

GPU Speed Of Light Throughput

All

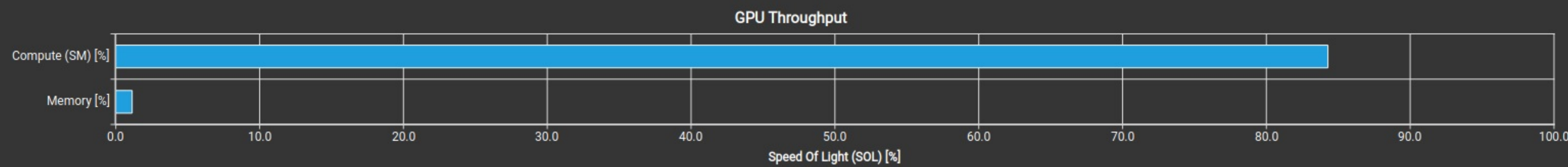
High-level overview of the throughput for compute and memory resources of the GPU. For each unit, the throughput reports the achieved percentage of utilization with respect to the theoretical maximum. Breakdowns show the throughput for each individual sub-metric of Compute and Memory to clearly identify 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 [%] | 84.28 | Duration [ms]            | 1.10       |
| Memory Throughput [%]       | 1.14  | Elapsed Cycles [cycle]   | 642,614    |
| L1/TEX Cache Throughput [%] | 1.02  | SM Active Cycles [cycle] | 623,294.78 |
| L2 Cache Throughput [%]     | 0.45  | SM Frequency [Mhz]       | 584.97     |
| DRAM Throughput [%]         | 1.14  | DRAM Frequency [Ghz]     | 4.99       |

