

Result

680 -brute_force_AL

Size

(1024, 1, 1)x(1024, 1, 1)

Time

1.10 ms

Cycles

642,480

GPU

0 -Tesla T4

SM Frequency

584.97 Mhz

Process

[14702] exe

Attributes

@

Summary

Details

Source

Context

Comments

Raw

Session

Compare

Tools

View

Export

Menu

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 [%]	84.28	Duration [ms]	1.10
Memory Throughput [%]	1.14	Elapsed Cycles [cycle]	642,480
L1/TEX Cache Throughput [%]	1.02	SM Active Cycles [cycle]	623,261.52
L2 Cache Throughput [%]	0.45	SM Frequency [Mhz]	584.97
DRAM Throughput [%]	1.14	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 [Kernel Profiling Guide](#).

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 2% of this device's fp32 peak performance and 0% of its fp64 peak performance. See the [Kernel Profiling 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

Peak Throughput (GOPS)

Arithmetic Intensity (FLOP/byte)

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]

40,000

Pass Groups

1

Maximum Buffer Size [bytes]

3

Dropped Samples [sample]

0

SM

Average Active Warps Per Cycle

Total Active Warps Per Cycle

SM Active Cycles

Executed Ipc Active

DRAM

DRAM Throughput

DRAM Read Bandwidth

DRAM Write Bandwidth

L1 Cache

Writeback Throughput

Hit Rate

Wavefronts (Data)

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]

2.49

SM Busy [%]

86.85

Executed Ipc Active [inst/cycle]

2.57

Issue Slots Busy [%]

64.15

Issued Ipc Active [inst/cycle]

2.57

Very High Utilization

ALU is the highest-utilized pipeline (86.8%) based on active cycles, taking into account the rates of its different instructions. It executes integer and logic operations. The pipeline is over-utilized and likely a performance bottleneck. Based on the number of executed instructions, the highest utilized pipeline (86.8%) 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 Profiling Guide](#) or hover over the pipeline name to understand the workloads handled by each pipeline. The [Instruction Mix](#) section shows the mix of executed instructions in this kernel. Check the [Warp State](#) section for which reasons cause warps to stall.

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]

3.66

Mem Busy [%]

0.45

L1/TEX Hit Rate [%]

98.98

Max Bandwidth [%]

1.14

L2 Hit Rate [%]

98.98

Mem Pipes Busy [%]

38.53

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.

Memory Chart

Kernel

Global

Local

Texture

Surface

Shared

L1/TEX Cache

L2 Cache

System Memory

Device Memory

Peer Memory

Shared Memory

Instructions

Requests

Wavefronts

% Peak

Bank Conflicts

L1/TEX Cache

Instructions

Requests

Wavefronts

% Peak

Hit Rate

Sectors

Sectors/Req

Hit Rate

Bytes

Sector Misses to L2

% Peak to L2

Returns to SM

L2 Cache

Requests

Sectors

Sectors/Req

% Peak

Hit Rate

Bytes

Sector Misses to Device

Sector Misses to System

Sector Misses to Peer

Device Memory

Sectors

% Peak

Bytes

Throughput

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]

7.75

No Eligible [%]

30.79

Issued Warps Per Scheduler [warp]

3.96

One or More Eligible [%]

69.21

Issued Warp Per Scheduler [cycle]

0.69

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]

11.19

Avg. Active Threads Per Warp

24.15

Warp Cycles Per Executed Instruction [cycle]

11.20

Avg. Not Predicated Off Threads Per Warp

21.70

Not Selected Stalls

On average, each warp of this kernel spends 3.9 cycles being stalled waiting for the micro scheduler to select the warp to issue. Not selected warps are eligible warps that were not picked by the scheduler to issue that cycle as another warp was selected. A high number of not selected warps typically means you have sufficient warps to cover warp latencies and you may consider reducing the number of active warps to possibly increase cache coherence and data locality. This stall type represents about 34.4% of the total average of 11.2 cycles between issuing two instructions.

