

Report

Current

Result

Size

Time

Cycles

GPU

SM Frequency

Process

Attributes

Baseline 1

jacob2D_pointer_swap_1024_profiling

572 - jacob1_kernel

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

54.50 us

81,671

0 - NVIDIA GeForce RTX 5060 Ti

2.35 GHz

[354543] jacob2D_warp_spl_1024.exe

[435743] jacob2D_pointer_swap_1024.exe

Summary

Details

Source

Context

Comments

Raw

Session

Compare

Tools

View

Export

Menu

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 contributors. 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 [%]

35.41 (+18.92%)

Duration [us]

34.50 (+136.40%)

27.66 (+57.68%)

Elapsed Cycles [cycle]

81,671 (+143.55%)

25.24 (+77.71%)

SM Active Cycles [cycle]

77,541.89 (+157.36%)

23.16 (+49.51%)

SM Frequency [GHz]

2.35 (+3.19%)

27.66 (+57.68%)

DRAM Frequency [GHz]

13.77 (+0.14%)

Latency Issue

This workload exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of this device. Achieved compute throughput and/or memory bandwidth below 60.0% of peak typically indicate latency issues. Look at the Scheduler Statistics and Warp State sections for potential reasons.

Key Performance Indicators

Roofline Analysis

The ratio of peak float (FP32) to double (FP64) performance on this device is 64:1. The workload achieved close to 1% 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 [%]

0.0

10.0

20.0

30.0

40.0

50.0

60.0

70.0

80.0

90.0

100.0

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 [us]

3 (+0.00%)

Pass Groups

1 (+0.00%)

Maximum Buffer Size [MiByte]

8.26 (+0.00%)

-

Executed IPC Active

1.54 inst/cycle

0

Blocks Launched

416 block

0

Blocks Active

1.48M block

0

Warps Launched

3.46k warp

0

Warps Active

1.33M warp

0

CGAs Launched

0

CGAs Active

0

SM

41.1 %

0

SM Throughput

0

SM ALU Heavy

100 %

0

SM ALU Size 64B

100 %

0

SM FMA

100 %

0

SM FMA Heavy

100 %

0

SM Tensor

0

SM Tensor HMMA

100 %

0

SM Tensor IMMA

100 %

0

SM Uniform

100 %

0

SM XU

0

SM Bytes Shared

730G Gbyte/s

0

SM DCC Hit Miss

17.8k cycle

0

L1

L1 LSU Writeback Throughput

100 %

0

L1 TEX Writeback Throughput

100 %

0

L1 Hit Miss

0

L1 Lookup Hit %

100 %

0

L1 Lookup Miss %

100 %

0

L1 Lookup Hit Miss

44.6k sector

0

SMEM Bank Conflicts

19.4k

0

L1 Wavefronts (Data)

0

L1 Wavefronts %

100 %

0

L1 Wavefronts

114k

0

L2

L2 Throughput

26.5 %

0

L2 Sectors %

100 %

0

L2 Sectors

81.5k sector

0

L2 to XBAR Active

100 %

0

XBAR to L2 Active

100 %

0

L2 Hit Miss

0

L2 Hit Miss

41.1k sector

0

L2 Hit Rate CE

0 sector

L2 Hit Rate GCC

100 sector

0

L2 Hit Rate GPC

100 sector

0

L2 Hit Rate HUB

100 sector

0

L2 Hit Rate TEX Atom

0 sector

L2 Hit Rate TEX Read

57.2 sector

0

L2 Hit Rate TEX Write

56.1 sector

0

SysL2

SysL2 Throughput

3.37 %

0

SysL2 Sectors %

100 %

0

SysL2 Sectors

416 sector

0

SysL2 Atomic Input Active

100 %

0

SysL2 Hit Miss

416 sector

0

SysL2 to XBAR Active

100 %

0

XBAR to SysL2 Active

100 %

0

PCIE

PCIE Throughput

100 %

0

PCIE Bytes Read

851M Mbyte/s

0

PCIE Bytes Write

4.3G Gbyte/s

0

DRAM

DRAM Bandwidth

100 %

0

DRAM Bytes

137G Gbyte/s

0

DRAM Sectors Read

12.9k sector

0

DRAM Sectors Write

0 sector

Workload Execution

0

jacob1_warp_specialized

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]

1.42 (+26.53%)

SM Busy [%]

35.41 (+26.53%)

Executed IPC Active [inst/cycle]

1.48 (+19.94%)

Issue Slots Busy [%]

35.41 (+26.53%)

Issued IPC Active [inst/cycle]

1.48 (+19.94%)

Low Utilization

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.

