

Report

Current

Baseline 1

gemm_fp32_warp_spl_8192_profiling

gemm_fp32_shared_mem_8192_profiling

Res

574 - gemm_warp_specialized_16x16

Size

(256, 256, 16x16, 16, 1)

Time

869.70 ms

Cycle

0 - NVIDIA GeForce RTX 5060 Ti

GPU

0 - NVIDIA GeForce RTX 5060 Ti

SM Frequency

2.41 GHz

Process

[825162] gemm_fp32_warp_spl_8192.exe

Attributes

[825059] gemm_fp32_shared_mem_8192.exe

Summary

Details

Source

Context

Comment

Raw

Session

GPU Speed of Light Throughput

GPU Throughput Chart

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 profile chart.

Compute (SM) Throughput [%]	99.81 (+29.25%)	Duration [ms]	869.70 (+21.08%)
Memory Throughput [%]	99.81 (+29.25%)	Elapsed Cycles [cycle]	2,09,20,01,743.14 (+23.08%)
L1/TEX Cache Throughput [%]	99.83 (+29.19%)	SM Active Cycles [cycle]	2,09,20,01,743.14 (+23.14%)
L2 Cache Throughput [%]	7.18 (-18.62%)	SM Frequency [GHz]	2.41 (+0.00%)
DRAM Throughput [%]	18.17 (-18.63%)	DRAM Frequency [GHz]	13.79 (+0.00%)

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 [Launch Statistics](#) section.

Roofline Analysis

The ratio of peak float (FP32) to double (FP64) performance on this device is 64:1. The workload achieved 9% of this device's FP32 peak performance and 0% of its FP64 peak performance. See the [Profiling Guide](#) for more details on roofline analysis.

GPU Throughput

Compute (SM) [%]

Memory [%]

Speed Of Light (SOL) [%]

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 [µs] 192 (+0.00%) # Pass Groups 1 (+0.00%)

Maximum Buffer Size [Mbyte] 53.28 (+23.00%)

Executed IPC Active

Blocks Launched

Blocks Active

Warps Launched

Warps Active

CGAs Launched

CGAs Active

SM

SM Throughput

SM ALU Heavy

SM ALU Size 64B

SM FMA

SM FMA Heavy

SM Tensor

SM Tensor HMMMA

SM Tensor MMA

SM Uniform

SM XU

SM Bytes Shared

SM DCC Hit Miss

L1

L1 LSU Writeback Throughput

L1 TEX Writeback Throughput

L1 Hit Miss

L1 Lookup Hit %

L1 Lookup Miss %

L1 Lookup Miss

SMEM Bank Conflicts

L1 Wavefronts (Data)

L1 Wavefronts %

L1 Wavefronts

L2

L2 Throughput

L2 Sectors %

L2 Sectors

L2 to XBAR Active

XBAR to L2 Active

L2 Hit Miss

L2 Hit Rate CE

L2 Hit Rate GDC

L2 Hit Rate GPC

L2 Hit Rate HUB

L2 Hit Rate TEX Atom

L2 Hit Rate TEX Read

L2 Hit Rate TEX Write

SysL2

SysL2 Throughput

SysL2 Sectors %

SysL2 Sectors

SysL2 Atomic Input Active

SysL2 Hit Miss

SysL2 to XBAR Active

XBAR to SysL2 Active

PCIe

PCIe Throughput

PCIe Bytes Read

PCIe Bytes Write

DRAM

DRAM Bandwidth

DRAM Bytes

DRAM Sectors Read

DRAM Sectors Write

Workload Execution

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]	1.27 (+16.50%)	SM Busy [%]	35.47 (+28.93%)
Executed IPC Active [Inst/cycle]	1.27 (+16.44%)	Issue Slots Busy [%]	31.72 (+16.50%)
Issued IPC Active [Inst/cycle]	1.27 (+16.44%)		

Low Utilization

Est. Local Speedup: 86.20%

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

Pipe Utilization (% of elapsed cycles)

Pipe Utilization (% of peak instructions executed over elapsed cycles)

Utilization [%]

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]	80.17 (+18.63%)	Mem Busy [%]	50.04 (+8.45%)
L1/TEX Hit Rate [%]	0.00 (+inf%)	Max Bandwidth [%]	99.81 (+29.25%)
L2 Hit Rate [%]	49.48 (+0.26%)	Mem Pipes Busy [%]	99.81 (+29.25%)
L2 Compression Input Sectors [sector]	0 (+0.00%)	Local Memory Spilling Requests	0 (+0.00%)
L2 Compression Ratio	0 (+0.00%)	Local Memory Spilling Request Overhead [%]	0 (+0.00%)
L2 Compression Success Rate [%]	0 (+0.00%)	L2 Persisting Size [Mbyte]	6.29 (+0.00%)

Memory Chart

Values: Transfer Size Inactivity Greyed Out

Kernel

Global

Local

Texture

Surface

Load Global Store Shared

TMA

Shared

DSMEM

ICC

DCC

IBC

GCC

System L2 Cache

Device L2 Cache

Distributed Shared Memory

De-compressor Rate

Compressor Rate

System Memory

Peer Memory

Device Memory

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]	7.99 (+0.00%)	No Eligible [%]	68.27 (+6.16%)
Eligible Warps Per Scheduler [warp]	0.82 (+16.40%)	One or More Eligible [%]	31.84 (+0.50%)

Issue Slot Utilization

Est. Local Speedup: 0.11%

Every scheduler is capable of issuing one instruction per cycle but for this workload each scheduler only issues an instruction every 2.2 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 7.99 active warps per scheduler, but only an average of 0.82 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](#) section can help, too.