Est. Local Speedup: 34.4%

Math Pipe Throttle Stalls

On average, each warp of this kernel spends 3.4 cycles being stalled waiting for the execution pipe to be available. This stall occurs when all active warps execute their next instruction on a specific, oversubscribed math pipeline. Try to increase the number of active warps to hide the existent latency or try changing the instruction mix to utilize all available pipelines in a more balanced way. This stall type represents about 30.8% of the total average of 11.2 cycles between issuing two instructions.

Est. Local Speedup: 30.79%

Warp Stall

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

Thread Divergence

Instructions are executed in warps, which are groups of 32 threads. Optimal instruction throughput is achieved if all 32 threads of a warp execute the same instruction. The chosen launch configuration, early thread completion, and divergent flow control can significantly lower the number of active threads in a warp per cycle. This kernel achieves an average of 24.1 threads being active per cycle. This is further reduced to 21.7 threads per warp due to predication. The compiler may use predication to avoid an actual branch. Instead, all instructions are scheduled, but a per-thread condition code or predicate controls which threads execute the instructions. Try to avoid different execution paths within a warp when possible.

Est. Speedup: 27.11%

Warp State (All Cycles)

Warp States

Cycles per Instruction

Instruction Statistics

Statistics of the executed low-level assembly instructions (OAS). The instruction mix provides insight into the types and frequency of the executed instructions. A diverse 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]

63,963,136

Avg. Executed Instructions Per Scheduler [inst]

399,769.60

Issued Instructions [inst]

63,973,086

Avg. Issued Instructions Per Scheduler [inst]

399,831.79

FP32 Non-Fused Instructions

This kernel executes 0 fused and 491,520 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). Check the Source page to identify where this kernel executes FP32 instructions.

Est. Speedup: 8.97%

Executed Instruction Mix

Opcodes

Executed Warp-Local 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 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

CachePreferenceNone

Registers Per Thread [register/thread]

17

Static Shared Memory Per Block [byte/block]

0

Block Size

1,024

Dynamic Shared Memory Per Block [byte/block]

3.07

Threads [thread]

1,048,576

Driver Shared Memory Per Block [byte/block]

0

Waves Per SM

25.60

Shared Memory Configuration Size [kbyte]

32.77

Uses Green 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 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]

2

Theoretical Active Warps per SM [warp]

32

Block Limit Shared Mem [block]

10

Achieved Occupancy [%]

89.84

Block Limit Warps [block]

1

Achieved Active Warps per SM [warp]

28.75

Block Limit SM [block]

16

Achieved Occupancy

The difference between calculated theoretical (100.0%) and measured achieved occupancy (89.8%) 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.

Est. Local Speedup: 10.16%

Impact of Varying Register Count Per Thread

Warp Occupancy

Registers Per Thread

Impact of Varying Block Size

Warp Occupancy

Block Size

Impact of Varying Shared Memory Usage Per Block

Warp Occupancy

Shared Memory Per Block

GPU and Memory Workload Distribution

Analysis of workload distribution in active cycles of GPU, SM, SMSP, L1 & L2 caches, and DRAM.

Average SM Active Cycles [cycle]

623,261.55

Average L1 Active Cycles [cycle]

623,261.55

Average L2 Active Cycles [cycle]

44,179.50

Average SMSP Active Cycles [cycle]

577,681.41

Average DRAM Active Cycles [cycle]

563,388

Total SM Active Cycles [cycle]

25,691.080

Total L2 Elapsed Cycles [cycle]

25,691.080

Total L1 Elapsed Cycles [cycle]

30,048.064

Total SMSP Elapsed Cycles [cycle]

102,764.220

Total DRAM Elapsed Cycles [cycle]

43,871.232

Workload Distribution

SM Active Cycles

SMSP Active Cycles

L1 Active Cycles

L2 Active Cycles

DRAM Active Cycles

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 scheduled warps fail to issue every cycle.

Branch Instructions [inst]

5,357,568

Branch Efficiency [%]

96.36

Branch Instructions Ratio [%]

0.08

Avg. Divergent Branches

1.024

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.