

Report

Baseline 1

Result

jacob2D_shared_mem_1024_profiling

Size

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

Time

27.10 us

Cycles

63,792

GPU

0 - NVIDIA GeForce RTX 5090 Ti

SM Frequency

2.34 GHz

Process

[345743] jacob2D_shared_mem_swap_1024.exe

Attributes

Summary

Details

Source

Context

Comments

Raw

Session

Compare

Tools

View

Export

Fullscreen

GPU Speed of Light Throughput

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 bar chart.

Compute (SM) Throughput [%]

50.76 (+16.22%)

Duration [ns]

27.10 (+85.75%)

Memory Throughput [%]

50.76 (+22.48%)

Elapsed Cycles [cycle]

63,792 (+69.22%)

L1/TEX Cache Throughput [%]

53.54 (+10.93%)

SM Active Cycles [cycle]

60,114.81 (+99.52%)

L2 Cache Throughput [%]

36.15 (+21.17%)

SM Frequency [GHz]

2.34 (+2.52%)

DRAM Throughput [%]

35.15 (+46.22%)

DRAM Frequency [GHz]

13.77 (+0.12%)

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 [Launch Parameters](#) and [Scheduler Statistics](#) for potential reasons.

Key Performance Indicators

GPU Throughput Chart

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 [ms]

3 (+0.00%)

Pass Groups

1 (+0.00%)

Maximum Buffer Size [Mbyte]

4.13 (+50.00%)

PM Sampling Data

Sampling interval is larger than 10% of the workload duration, which likely results in very few collected samples. For better results, use the -pm-sampling-interval option to reduce the sampling interval. Use -pm-sampling-buffer-size to increase the sampling buffer size for the smaller interval, or don't set a fixed buffer size and let the tool adjust it automatically.

Compute Workload Analysis

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 Active [Inst/cycle]

1.11 (+0.89%)

SM Busy [%]

27.74 (+0.89%)

Executed lpc Active [Inst/cycle]

1.17 (+5.40%)

Issue Slots Busy [%]

27.74 (+0.89%)

Issued lpc Active [Inst/cycle]

1.17 (+5.40%)

Low Utilization

Est. Local Speedup: 94.80%

All compute pipelines are under-utilized. Either this workload is very small or it doesn't issue enough warps per scheduler. Check the [Launch Parameters](#) and [Scheduler Statistics](#) sections for further details.

Memory Workload Analysis

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 uncoalesced global stores.

Memory Throughput [Gbyte/s]

154.90 (+46.15%)

Mem Busy [%]

41.15 (+10.17%)

L1/TEX Hit Rate [%]

0.40 (+99.31%)

Max Bandwidth [%]

50.76 (+22.34%)

Cache Rate [%]

51.79 (+68.03%)

Perfmem Cluster Rate

0 (+0.00%)

Perfmem Spilling Requests

0 (+0.00%)

L2 Compression Input Sectors [sector]

0 (+0.00%)

Local Memory Spilling Request Overhead [%]

0 (+0.00%)

L2 Compression Rate [%]

0 (+0.00%)

L2 Persting Size [Mbyte]

6.29 (+0.00%)

L2 Compression Success Rate [%]

0 (+0.00%)

The memory access pattern for global loads from L1/TEX might not be optimal. On average, only 18.9 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 loads.

Scheduler Statistics

Summary of the statistics of the scheduler 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]

10.54 (+22.41%)

No Eligible [%]

70.59 (+2.18%)

Eligible Warps Per Scheduler [warp]

0.45 (+4.00%)

One or More Eligible [%]

29.41 (+4.67%)

Issued Warp Per Scheduler

0.29 (+4.87%)

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 spend 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]

35.85 (+28.64%)

Avg. Active Threads Per Warp

26.75 (+16.30%)

Warp Cycles Per Executed Instruction [cycle]

25,32,736 (+88.74%)

Avg. Non-Prefetched Off-Threads Per Warp

25.62 (+17.84%)

Long Scoreboard Stalls

On average, each warp of this workload spends 19.2 cycles being stalled waiting for a scoreboard dependency on a L1/TEX (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 L1/TEX 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 33.5% of the total average of 55.8 cycles between issuing two instructions.

Est. Local Speedup: 49.24%

Warp Stall

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

Warp State (All Cycles)

Stall Long Scoreboard

Stall Barrier

Stall Wait

Stall Short Scoreboard

Stall MID Throttle

Selected

Stall Branch Resolving

Stall Drain

Stall Not Issued

Stall No Instruction

Stall Dispatch Stall

Stall Math Pipe Throttle

Stall LG Throttle

Stall Member

Stall Misc

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 instructions typically 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]

25,32,736 (+88.74%)

Avg. Executed Instructions Per Scheduler [Inst]

17,588.44 (+88.74%)

Issued Instructions [Inst]

1.02 (+0.00%)

Avg. Issued Instructions Per Scheduler [Inst]

32,640 (+88.74%)

FP32 Non-Fused Instructions

This kernel executes 0 fused and 16,840 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).

Est. Local Speedup: 1.94%

Executed Instruction Categories

Integer

Load/Store

Control

Uniform Datapath

Miscellaneous

Floating Point

Movement

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 (+0.00%)

Function Cache Configuration

CACHEPERFNone (CACHEPERFNone)

Cluster Size

16 (+0.00%)

Function Scheduler Policy

PolicyShared (PolicyShared)

Registers Per Thread [register/thread]

18 (10.00%)

Cluster Scheduling Policy

PolicyShared (PolicyShared)

Static Shared Memory Per Block [byte/block]

1.30 (+0.00%)

Block Size

256 (+0.00%)

Dynamic Shared Memory Per Block [byte/block]

0 (+0.00%)

Threads [thread]

10,48,576 (+0.00%)

Shared Shared Memory Per Block [byte/block]

1.02 (+0.00%)

Waves Per SM

18.56 (+0.00%)

Shared Memory Configuration Size [byte]

32,77 (+100.00%)

Uses Green Context

0 (+0.00%)

Block Size

1,024 (+0.00%)

Uses 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 Occupancy [%]

100 (+0.00%)

Block Limit Shared Mem [block]

13 (+18.75%)

Theoretical Active Warps per SM [warp]

48 (+0.00%)

Block Limit Warps [block]

6 (+0.00%)

Achieved Occupancy [%]

89.91 (+20.94%)

Block Limit Warps [block]

24 (+0.00%)

Achieved Active Warps per SM [warp]

42.73 (+20.94%)

Block Limit Warps [block]

24 (+0.00%)

Cluster Occupancy [%]

0 (+0.