Est. Local Speedup: 84.86%

Key Performance Indicators

Pipe Utilization (% of elapsed cycles)

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

ALU

FMA

FP64

TMA

Tensor (All)

Tensor (FP)

Tensor (NT)

LSU

ADU

ALU

FMA

XU

CBU

FMA (FP16)

Uniform

FP64

FP64 (DMMA)

FP64 (FP64)

TEX

TMA

Tensor (FP)

Tensor (NT)

WorkID/CLC

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 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]

121.91 (+57.62%)

Mem Busy [%]

21.97 (+1.17%)

L1/TEX Hit Rate [%]

8.98 (+84.50%)

Max Bandwidth [%]

27.66 (+57.68%)

L1/TEX Hit Rate [%]

53.25 (+99.26%)

Mem Pipes Busy [%]

24.09 (+44.84%)

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 [MiByte]

6.29 (+0.00%)

L1TEX Global Load Access Pattern

Est. Speedup: 11.96%

The memory access pattern for global loads from L1TEX might not be optimal. On average, only 16.1 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

L1TEX Global Store Access Pattern

Est. Speedup: 0.05%

The memory access pattern for global stores to L1TEX might not be optimal. On average, only 31.9 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 stores.

Key Performance Indicators

Shared Load Bank Conflicts

Est. Speedup: 13.13%

The memory access pattern for shared loads might not be optimal and causes on average a 2.1 - way bank conflict across all 163840 shared load requests. This results in 177186 bank conflicts, which represent 52.01% of the overall 340651 wavefronts for shared loads. Check the Source Counters section for uncoalesced shared loads.

Key Performance Indicators

Memory Chart

Values: Transfer Size Inactivity: Greyed Out

Kernel

Global

Local

Texture

Surface

Load Global Store Shared

TMA

Shared

DSMEM

ICC

DCC

IDC

System L2 Cache

Device L2 Cache

System 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 policy latency hiding.

Active Warps Per Scheduler [warp]

10.09 (+17.12%)

No Eligible [%]

62.61 (+9.37%)

Eligible Warps Per Scheduler [warp]

41.42,204 (+208.68%)

Avg. Executed Instructions Per Scheduler [inst]

28,765.31 (+208.68%)

Issued Warp Per Scheduler

0.58 (+20.77%)

One or More Eligible [%]

37.39 (+20.94%)

Issue Slot Utilization

Est. Local Speedup: 62.61%

Every scheduler is capable of issuing one instruction per cycle, but for this workload each scheduler only issues an instruction every 2.7 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 11.08 active warps per scheduler, but only an average of 0.58 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]

26.98 (+3.16%)

Avg. Active Threads Per Warp

31.33 (+1.97%)

Warp Cycles Per Executed Instruction [cycle]

26.98 (+3.16%)

Avg. Not Predicted Off Threads Per Warp

29.18 (+6.42%)

Barrier Stalls

Est. Speedup: 46.42%

On average, each warp of this workload spends 12.5 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 up 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 46.4% of the total average of 27.0 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 Barrier

Stall Long Scoreboard

Stall Wait

Stall Short Scoreboard

Selected

Stall Not Selected

Stall Branch Resolving

Stall No Instruction

Stall Dispatch Stal

Stall Math Pipe Throttle

Stall MIO Throttle

Stall Drain

Stall Misc

Stall LG Throttle

Stall Member

Stall Sleeping

Stall Tex Throttle

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]

41,42,204 (+208.68%)

Avg. Executed Instructions Per Scheduler [inst]

28,765.31 (+208.68%)

Issued Instructions [inst]

41,42,204 (+208.68%)

Avg. Issued Instructions Per Scheduler [inst]

28,765.31 (+208.68%)

FP32 Non-Fused Instructions

Est. Speedup: 3.96%

This kernel executes 0 fused and 163840 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 50% (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

4096 (+0.00%)

Preferred Cluster Configuration

Cache/Prefer/None (Cache/Prefer/None)

Cluster Size

0 (+0.00%)

