

Report

Result

Size

Time

Cycles

GPU

SM Frequency

Process

Attributes

Baseline 1

gemm_ws_1024_2560_profiling

gemm_am_1024_profiling

574 - gemm_kernel

574 - gemm_kernel

(64, 64, 1)x(16, 16, 1)

1.93 ms

46,54,521

0 - NVIDIA GeForce RTX 5090 Ti

2.40 GHz

[110593] gemm_ws_1024_2560

Summary

Details

Source

Context

Comments

Raw

Session

Compare

Tools

View

Export

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 pipeline chart.

Compute (SM) Throughput [%]	95.45	(+0.47%)	Duration [ms]	1.93	(+59.01%)
Memory Throughput [%]	95.45	(+0.47%)	Elapsed Cycles [cycle]	46,54,521	(+59.01%)
L1/TEX Cache Throughput [%]	95.71	(+0.70%)	SM Active Cycles [cycle]	46,37,211.75	(+59.43%)
L2 Cache Throughput [%]	12.66	(+4.17%)	SM Frequency [GHz]	2.40	(+0.03%)
DRAM Throughput [%]	1.48	(+37.04%)	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 [Compute Workload Analysis](#) section.

Routine 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 routine analysis.

GPU Throughput

Compute (SM) [%]	95.45	(+0.47%)
Memory [%]	95.45	(+0.47%)

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]

3

(+0.00%)

Pass Groups

1

(+0.00%)

Maximum Buffer Size [Mbyte]

16.52

(+0.00%)

Executed IPC Active

192 block

1.57M block

1.73k warp

1.56M warp

CGAs Launched

0

CGAs Active

0

SM

SM Throughput

54.7 %

SM ALU Heavy

100 %

SM ALU Size 64B

100 %

SM FMA

100 %

SM FMA Heavy

100 %

SM Tensor

100 %

SM Tensor HMMMA

100 %

SM Tensor IMMA

100 %

SM Uniform

100 %

SM XU

100 %

SM Bytes Shared

4.767 Tbyte/s

SM DCC Hit Miss

28.1k cycle

L1

L1 LSU Writeback Throughput

100 %

L1 TEX Writeback Throughput

100 %

L1 Hit Miss

0

L1 Lookup Hit %

100 %

L1 Lookup Miss %

100 %

L1 Lookup Hit Miss

30k sector

SMEM Bank Conflicts

7.99k

L1 Wavefronts (Data)

100 %

L1 Wavefronts %

271k

L2

L2 Throughput

15.2 %

L2 Sectors %

100 %

L2 Sectors

59.5k sector

L2 to XBAR Active

100 %

XBAR to L2 Active

100 %

L2 Hit Miss

0

L2 Hit Miss

30k sector

L2 Hit Rate CE

0 sector

L2 Hit Rate GCC

100 sector

L2 Hit Rate GPC

100 sector

L2 Hit Rate HUB

100 sector

L2 Hit Rate TEX Atom

100 sector

L2 Hit Rate TEX Read

100 sector

L2 Hit Rate TEX Write

100 sector

SysL2

SysL2 Throughput

3.39 %

SysL2 Sectors %

100 %

SysL2 Sectors

416 sector

SysL2 Atomic Input Active

100 %

SysL2 Hit Miss

416 sector

SysL2 to XBAR Active

100 %

XBAR to SysL2 Active

100 %

PCIE

PCIE Throughput

100 %

PCIE Bytes Read

1M Mbyte/s

PCIE Bytes Write

30 Gbyte/s

DRAM

DRAM Bandwidth

100 %

DRAM Bytes

1090 Gbyte/s

DRAM Sectors Read

10.2k sector

DRAM Sectors Write

0 sector

Workload Execution

gemm_warp_specialized_16x16

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]	1.95	(+36.40%)	SM Busy [%]	48.74	(+35.03%)
Executed IPC Active [Inst/cycle]	1.96	(+36.08%)	Issue Slots Busy [%]	48.74	(+36.40%)

Low Utilization

Est. Local Speedup: 83.78%

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)

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 Bus), exhausting the available communication bandwidth between those units (Max Bandwidth), or by reaching the maximum number of active warps (Max Warps). Detailed chart of the memory units. Detailed tables with data for each memory unit.

Memory Throughput [Gbyte/s]	6.51	(37.04%)	Mem Busy [%]	53.23	(12.70%)
L1/TEX Hit Rate [%]	1.25	(77.25%)	Max Bandwidth [%]	95.45	(+0.47%)
L2 Hit Rate [%]	97.98	(+0.44%)	Mem Pipes Busy [%]	95.45	(+0.47%)
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%)

Shared Store Bank Conflicts

Est. Speedup: 39.23%

The memory access pattern for shared stores might not be optimal and causes on average a 1.7 - way bank conflict across all 6324224 shared store requests. This results in 4392063 bank conflicts, which represent 40.99% of the overall 10714859 warfronts for shared stores. Check the [Source Counters](#) section for unallocated shared stores.