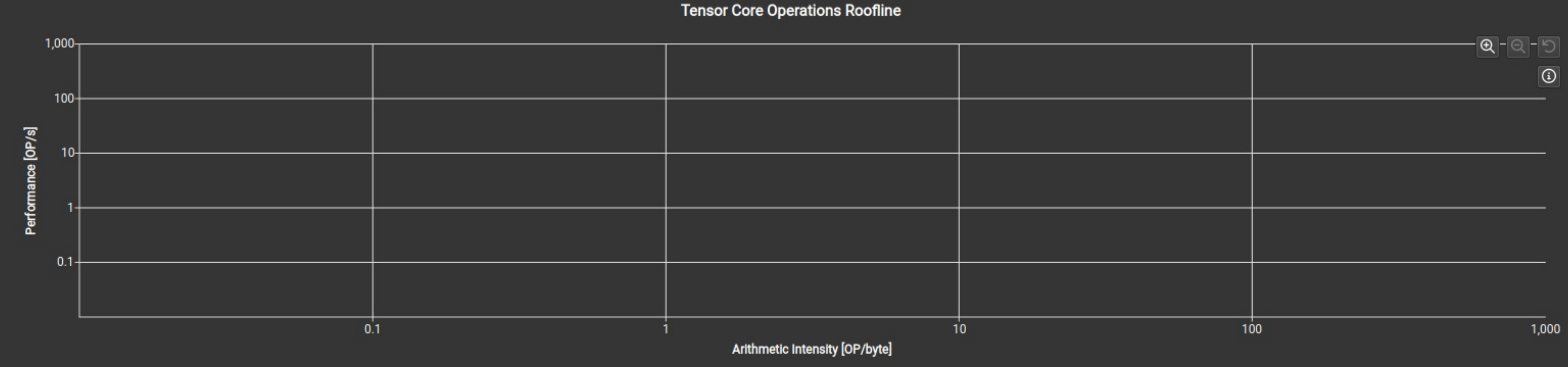
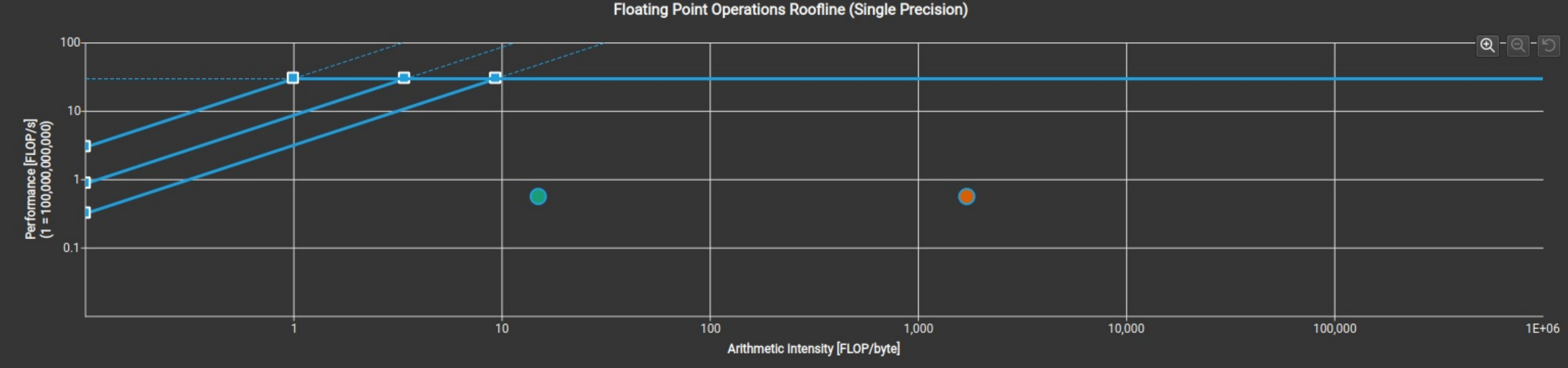
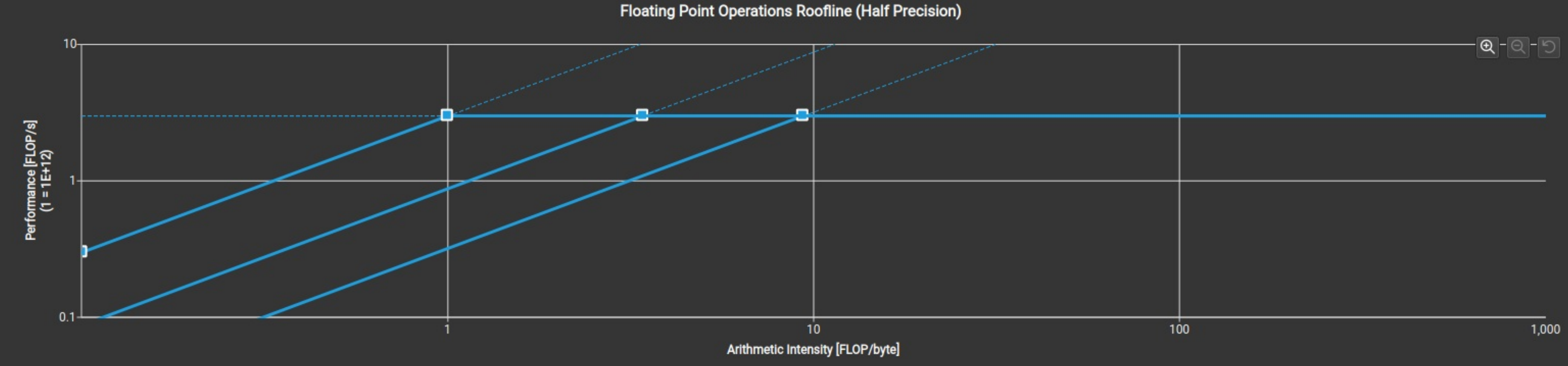
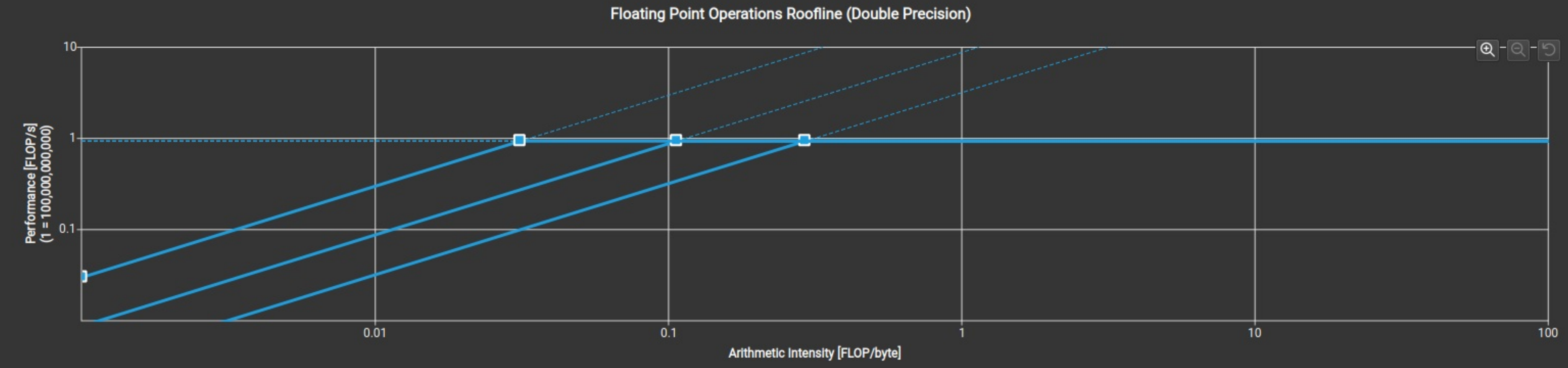
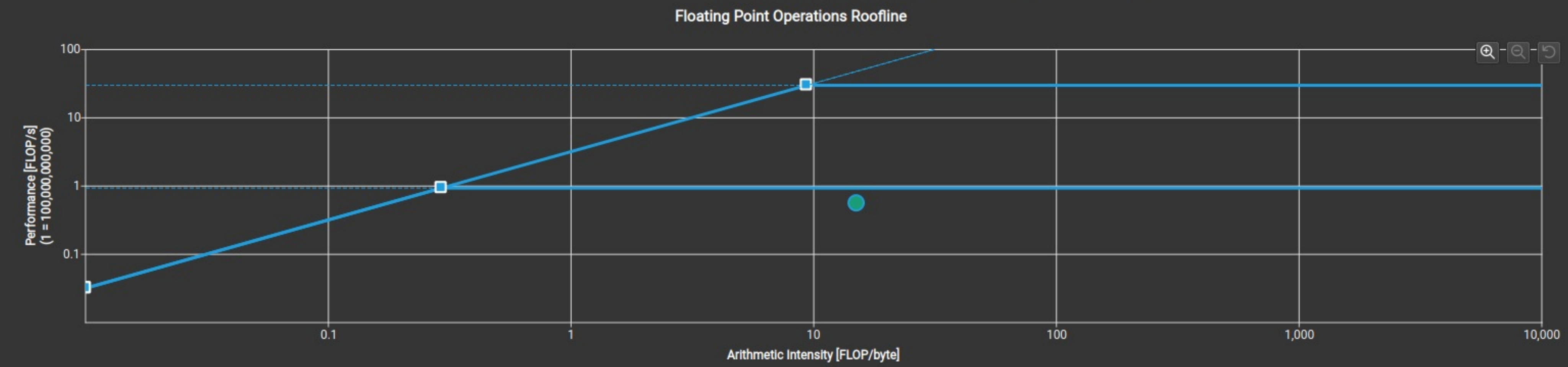
- High Throughput

