

▼ Details

Result: 0 - 535 - global_mem_mmmt

▼ Clear Baselines

▼ Apply Rules

▼ Occupancy Calculator

▼ Save as PDF

Report

Result

Time

Cycles

Reqs

GPU

SM Frequency

CC

Process

Current

mmul_3000

S35 - global_mem_mmmt (375, 6...

27.91 msecsd

41,855,661

40

0 - NVIDIA GeForce RTX 3070

1.50 cycle/msecsd

8.6

[17568] HW3.exe

t16x16_shared

mmul_3000

S38 - shared_mem_mmmt (188, 1...

31.21 msecsd

46,806,961

27

0 - NVIDIA GeForce RTX 3070

1.50 cycle/msecsd

8.6

[17568] HW3.exe

t16x16_global

mmul_3000

S35 - global_mem_mmmt (188, 1...

27.41 msecsd

41,108,780

40

0 - NVIDIA GeForce RTX 3070

1.50 cycle/msecsd

8.6

[17568] HW3.exe

132x32_global

132x32_shared

18x8_shared

cta32x8_global

cta64x4_global

cta8x32_global

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 [%]
Memory Throughput [%]
L1/TEX Cache Throughput [%]
L2 Cache Throughput [%]
DRAM Throughput [%]

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 [Kernel Profiling Guide](#) section.

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 64:1. The kernel achieved 6% 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.

GPU Throughput

Compute Throughput Breakdown

Memory Throughput Breakdown

Floating Point Operations Roofline

Performance [FLOP/s]
(N = 1E+12)
Arithmetic Intensity [FLOP/byte]

Compute Workload Analysis

All

Low Utilization

All compute pipelines are under-utilized. Either this kernel is very small or it doesn't issue enough warps per scheduler. Check the [Launch Statistics](#) and [Scheduler Statistics](#) sections for further details.

Pipe Utilization (% of active cycles)

Pipe Utilization (% of peak instructions executed)

Memory Workload Analysis

All

L2 Load Access Pattern

The memory access pattern for loads from L1TEX to L2 is not optimal. The granularity of an L1TEX request to L2 is a 128 byte cache line. That is 4 consecutive 32-byte sectors per L2 request. However, this kernel only accesses an average of 1.0 sectors out of the possible 4 sectors per cache line. Check the [Source Counters](#) section for uncoalesced loads and try to minimize how many cache lines need to be accessed per memory request.

L2 Store Access Pattern

The memory access pattern for stores from L1TEX to L2 is not optimal. The granularity of an L1TEX request to L2 is a 128 byte cache line. That is 4 consecutive 32-byte sectors per L2 request. However, this kernel only accesses an average of 1.0 sectors out of the possible 4 sectors per cache line. Check the [Source Counters](#) section for uncoalesced stores and try to minimize how many cache lines need to be accessed per memory request.

Memory Chart

Shared Memory

L1/TEX Cache

L2 Cache

Instruction Statistics

All

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.

Eligible Warps Per Scheduler [warp]
Active Warps Per Scheduler [warp]
Issued Warp Per Scheduler

Issue Slot Utilization

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 3.7 cycles. This might leave hardware resources underutilized and may lead to less optimal performance. One of the maximum of 12 warps per scheduler, this kernel allocates an average of 11.02 warps per scheduler, but only an average of 1.31 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 instructions 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 Statistics](#) and [Source Counters](#) sections can help, too.

Warps Per Scheduler

Warp State Statistics

Warp State (All Cycles)

Instruction Statistics

NVLink Topology

NVLink Tables

Launch Statistics

Occupancy

Source Counters

Warp Stall Sampling (All Cycles)

Most Instructions Executed