

The report contains imported source files.

GPU Speed of Light Throughput

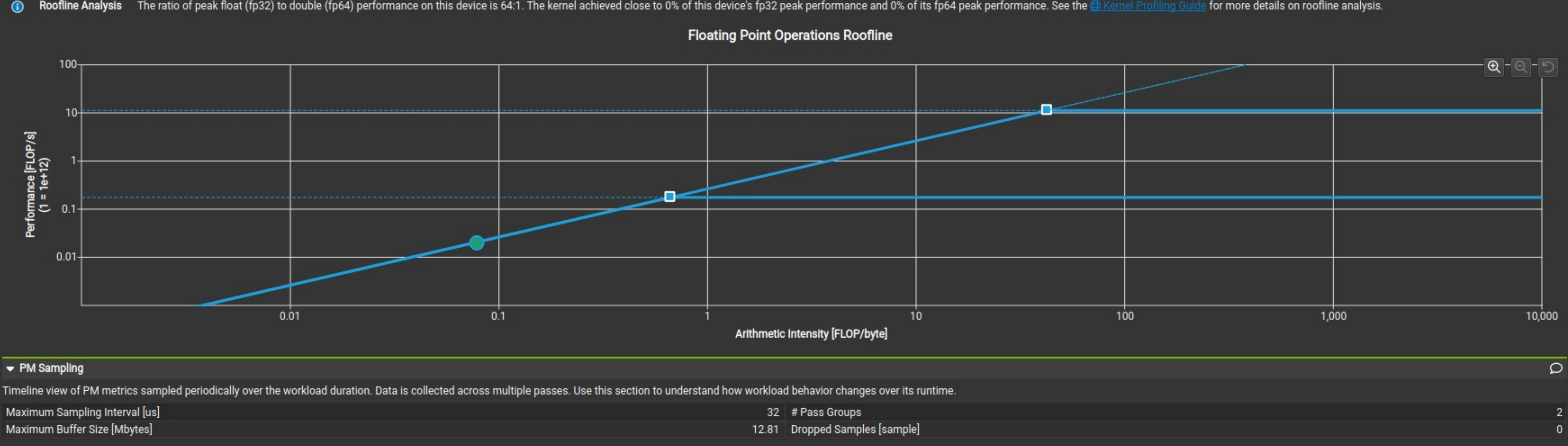
GPU Throughput Rooflines

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 contributors. High-level overview of the utilization for compute and memory resources of the GPU presented as a roofline chart.

|                             |       |                          |             |
|-----------------------------|-------|--------------------------|-------------|
| Compute (SM) Throughput [%] | 91.89 | Duration [ms]            | 10.44       |
| Memory Throughput [%]       | 91.89 | Elapsed Cycles [cycle]   | 19106262    |
| L1/TEX Cache Throughput [%] | 15.97 | SM Active Cycles [cycle] | 19061707.04 |
| L2 Cache Throughput [%]     | 30.72 | SM Frequency [GHz]       | 1.83        |
| DRAM Throughput [%]         | 91.89 | DRAM Frequency [GHz]     | 8.24        |