The kernel is utilizing greater than 80.0% of the available compute or memory performance of the device. To further improve performance, work 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 32.1. The kernel achieved 2% 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: Pipe Alu Cycles Active [%]                   | 84.28 | DRAM: Cycles Active [%]                                      | 1.14                        |  |  |
| SM: Issue Active [%]                             | 62.26 | DRAM: Dram Sectors [%]                                       | 0.83                        |  |  |
| SM: Inst Executed [%]                            | 62.25 | L1: M L1tex2xbar Req Cycles Active [%]                       | 0.51                        |  |  |
| SM: Inst Executed Pipe Adu [%]                   | 38.54 | L1: Lsuin Requests [%]                                       | 0.51                        |  |  |
| IDC: Request Cycles Active [%]                   | 22.02 | L2: T Sectors [%]  | 0.45                        |  |  |
| SM: Pipe Fma Cycles Active [%]                   | 17.41 | L2: Xbar2lts Cycles Active [%]                               | 0.44                        |  |  |
| SM: Mio Inst Issued [%]                          | 13.02 | L1: Data Pipe Lsu Wavefronts [%]                             | 0.26                        |  |  |
| SM: Mio Pq Read Cycles Active [%]                | 11.14 | L2: D Sectors [%]  | 0.19                        |  |  |
| SM: Mio Pq Write Cycles Active [%]               | 11.14 | L1: Lsu Writeback Active [%]                                 | 0.13                        |  |  |
| SM: Inst Executed Pipe Cbu Pred On Any [%]       | 8.42  | GPU: Compute Memory Access Throughput Internal Activity [%]  | 0.11                        |  |  |
| SM: Mio2rf Writeback Active [%]                  | 5.63  | L2: T Tag Requests [%]                                       | 0.11                        |  |  |
| SM: Inst Executed Pipe Lsu [%]                   | 0.51  | L1: Data Bank Reads [%]                                      | 0.06                        |  |  |
| SM: Inst Executed Pipe Uniform [%]               | 0.06  | L1: Data Bank Writes [%]                                     | 0.06                        |  |  |
| SM: Memory Throughput Internal Activity [%]      | 0     | L2: Lts2xbar Cycles Active [%]                               | 0.00                        |  |  |
| SM: Pipe Tensor Cycles Active [%]                | 0     | L2: D Sectors Fill Device [%]                                | 0.00                        |  |  |
| SM: Pipe Shared Cycles Active [%]                | 0     | L1: Texin Sm2tex Req Cycles Active [%]                       | 0.00                        |  |  |
| SM: Pipe Fp64 Cycles Active [%]                  | 0     | L1: F Wavefronts [%]   | 0.00                        |  |  |
| SM: Instruction Throughput Internal Activity [%] | 0     | L1: M Xbar2l1tex Read Sectors [%]                            | 0                           |  |  |
| SM: Inst Executed Pipe Xu [%]                    | 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                           |  |  |
| 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                  | 23,946.76           |
| Src:fp16 Dst:fp16           | 0            | 0                    | 0                | 0      | 40,960                  | 23,946.76           |
| Src:int1                    | 0            | 0                    | 0                | 0      | 81,920                  | 47,893.53           |
| Src:int4                    | 0            | 0                    | 0                | 0      | 81,920                  | 47,893.53           |
| Src:int8                    | 0            | 0                    | 0                | 0      | 81,920                  | 47,893.53           |