Key Performance Indicators

Memory Chart

Values: Transfer Size

Inactivity: Greyed Out

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 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 failed issue slots indicates poor latency hiding.

Active Warps Per Scheduler [warp]	11.75	(+0.04%)	One or More Eligible [%]	51.12	(20.23%)
Eligible Warps Per Scheduler [warp]	1.11	(21.59%)	Any Issued Instructions Per Scheduler [Inst]	48.88	(+36.08%)

Issue Slot Utilization

Est. Speedup: 4.55%

Every scheduler is capable of issuing one instruction per cycle, but for this workload each scheduler only issues an instruction every 2.0 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 11.75 active warps per scheduler, but only an average of 1.11 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

Warps 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 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.

Barrier Stalls	24.04	(26.49%)	Avg. Active Threads Per Warp	32	(+0.00%)
Barrier Stalls	24.04	(26.49%)	Avg. Not Predicted Off Threads Per Warp	31.12	(2.72%)

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

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]	32,63,65,280	(+116.98%)	Avg. Executed Instructions Per Scheduler [Inst]	22,66,463.93	(+116.98%)
Issued Instructions [Inst]	32,63,65,280	(+116.98%)	Avg. Issued Instructions Per Scheduler [Inst]	22,66,463.93	(+116.98%)

FP32 Non-Fused Instructions

This kernel executes 33554432 fused and 33587200 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). Executed Speedup: 4.06%

Key Performance Indicators

Executed Instruction Categories

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 GPU resources needed to execute the kernel. Choosing an efficient launch configuration maximizes device utilization.

Grid Size	4,096	(+0.00%)	Function Cache Configuration	CachePreference	(CachePreference)
Cluster Size	0	(+0.00%)	Preferred Cluster Size	0	(+0.00%)
Registers Per Thread [Register/thread]	40	(+11.11%)	Cluster Scheduling Policy	PolicySpread	(PolicySpread)
Static Shared Memory Per Block [Kbyte/block]	5.38	(+162.50%)	Block Size	256	(+0.00%)
Dynamic Shared Memory Per Block [Kbyte/block]	0	(+0.00%)	Threads [thread]	10,48,576	(+0.00%)
Driver Shared Memory Per Block [Kbyte/block]	1.02	(+0.00%)	Waves Per SM	18.96	(+0.00%)
Shared Memory Configuration Size [Kbyte]	102.40	(+56.25%)	Uses Green Context	0	(+0.00%)
Block Size	1,024	(+0.00%)	# SMs [SM]	510,990	(+0.00%)
# TPCs	18	(+0.00%)	Enabled TPCs	36	(+0.00%)

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 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	(+0.00%)	Block Limit Registers [block]	6	(+0.00%)
Theoretical Active Warps Per SM [warp]	48	(+0.00%)	Block Limit Shared Mem [block]	16	(23.81%)
Average L2 Active Cycles [cycle]	91.88	(+68.93%)	Block Limit Warps [block]	5	(10.00%)
Achieved Active Warps Per SM [warp]	47.01	(+0.00%)	Block Limit SM [block]	24	(+0.00%)
Cluster Occupancy [%]	0	(+0.00%)	Block Limit Barriers [block]	24	(+0.00%)
Max Active Clusters [cluster]	0	(+0.00%)	Max Cluster Size [block]	8	(+0.00%)
Overall GPU Occupancy [%]	0	(+0.00%)			

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]	46,37,211.75	(+59.43%)	Average L1 Active Cycles [cycle]	46,37,211.75	(+59.43%)
Average L2 Active Cycles [cycle]	16,73,96	(+10.12%)	Average SMP Active Cycles [cycle]	20,97,159	(+59.43%)
Average DRAM Active Cycles [cycle]	3,93,816	(+10.12%)	Total SM Elapsed Cycles [cycle]	16,73,96,208	(+59.05%)
Total L1 Elapsed Cycles [cycle]	16,73,96,208	(+59.05%)	Total L2 Elapsed Cycles [cycle]	6,65,91,104	(+59.05%)
Total SMSP Elapsed Cycles [cycle]	66,95,84,832	(+59.05%)	Total DRAM Elapsed Cycles [cycle]	10,67,04,996	(+59.01%)

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 each stall reason. Only focus on stalls if the schedulers fail to issue every cycle.

Branch Instructions [Inst]	2,97,12,384	(+1,283.38%)	Branch Efficiency [%]	100	(+0.00%)
Branch Instructions Ratio [%]	0.09	(+523.82%)	Warp Divergence Branches [branch]	0	(+0.00%)

Unallocated Shared Accesses

This kernel has unallocated shared accesses resulting in a total of 4194304 excessive warfronts (5% of the total 79757312 warfronts). Check the L1 Wavefronts Shared Excessive table for the primary source locations. The [CUDA Best Practices Guide](#) has an example on optimizing shared memory accesses.

Key Performance Indicators

L1 Wavefronts Shared Excessive