**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 DRAM in the [Memory Workload Analysis](#) section.

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

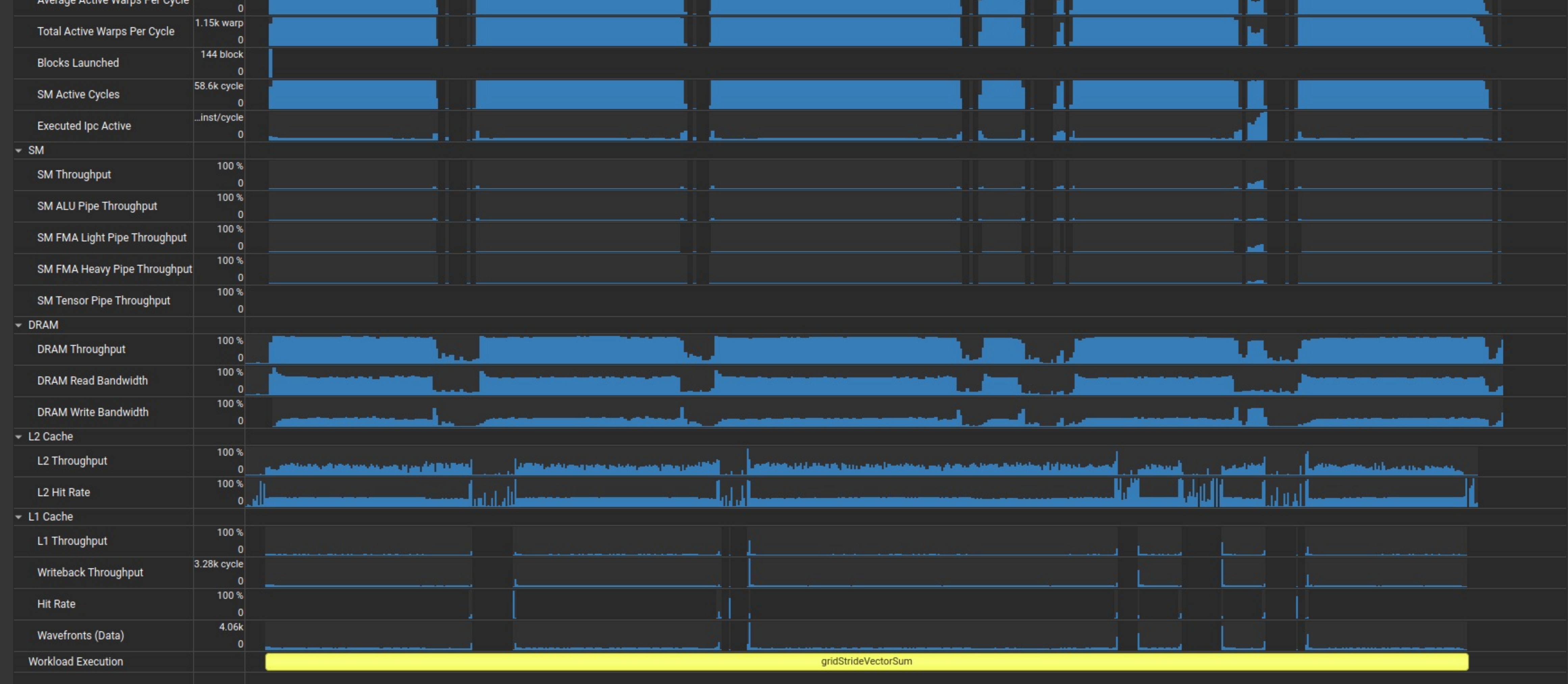


PM Sampling

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 [us] 22 # Pass Groups 2

Maximum Buffer Size [Mbytes] 12.81 Dropped Samples [sample] 0



Compute Workload Analysis

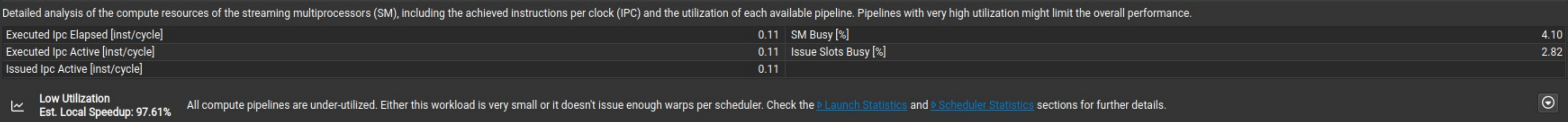
Detailed analysis of the compute resources of the streaming multiprocessors (SM), including the achieved instructions per clock (IPC) and the utilization of each available pipeline. Pipelines with very high utilization might limit the overall performance.

Executed IPC Elapsed [inst/cycle] 0.11 SM Busy [%] 4.10

Executed IPC Active [inst/cycle] 0.11 Issue Slots Busy [%] 2.92

Issued IPC Active [inst/cycle] 0.11

**Low Utilization** Est. Local Speedup: 97.61% All compute pipelines are under-utilized. Either this workload is very small or it doesn't issue enough warps per scheduler. Check the [Launch Statistics](#) and [Scheduler Statistics](#) sections for further details.



Memory Workload Analysis

Memory Chart

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 [byte/s] 242.39 Mem Busy [%] 30.72

