

Result

Size

Time

Cycles

GPU

SM Frequency

Process

Attributes

Current

532 - gemm\_kernel

(16, 64, 1)x(32, 8, 1)

1.03 ms

1,665,875

0 - NVIDIA GeForce RTX 4060 Laptop GPU

1.62 Ghz

[35826] gemm

Summary

Details

Source

Context

Comments

Raw

Session

Compare

Tools

View

Export

GPU Speed Of Light Throughput

Roofline Single Precision

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 [%]	63.11	Duration [ms]	1.03
Memory Throughput [%]	66.60	Elapsed Cycles [cycle]	1,665,875
L1/TEX Cache Throughput [%]	87.73	SM Active Cycles [cycle]	1,553,291.33
L2 Cache Throughput [%]	66.60	SM Frequency [Ghz]	1.62
DRAM Throughput [%]	1.20	DRAM Frequency [Ghz]	7.99

Balanced Throughput

Compute and Memory are well-balanced: To reduce runtime, both computation and memory traffic must be reduced. Check both the [Compute Workload Analysis](#) and [Memory Workload Analysis](#) sections.

Roofline Analysis

The ratio of peak float (FP32) to double (FP64) performance on this device is 64:1. The workload achieved 4% of this device's FP32 peak performance and 0% of its FP64 peak performance. See the [Profiling Guide](#) for more details on roofline analysis.

Floating Point Operations Roofline (Single Precision)

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]	1	# Pass Groups	2
Maximum Buffer Size [Mbyte]	8.39	-	-

Compute Workload Analysis

Pipe Utilization (Elapsed Cycles)

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]	0.61	SM Busy [%]	21.06
Executed ipc Active [inst/cycle]	0.66	Issue Slots Busy [%]	15.32
Issued ipc Active [inst/cycle]	0.66		

Low Utilization

Est. Local Speedup: 94.54%

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.

Key Performance Indicators

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 [Gbyte/s]	3.06	Mem Busy [%]	35.99
L1/TEX Hit Rate [%]	93.95	Max Bandwidth [%]	66.60
L2 Hit Rate [%]	99.61	Mem Pipes Busy [%]	63.11
L2 Compression Input Sectors [sector]	0	Local Memory Spilling Requests	0
L2 Compression Ratio	0	Local Memory Spilling Request Overhead [%]	0
L2 Compression Success Rate [%]	0	L2 Persisting Size [Mbyte]	6.29

L1TEX Global Load Access Pattern

Est. Speedup: 15.33%

The memory access pattern for global loads from L1TEX might not be optimal. On average, only 26.4 of the 32 bytes transmitted per sector are utilized by each thread. This could possibly be caused by a stride between threads. Check the [Source Counters](#) section for uncoalesced global loads.

Key Performance Indicators

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]	11.72	No Eligible [%]	83.71
Eligible Warps Per Scheduler [warp]	0.22	One or More Eligible [%]	16.29
Issued Warp Per Scheduler	0.16		

Issue Slot Utilization

Est. Local Speedup: 33.40%

Every scheduler is capable of issuing one instruction per cycle, but for this workload each scheduler only issues an instruction every 6.1 cycles. This might leave hardware resources underutilized and may lead to less optimal performance. Out of the maximum of 12 warps per scheduler, this workload allocates an average of 11.72 active warps per scheduler, but only an average of 0.22 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 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 Statistics](#) and [Source Counters](#) sections can help, too.

Key Performance Indicators

Warp State Statistics

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]	71.90	Avg. Active Threads Per Warp	32
Warp Cycles Per Executed Instruction [cycle]	71.92	Avg. Not Predicated Off Threads Per Warp	31.94

Long Scoreboard Stalls

Est. Speedup: 33.40%

On average, each warp of this workload spends 65.6 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 91.2% of the total average of 71.9 cycles between issuing two instructions.

Key Performance Indicators

Warp Stall

Check the [Warp Stall Sampling \(All Samples\)](#) table for the top stall locations in your source based on sampling data. The [Profiling Guide](#) provides more details on each stall reason.

Instruction Statistics

Opcode Category Chart

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

Executed Instructions [inst]	24,494,080	Avg. Executed Instructions Per Scheduler [inst]	255,146.67
Issued Instructions [inst]	24,500,913	Avg. Issued Instructions Per Scheduler [inst]	255,217.84

FP32 Non-Fused Instructions

Est. Speedup: 1.37%

This kernel executes 4194304 fused and 4202496 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 25% (relative to its current performance).

Key Performance Indicators

NVLink Topology

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

NVLink Tables

Detailed tables with properties for each NVLink.

NUMA Affinity

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

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	1,024	Function Cache Configuration	CachePreferNone
Registers Per Thread [register/thread]	26	Static Shared Memory Per Block [byte/block]	0
Block Size	256	Dynamic Shared Memory Per Block [byte/block]	0
Threads [thread]	262,144	Driver Shared Memory Per Block [Kbyte/block]	1.02
Waves Per SM	7.11	Shared Memory Configuration Size [Kbyte]	16.38
Uses Green Context	0	Stack Size	1,024
# SMs [SM]	24	# TPCs	12
Enabled TPC IDs	all	-	-

Occupancy

% Occupancy Graphs

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 [%]	98.48	Block Limit Warps [block]	6
Achieved Active Warps Per SM [warp]	47.27	Block Limit SM [block]	24

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]	1,553,291.33	Average L1 Active Cycles [cycle]	1,553,291.33
Average L2 Active Cycles [cycle]	1,555,640.38	Average SMSP Active Cycles [cycle]	1,566,327.11
Average DRAM Active Cycles [cycle]	98,500	Total SM Elapsed Cycles [cycle]	39,978,160
Total L1 Elapsed Cycles [cycle]	39,978,160	Total L2 Elapsed Cycles [cycle]	26,395,168
Total SMSP Elapsed Cycles [cycle]	159,912,640	Total DRAM Elapsed Cycles [cycle]	32,878,592

SMs Workload Imbalance

Est. Speedup: 6.20%

One or more SMs have a much higher number of active cycles than the average number of active cycles. Maximum instance value is 6.65% above the average, while the minimum instance value is 2.52% below the average.

Key Performance Indicators

SMSPs Workload Imbalance

Est. Speedup: 5.17%

One or more SMSPs have a much higher number of active cycles than the average number of active cycles. Maximum instance value is 5.50% above the average, while the minimum instance value is 3.17% below the average.

Key Performance Indicators

L1 Slices Workload Imbalance

Est. Speedup: 6.20%

One or more L1 Slices have a much higher number of active cycles than the average number of active cycles. Maximum instance value is 6.65% above the average, while the minimum instance value is 2.52% below the average.

Key Performance Indicators

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]	327,680	Branch Efficiency [%]	100
Branch Instructions Ratio [%]	0.01	Avg. Divergent Branches [branches]	0

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.