Key Performance Indicators

Warps Per Scheduler

GPU Maximum Warps Per Scheduler

Theoretical Warps Per Scheduler

Active Warps Per Scheduler

Eligible Warps Per Scheduler

Issued Warp Per Scheduler

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 predicted 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]	25.19 (+14.20%)	Avg. Active Threads Per Warp	32 (+0.00%)
Warp Cycles Per Executed Instruction [cycle]	25.19 (+14.20%)	Avg. Not Predicted Off Threads Per Warp	31.84 (+0.50%)

Mio Throttle Stalls

Est. Speedup: 0.11%

On average, each warp of this workload spends 10.1 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 40.1% of the total average of 25.2 cycles between issuing two instructions.

Key Performance Indicators

Barrier Stalls

On average, each warp of this workload spends 8.7 cycles being stalled waiting for sibling warps at a CTA barrier. A high number of warps waiting at a barrier is commonly caused by diverging code paths before a barrier. This causes some warps to wait a long time until other warps reach the synchronization point. Whenever possible, try to divide the work into blocks of uniform workloads. If the block size is 512 threads or greater, consider splitting it into smaller groups. This can increase eligible warps without affecting occupancy, unless shared memory becomes a new occupancy limiter. Also, try to identify which barrier instruction causes the most stalls, and optimize the code executed before that synchronization point first. This stall type represents about 34.6% of the total average of 25.2 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.

Warp State (All Cycles)

Stall MIO Throttle

Stall Barrier

Stall Not Selected

Stall Wait

Stall Long Scoreboard

Selected

Stall Short Scoreboard

Stall No Instruction

Stall Dispatch Stall

Stall Branch Resolving

Stall Math Pipe Throttle

Stall Drain

Stall LG Throttle

Stall Misc

Stall Member

Stall Sleeping

Stall Tex Throttle

Cycles per Instruction

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]	95,58,118,96,704 (+43.38%)	Avg. Executed Instructions Per Scheduler [Inst]	66,37,63,171.56 (+43.38%)
Issued Instructions [Inst]	95,58,118,96,704 (+43.38%)	Avg. Issued Instructions Per Scheduler [Inst]	66,37,63,171.56 (+43.38%)

FP32 Non-Fused Instructions

Est. Speedup: 3.45%

This kernel executes 17179869184 fused and 17181966336 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

Executed Instruction Categories

Load/Store

Floating Point

Integer

Control

Movement

Uniform Datapath

Miscellaneous

Conversion

Executed Warp-Level Instructions/Opcode

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 Configuration

Summary of the configurations 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	65,536 (+0.00%)	Function Cache Configuration	CachePreferNone (CachePreferNone)
Cluster Size	0 (+0.00%)	Preferred Cluster Size	0 (+0.00%)
Registers Per Thread [register/thread]	40 (+5.26%)	Cluster Scheduling Policy	PolicySpread (PolicySpread)
Static Shared Memory Per Block [kbyte/block]	20,99 (+156.25%)	Block Size	256 (-75.00%)
Dynamic Shared Memory Per Block [byte/block]	0 (+0.00%)	Threads [thread]	1,67,77,216 (-75.00%)
Driver Shared Memory Per Block [kbyte/block]	1,02 (+0.00%)	Waves Per SM	455.11 (-75.00%)
Shared Memory Configuration Size [kbyte]	102,46 (+508.00%)	Users Green Context	24 (+0.00%)
Stack Size	1,024 (+0.00%)	SMs [SM]	392,06,75,440 (+12.72%)
# TPCs	18 (+0.00%)	Enabled TPC IDs	8 (+0.00%)

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 actually in use. Higher occupancy does not always result in higher performance. Lower 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 indicate highly imbalanced workloads.

Theoretical Occupancy [%]	66.67 (+0.00%)	Block Limit Shared Mem [block]	6 (+500.00%)
Theoretical Active Warps per SM [warp]	32 (+0.00%)	Block Limit Warps [block]	6 (+500.00%)
Achieved Occupancy [%]	66.60 (+0.10%)	Block Limit SM [block]	24 (+0.00%)
Achieved Active Warps Per SM [warp]	31.97 (+0.10%)	Block Limit Barriers [block]	8 (+0.00%)
Cluster Occupancy [%]	0 (+0.00%)	Max Cluster Size [block]	8 (+0.00%)
Max Active Clusters [cluster]	0 (+0.00%)		
Overall GPU Occupancy [%]	0 (+0.00%)		

Key Performance Indicators

GPU and Memory Workload Distribution

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

Average SM Active Cycles [cycle]	2,09,20,01,743.14 (+23.14%)	Average L1 Active Cycles [cycle]	2,09,20,01,743.14 (+23.14%)
Average L2 Active Cycles [cycle]	1,88,84,97,483.94 (+38.84%)	Average SMP Active Cycles [cycle]	2,09,20,00,198.78 (+23.15%)
Average DRAM Active Cycles [cycle]	2,187,78,248 (+0.15%)	Total SM Elapsed Cycles [cycle]	75,33,06,60,660 (+23.08%)
Total L1 Elapsed Cycles [cycle]	75,33,06,60,660 (+23.08%)	Total L2 Elapsed Cycles [cycle]	30,24,06,75,440 (+12.72%)
Total SMP Elapsed Cycles [cycle]	3,01,32,26,42,640 (+23.08%)	Total DRAM Elapsed Cycles [cycle]	47,36,61,76,256 (+23.08%)

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]	195,87,39,968 (+260.62%)	Branch Efficiency [%]	100 (+0.00%)
Branch Instructions Ratio [%]	0.02 (+151.51%)	Avg. Divergent Branches [branches]	0 (+0.00%)

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.