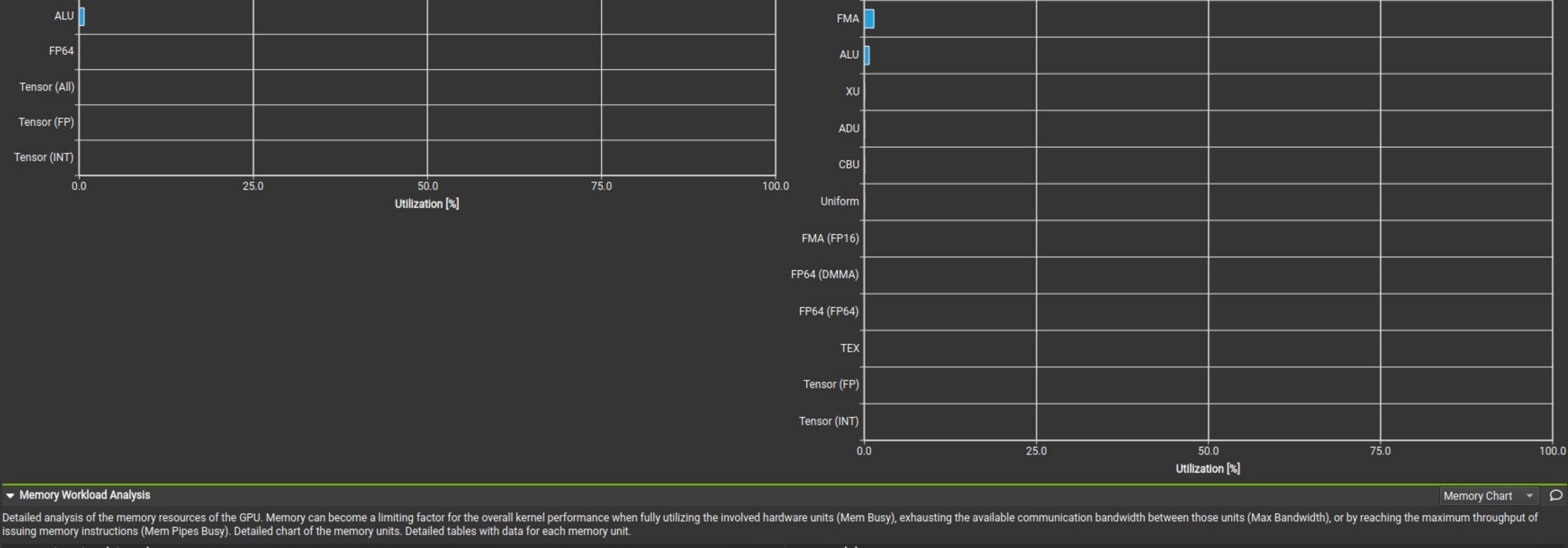
L1/TEX Hit Rate [%] 0 Max Bandwidth [%] 91.89

L2 Hit Rate [%] 31.51 Mem Pipes Busy [%] 8.16

L2 Compression Success Rate [%] 0 L2 Compression Ratio 0

Memory Chart

Show As: Transfer Size



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 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] 11.99 No Eligible [%] 97.15

Eligible Warps Per Scheduler [warp] 0.03 One or More Eligible [%] 2.85

Issued Warp Per Scheduler 0.03

**Issue Slot Utilization** Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 35.1 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 11.99 active warps per scheduler, but only an average of 0.03 warps were eligible per cycle. Eligible warps are eligible 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 Summary](#) and [Warp Stalls](#) sections can help, too.

Est. Local Speedup: 8.11%

Warp State Summary

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.

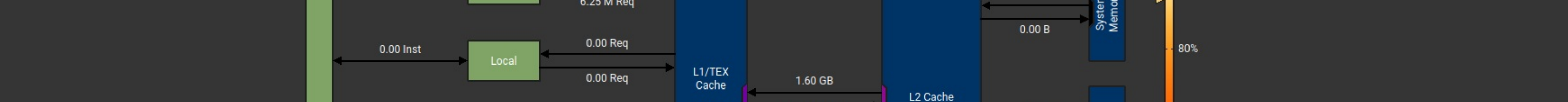
Warp Cycles Per Issued Instruction [cycle] 421.38 Avg. Active Threads Per Warp 32

Warp Cycles Per Executed Instruction [cycle] 421.37 Avg. Not Predicted Off Threads Per Warp 31.99

**Long Scoreboard Stalls** On average, each warp of this kernel spends 417.4 cycles being stalled waiting for a scoreboard dependency on a L1TEX (local, global, surface, texture) operation. Find the instruction producing the data being waited upon to identify the culprit. To reduce the number of cycles waiting on L1TEX data accesses verify the memory access patterns are optimal for the target architecture, attempt to increase cache hit rates by increasing data locality (coalescing), or by changing the cache configuration. Consider moving frequently used data to shared memory. This stall type represents about 99.1% of the total average of 421.3 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

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/Opcodes' and 'Executed Instructions' are measured differently and can diverge if cycles are spent in system calls.

Executed Instructions [inst] 51616304 Avg. Executed Instructions Per Scheduler [inst] 537669.83

Issued Instructions [inst] 51621580 Avg. Issued Instructions Per Scheduler [inst] 537724.79

**FP32 Non-Fused Instructions** This kernel executes 0 fused and 6250000 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.

Est. Speedup: 1.20%



NVLink Topology

NVLink Topology diagram shows logical NVLink connections with transmit/receive throughput.

The system does not have any NVLink connections.

NVLink Tables

Detailed tables with properties for each NVLink.

Logical NVLink Properties

The system does not have any NVLink connections.

NUMA Affinity

Non-uniform memory access (NUMA) affinities based on compute and memory distances for all GPUs.

NUMA ID Table

NUMA information is not available on the target system.

