

Result

680-brute_force_AL

Size

(1024, 1)x(1024, 1) 1.16 ms

Time

0-Teala T4 584.99 Mhz

Cycles

GPU

GPU

SM Frequency

Process

16274.exe

Attributes

Summary

Details

Source

Context

Comments

Raw

Session

Compare

Tools

View

Export

GPU Speed of Light Throughput

All

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 [%]

83.75

Duration [ms]

1.16

Memory Throughput [%]

58.70

Elapsed Cycles [cycle]

676,917

L1/TEX Cache Throughput [%]

61.78

SM Active Cycles [cycle]

661,657.32

L2 Cache Throughput [%]

51.81

SM Frequency [Mhz]

584.99

DRAM Throughput [%]

1.88

DRAM Frequency [Ghz]

4.99

High Throughput

The kernel 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](#) section

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 19% of this device's fp32 peak performance and 0% of its fp64 peak performance. See the [Kernel Roofline Guide](#) for more details on roofline analysis.

GPU Throughput

Compute (SM) [%]

Memory [%]

Speed of Light (SOL) [%]

Compute Throughput Breakdown

Memory Throughput Breakdown

Floating Point Operations Roofline

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 [cycle]

80,000

Pass Groups

1

Maximum Buffer Size [bytes]

3.62

Dropped Samples [sample]

0

SM

DRAM

L1 Cache

Workload Execution

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 Elapsed [Inst/cycle]

3.35

SM Busy [%]

85.49

Executed Ipc Active [Inst/cycle]

3.42

Issue Slots Busy [%]

85.49

Issued Ipc Active [Inst/cycle]

3.42

High Utilization

ALU is the highest-utilized pipeline (69.5%) based on active cycles, taking into account the rates of its different instructions. It executes integer and logic operations. The pipeline is well-utilized, but might become a bottleneck if more work is added. Based on the number of executed instructions, the highest utilized pipeline (69.5%) is ALU. It executes integer and logic operations. Comparing the two, the overall pipeline utilization appears to be caused by frequent, low-latency instructions. See the [Kernel Roofline](#) section for more details on roofline analysis.

Pipe Utilization (% of active cycles)

Pipe Utilization (% of peak instructions executed)

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 throughput of issuing memory instructions (Mem Pipes Busy). Detailed chart of the memory units. Detailed tables with data for each memory unit.

Memory Throughput [byte/s]

6.01

Mem Busy [%]

30.88

L1/TEX Hit Rate [%]

98.70

Max Bandwidth [%]

58.70

L2 Hit Rate [%]

100.26

Mem Pipes Busy [%]

58.70

Memory L2 Compression

The optional metric lts__average_gcomp_input_sector_success_rate_pct could not be found. Collecting it as an additional metric could enable the rule to provide more guidance.

L1/TEX Local Load Access Pattern

The memory access pattern for local loads from L1/TEX might not be optimal. On average, only 0.2 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 Code](#) section for uncoalesced local loads.

L1/TEX Local Store Access Pattern

The memory access pattern for local stores to L1/TEX might not be optimal. On average, only 0.2 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 Code](#) section for uncoalesced local stores.

Memory Chart

Kernel

Global

Local

Texture

Surface

Shared

L1/TEX Cache

L2 Cache

System Memory

Device Memory

Peer Memory

Shared Memory

Shared Load

Shared Load Matrix

Shared Store

Shared Atomic

Other

Total

L1/TEX Cache

L2 Cache

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]

7.85

No Eligible [%]

13.58

Active Warps Per Scheduler [warp]

3.31

One or More Eligible [%]

85.42

Issued Warp Per Scheduler

0.86

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 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]

9.09

Avg. Active Threads Per Warp

32

Warp Cycles Per Executed Instruction [cycle]

9.09

Avg. Not Predicted Off-Threads Per Warp

30.46

Warp State (All Cycles)

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]

90,505,216

Avg. Executed Instructions Per Scheduler [Inst]

565,657.60

Issued Instructions [Inst]

90,508,126

Avg. Issued Instructions Per Scheduler [Inst]

565,675.79

FP32 Non-Fused Instructions

This kernel executes 7045120 fused and 7045120 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). Check the Source page to identify where this kernel executes FP32 instructions.

Est. Speedup: 12.55%

Executed Instruction Mix

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.

NVLink 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

1,024

Function Cache Configuration

Cache/Prefer/None

Registers Per Thread [register/thread]

22

Static Shared Memory Per Block [bytes/block]

0

Block Size

1,024

Dynamic Shared Memory Per Block [bytes/block]

0

Threads [thread]

1,048,576

Driver Shared Memory Per Block [bytes/block]

0

Waves Per SM

25.60

Shared Memory Configuration Size [kbyte]

32.77

Uses From Context

0

SMs [SM]

40

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 latency, resulting in overall performance degradation. Large discrepancies between the theoretical and the achieved occupancy during execution typically indicates highly imbalanced workloads.

Theoretical Occupancy [%]

100

Block Limit Registers [block]

2

Theoretical Active Warps per SM [warp]

32

Block Limit Shared Mem [block]

16

Achieved Occupancy [%]

97.31

Block Limit Warps [block]

1

Achieved Active Warps per SM [warp]

31.14

Block Limit SM [block]

16

Impact of Varying Register Count Per Thread

Impact of Varying Block Size

Impact of Varying Shared Memory Usage Per Block

GPU and Memory Workload Distribution

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

Analysis SM Active Cycles [cycle]

661,657.32

Average L1 Active Cycles [cycle]

661,657.32

Analysis L2 Active Cycles [cycle]

289,170.84

Average SMP Active Cycles [cycle]

654,575.04

Analysis DRAM Active Cycles [cycle]

108,686

Total SM Elapsed Cycles [cycle]

27,018,608

Total L2 Elapsed Cycles [cycle]

31,657,260

Total SMP Elapsed Cycles [cycle]

46,161,920

Workload Distribution

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]

4,882,432

Branch Efficiency [%]

100

Branch Instructions Ratio [%]

0.05

Avg. Divergent 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.