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]   | 623,294.78  | Average L1 Active Cycles [cycle]   | 623,294.78 |
| Average L2 Active Cycles [cycle]   | 43,810.69   | Average SMSP Active Cycles [cycle] | 577,727.59 |
| Average DRAM Active Cycles [cycle] | 62,425.50   | Total SM Elapsed Cycles [cycle]    | 25,689,640 |
| Total L1 Elapsed Cycles [cycle]    | 25,689,640  | Total L2 Elapsed Cycles [cycle]    | 30,053,792 |
| Total SMSP Elapsed Cycles [cycle]  | 102,758,560 | Total DRAM Elapsed Cycles [cycle]  | 43,876,352 |

| Workload Distribution |            |         |         |            |
|-----------------------|------------|---------|---------|------------|
|                       | Average    | Min     | Max     | Sum        |
| SM Active Cycles      | 623,294.78 | 610,094 | 637,021 | 24,931,791 |
| SMSP Active Cycles    | 577,727.59 | 563,311 | 592,635 | 92,436,414 |
| L1 Active Cycles      | 623,294.78 | 610,094 | 637,021 | 24,931,791 |
| L2 Active Cycles      | 43,810.69  | 43,264  | 44,430  | 1,401,942  |
| DRAM Active Cycles    | 62,425.50  | 61,868  | 62,932  | 499,404    |
|                       |            |         |         |            |