on roofline analysis.



Compute Throughput Breakdown			Memory Throughput Breakdown		
SM: Pipe Fp64 Cycles Active [%]	44.33	L1: Lsuin Requests [%]	38.89		
SM: Inst Executed Pipe Lsu [%]	38.89	DRAM: Cycles Active [%]	30.51		
SM: Issue Active [%]	30.40	L1: Data Pipe Lsu Wavefronts [%]	25.39		
SM: Inst Executed [%]	30.39	DRAM: Dram Sectors [%]	22.32		
SM: Pipe Alu Cycles Active [%]	20.89	L1: M L1tex2xbar Req Cycles Active [%]	21.11		
SM: Pipe Fma Cycles Active [%]	20.28	L2: T Sectors [%]	18.37		
SM: Mio Inst Issued [%]	14.88	L2: Xbar2lts Cycles Active [%]	18.08		
SM: Inst Executed Pipe Xu [%]	13.67	L1: Lsu Writeback Active [%]	17.91		
SM: Mio Pq Write Cycles Active [%]	10.26	L2: D Sectors [%]	7.02		
SM: Mio2rf Writeback Active [%]	9.00	L1: Data Bank Reads [%]	6.80		
SM: Mio Pq Read Cycles Active [%]	3.38	GPU: Compute Memory Access Throughput Internal Activity [%]	4.81		
SM: Inst Executed Pipe Cbu Pred On Any [%]	2.42	L2: T Tag Requests [%]	4.79		
SM: Inst Executed Pipe Adu [%]	0.20	L1: Data Bank Writes [%]	2.68		
IDC: Request Cycles Active [%]	0.14	L1: M Xbar2l1tex Read Sectors [%]	1.57		
SM: Inst Executed Pipe Uniform [%]	0.06	L2: Lts2xbar Cycles Active [%]	1.35		
SM: Pipe Tensor Cycles Active [%]	0	L2: D Sectors Fill Device [%]	0.59		
SM: Inst Executed Pipe Fp16 [%]	0	L1: Texin Sm2tex Req Cycles Active [%]	0.00		
SM: Pipe Shared Cycles Active [%]	0	L1: F Wavefronts [%]	0.00		
SM: Memory Throughput Internal Activity [%]	0	GPU: Compute Memory Request Throughput Internal Activity [%]	0		
SM: Instruction Throughput Internal Activity [%]	0	L1: Tex Writeback Active [%]	0		
SM: Inst Executed Pipe Tex [%]	0	L2: D Atomic Input Cycles Active [%]	0		
SM: Inst Executed Pipe Ipa [%]	0	L2: D Sectors Fill Sysmem [%]	0		
		L1: Data Pipe Tex Wavefronts [%]	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	24,001.15
Src:fp16 Dst:fp16	0	0	0	0	40,960	24,001.15
Src:int1	0	0	0	0	81,920	48,002.31
Src:int4	0	0	0	0	81,920	48,002.31
Src:int8	0	0	0	0	81,920	48,002.31

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]	2,407,309.05	Average L1 Active Cycles [cycle]	2,407,309.05
Average L2 Active Cycles [cycle]	2,963,807.72	Average SMSP Active Cycles [cycle]	2,308,630.13
Average DRAM Active Cycles [cycle]	6,443,149.50	Total SM Elapsed Cycles [cycle]	98,902,008
Total L1 Elapsed Cycles [cycle]	98,902,008	Total L2 Elapsed Cycles [cycle]	115,446,560
Total SMSP Elapsed Cycles [cycle]	395,608,032	Total DRAM Elapsed Cycles [cycle]	168,921,088

Workload Distribution				
	Average	Min	Max	Sum
SM Active Cycles	2,407,309.05	2,356,849	2,466,719	96,292,362
SMSP Active Cycles	2,308,630.13	2,222,743	2,397,477	369,380,821
L1 Active Cycles	2,407,309.05	2,356,849	2,466,719	96,292,362
L2 Active Cycles	2,963,807.72	2,846,578	3,063,791	94,841,847
DRAM Active Cycles	6,443,149.50	6,318,868	6,588,460	51,545,196