

Result

Size

Time

Cycles

GPU

SM Frequency

Process

Attributes

Current

574 - gemm_kernel_fp32

(64, 64, 1)x(32, 32, 1)

21.92 ms

35,509,003

0 - NVIDIA GeForce RTX 4060 Laptop GPU

1.62 Ghz

[36063] gemm_fp32

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 [%]	86.74	Duration [ms]	21.92
Memory Throughput [%]	86.74	Elapsed Cycles [cycle]	35,509,003
L1/TEX Cache Throughput [%]	87.04	SM Active Cycles [cycle]	35,388,244.50
L2 Cache Throughput [%]	6.06	SM Frequency [Ghz]	1.62
DRAM Throughput [%]	1.09	DRAM Frequency [Ghz]	7.99

High Throughput

This workload 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 64:1. The workload achieved 8% 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]	8	# Pass Groups	2
Maximum Buffer Size [Mbyte]	8.52	-	-

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]	0.95	SM Busy [%]	30.25
Executed Ipc Active [inst/cycle]	0.95	Issue Slots Busy [%]	23.76
Issued Ipc Active [inst/cycle]	0.95		

Low Utilization

Est. Local Speedup: 91.11%

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]	2.79	Mem Busy [%]	52.11
L1/TEX Hit Rate [%]	0.77	Max Bandwidth [%]	86.74
L2 Hit Rate [%]	97.74	Mem Pipes Busy [%]	86.74
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

Shared Store Bank Conflicts

Est. Speedup: 22.10%

The memory access pattern for shared stores might not be optimal and causes on average a 1.3 - way bank conflict across all 16777216 shared store requests. This results in 5709774 bank conflicts, which represent 25.39% of the overall 22486990 wavefronts for shared stores. Check the [Source Counters](#) section for uncoalesced shared stores.

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]	8.00	No Eligible [%]	76.15
Eligible Warps Per Scheduler [warp]	1.10	One or More Eligible [%]	23.85
Issued Warp Per Scheduler	0.24		

Issue Slot Utilization

Est. Local Speedup: 13.26%

Every scheduler is capable of issuing one instruction per cycle, but for this workload each scheduler only issues an instruction every 4.2 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 8.00 active warps per scheduler, but only an average of 1.10 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]	33.55	Avg. Active Threads Per Warp	32
Warp Cycles Per Executed Instruction [cycle]	33.55	Avg. Not Predicated Off Threads Per Warp	31.98

Mio Throttle Stalls

Est. Speedup: 13.26%

On average, each warp of this workload spends 21.2 cycles being stalled waiting for the MIO (memory input/output) instruction queue to be not full. This stall reason is high in cases of extreme utilization of the MIO pipelines, which include special math instructions, dynamic branches, as well as shared memory instructions. When caused by shared memory accesses, trying to use fewer but wider loads can reduce pipeline pressure. This stall type represents about 63.1% of the total average of 33.5 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]	810,024,960	Avg. Executed Instructions Per Scheduler [inst]	8,437,760
Issued Instructions [inst]	810,034,164	Avg. Issued Instructions Per Scheduler [inst]	8,437,855.88

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	4,096	Function Cache Configuration	CachePreferNone
Registers Per Thread [register/thread]	38	Static Shared Memory Per Block [Kbyte/block]	8.19
Block Size	1,024	Dynamic Shared Memory Per Block [byte/block]	0
Threads [thread]	4,194,304	Driver Shared Memory Per Block [Kbyte/block]	1.02
Waves Per SM	170.67	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 [%]	66.67	Block Limit Registers [block]	1
Theoretical Active Warps per SM [warp]	32	Block Limit Shared Mem [block]	1
Achieved Occupancy [%]	66.64	Block Limit Warps [block]	1
Achieved Active Warps Per SM [warp]	31.99	Block Limit SM [block]	24

Theoretical Occupancy

Est. Speedup: 13.26%

The 8.00 theoretical warps per scheduler this kernel can issue according to its occupancy are below the hardware maximum of 12. This kernel's theoretical occupancy (66.7%) is limited by the number of required registers, and the number of warps within each block.

Key Performance Indicators

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]	35,388,244.50	Average L1 Active Cycles [cycle]	35,388,244.50
Average L2 Active Cycles [cycle]	19,674,549.56	Average SMSP Active Cycles [cycle]	35,384,973.42
Average DRAM Active Cycles [cycle]	1,908,432	Total SM Elapsed Cycles [cycle]	852,214,496
Total L1 Elapsed Cycles [cycle]	852,214,496	Total L2 Elapsed Cycles [cycle]	562,798,384
Total SMSP Elapsed Cycles [cycle]	3,408,857,984	Total DRAM Elapsed Cycles [cycle]	700,798,976

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]	8,781,824	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.