Function Cluster Size

0 (+0.00%)

Cluster Scheduling Policy

Policy/Spread (Policy/Spread)

Static Shared Memory Per Thread [register/thread]

35 (+75.00%)

Block Size

1.30 (+44%)

Threads [thread]

255 (+20.00%)

Dynamic Shared Memory Per Block [Kbyte/block]

9 (+0.00%)

Block Limit Shared Mem [block]

10,485,576 (+143.44%)

Driver Shared Memory Per Block [Kbyte/block]

1.02 (+0.00%)

Waves Per SM

18.96 (+0.00%)

Shared Memory Configuration Size [Kbyte]

32.77 (+100.00%)

Uses Green Context

0 (+0.00%)

Stack Size

1,024 (+0.00%)

SMs [SM]

36 (+0.00%)

TPCs

18 (+0.00%)

Enabled TPC IDs

all (all)

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 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 Active Warps per SM [warp]

108 (+0.00%)

Block Limit Registers [block]

6 (+40.00%)

Achieved Occupancy [%]

83.39 (+13.31%)

Block Limit Warps [block]

13 (18.75%)

Cluster Occupancy [%]

40.03 (+13.31%)

Block Limit SM [block]

6 (+0.00%)

Max Active Clusters [cluster]

0 (+0.00%)

Block Limit Barriers [block]

24 (+0.00%)

Overall GPU Occupancy [%]

0 (+0.00%)

Max Cluster Size [block]

8 (+0.00%)

Uncoalesced Global Accesses

Est. Speedup: 16.61%

The difference between calculated theoretical (100.0%) and measured actual occupancy (83.4%) can be the result of warp scheduling overheads or workload imbalances during the kernel execution. Load imbalances can occur between warps within a block as well as across blocks of the same kernel. See the CUDA Best Practices Guide for more details on optimizing occupancy.

Key Performance Indicators

GPU and Memory Workload Distribution

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

Average SM Active Cycles [cycle]

77,541.89 (+157.36%)

Average L1 Active Cycles [cycle]

77,541.89 (+157.36%)

Average L2 Active Cycles [cycle]

69,192.88 (+161.72%)

Average SMP Active Cycles [cycle]

76,924.76 (+155.25%)

Average DRAM Active Cycles [cycle]

1,31,424 (+0.1%)

Total SM Elapsed Cycles [cycle]

29,24,448 (+143.96%)

Total L1 Elapsed Cycles [cycle]

29,240 (+0.00%)

Total L2 Elapsed Cycles [cycle]

11,65,776 (+143.44%)

Total SMP Elapsed Cycles [cycle]

1,16,57,792 (+143.96%)

Total DRAM Elapsed Cycles [cycle]

19,00,544 (+156.73%)

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]

3,61,596 (+452.29%)

Branch Efficiency [%]

97.26 (+66%)

Branch Instructions Ratio [%]

0.09 (+78.92%)

Avg. Divergent Branches [branches]

36.42 (+44%)

Uncoalesced Global Accesses

Est. Speedup: 33.30%

This kernel has uncoalesced global accesses resulting in a total of 161220 excessive sectors (35% of the total 459716 sectors). Check the L2 Theoretical Sectors Global Excessive table for the primary source locations. The CUDA Programming Guide has additional information on reducing uncoalesced device memory accesses.

Key Performance Indicators

L2 Theoretical Sectors Global Excessive

Location

Value

Value (%)

0x748167757860 in jacob1_warp_specialized

1,61,220

100

0x748167757860 in jacob1_warp_specialized

0

0

0x748167757860 in jacob1_warp_specialized

0

0

0x748167757860 in jacob1_warp_specialized

0

0

0x748167757860 in jacob1_warp_specialized

0

0

Uncoalesced Shared Accesses

Est. Speedup: 41.87%

This kernel has uncoalesced shared accesses resulting in a total of 163200 excessive wavefronts (44% of the total 372096 wavefronts). 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

Location

Value

Value (%)

0x748167757860 in jacob1_warp_specialized

32,640

20

0x748167757860 in jacob1_warp_specialized

32,640

20

0x748167757860 in jacob1_warp_specialized

32,640

20

0x748167757860 in jacob1_warp_specialized

32,640

20

0x748167757860 in jacob1_warp_specialized

32,640

20

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.