Launch Statistics

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 144 Partition Cache Configuration Cache/PreferNone

Registers Per Thread [register/thread] 26 Static Shared Memory Per Block [byte/block] 0

Block Size 256 Dynamic Shared Memory Per Block [kbyte/block] 0

Threads [thread] 36864 Driver Shared Memory Per Block [kbyte/block] 1.02

Waves Per SM 1 Shared Memory Configuration Size [kbyte] 16.38

Uses Green Context 0 # SMs [SM] 24

Occupancy

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.

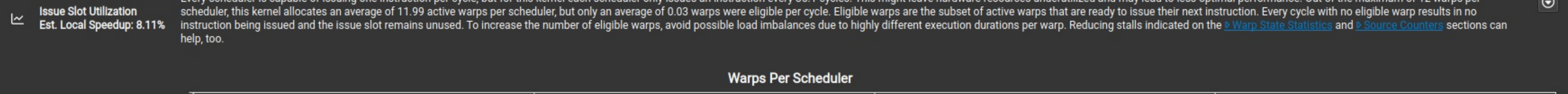
Theoretical Occupancy [%] 100 Block Limit Registers [block] 8

Theoretical Active Warps per SM [warp] 48 Block Limit Shared Mem [block] 16

Achieved Occupancy [%] 98.73 Block Limit Warps [block] 6

Achieved Active Warps per SM [warp] 47.39 Block Limit SM [block] 24

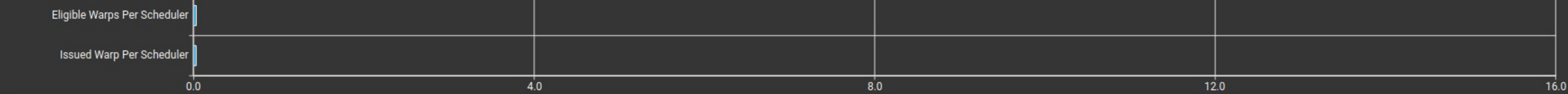
Impact of Varying Register Count Per Thread



Impact of Varying Block Size



Impact of Varying Shared Memory Usage Per Block



GPU and Memory Workload Distribution

Analysis of workload distribution in active cycles of SM, SMP, SMSP, L1 & L2 caches, and DRAM

SM Active Cycles [cycle] 19061707.04 Average L1 Active Cycles [cycle] 19061707.04

SMSP Active Cycles [cycle] 17244899.08 Average SMSP Active Cycles [cycle] 18895706.11

L1 Active Cycles [cycle] 79083516 Total SM Elapsed Cycles [cycle] 458550192

L2 Active Cycles [cycle] 19042522 Total L2 Elapsed Cycles [cycle] 204844116

DRAM Active Cycles [cycle] 17244899.08 Total DRAM Elapsed Cycles [cycle] 344236032

Workload Distribution

|                    | Average     | Min      | Max       | Sum        |
|--------------------|-------------|----------|-----------|------------|
| SM Active Cycles   | 19061707.04 | 19017568 | 19104252  | 457480969  |
| SMSP Active Cycles | 18895706.11 | 18824860 | 18954863  | 1613987787 |
| L1 Active Cycles   | 19042522    | 19017568 | 19104252  | 457480969  |
| L2 Active Cycles   | 17244899.08 | 17224729 | 205986669 | 344236032  |
| DRAM Active Cycles | 79083516    | 79044800 | 79103504  | 316334064  |

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] 1569424 Branch Efficiency [%] 100

Branch Instructions Ratio [%] 0.03 Avg. Divergent Branches 0

Warp Stall Sampling (All Samples)

| Location                                       | Value  | Value [%] | Location                                       | Value  | Value [%] |
|--|--------|-----------|--|--------|-----------|
| vectorSum.cu:35 (0x500674d4) in gridStrideVect | 107320 | 29        | vectorSum.cu:34 (0x50067500) in gridStrideVect | 150117 | 39        |
| vectorSum.cu:35 (0x500674e0) in gridStrideVect | 107208 | 28        | vectorSum.cu:35 (0x50067500) in gridStrideVect | 150211 | 39        |
| vectorSum.cu:35 (0x500674f0) in gridStrideVect | 107096 | 28        | vectorSum.cu:35 (0x50067500) in gridStrideVect | 150211 | 39        |
| vectorSum.cu:35 (0x50067500) in gridStrideVect | 106685 | 28        | vectorSum.cu:34 (0x50067500) in gridStrideVect | 150211 | 39        |
| vectorSum.cu:35 (0x500674e0) in gridStrideVect | 993    | 0         | vectorSum.cu:34 (0x50067500) in gridStrideVect | 150211 | 39        |

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 individual sections to focus on selected performance aspects and make profiling faster.