

Result

660 - brute\_force\_AL

Size

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

Time

1.51 ms

Cycles

885,956

GPU

0 - Tesla T4

SM Frequency

584.98 Mhz

Process

[18142] 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 [%]	84.55	Duration [ms]	1.51
Memory Throughput [%]	84.55	Elapsed Cycles [cycle]	885,956
L1/TEX Cache Throughput [%]	86.46	SM Active Cycles [cycle]	866,823.75
L2 Cache Throughput [%]	2.91	SM Frequency [Mhz]	584.98
DRAM Throughput [%]	1.21	DRAM Frequency [Ghz]	5.00

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 [P-Compute](#) section.

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32.1. The kernel achieved 15% 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

Performance [FLOP/s] (1-100,000,000,000)

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]160,000

# Pass Groups1

Maximum Buffer Size [Mbytes]2.25

Dropped Samples [sample]0

SM

Average Active Warps Per Cycle62.45 warp

Total Active Warps Per Cycle1.25k warp

SM Active Cycles109k cycle

Executed Ipc Active2.57 inst/cycle

Executed Ipc Active0

DRAM

DRAM Throughput100 %

DRAM Read Bandwidth100 %

DRAM Write Bandwidth100 %

L1 Cache

Writeback Throughput44.3k cycle

Hit Rate100 %

Wavefronts (Data)48.3k

Workload Executionbrute\_force\_AL

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 [%]63.72

Executed Ipc Active [inst/cycle]2.55

Issue Slots Busy [%]63.72

Issued Ipc Active [inst/cycle]2.55

Balanced

ALU is the highest-utilized pipeline (52.0%) based on active cycles, taking into account the rates of its different instructions. It executes integer and logic operations. It is well-utilized, but should not be a bottleneck.

Pipe Utilization (% of active cycles)

ALU

FMA

FP64

Shared (FP16+Tensor)

Tensor (All)

Tensor (FP)

Tensor (INT)

Utilization [%]

Pipe Utilization (% of peak instructions executed)

LSU

ALU

XU

FMA

ADU

CBU

FP16

FP64

TEX

Tensor (FP)

Tensor (INT)

Uniform

Utilization [%]

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 [Gbytes/s]3.88

Mem Busy [%]45.20

L1/TEX Hit Rate [%]98.92

Max Bandwidth [%]84.55

L2 Hit Rate [%]97.57

Mem Pipes Busy [%]84.55

Memory L2 Compression

The optional metric l2s\_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.

L1TEX Global Load Access Pattern

Est. Speedup: 2.40x

The memory access pattern for global loads from L1TEX might not be optimal. On average, only 31.1% of the 32 bytes transmitted per sector are utilized by each thread. This could possibly be caused by a stride between threads. Check the [P-Access Columns](#) section for uncoalesced global loads.

L1TEX Local Load Access Pattern

Est. Speedup: 85.71%

The memory access pattern for local loads from L1TEX 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 [P-Access Columns](#) section for uncoalesced local loads.

L1TEX Local Store Access Pattern

Est. Speedup: 85.71%

The memory access pattern for local stores to L1TEX 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 [P-Access Columns](#) section for uncoalesced local stores.

Memory Chart

Show As: Transfer Size

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

Sectors

Sectors/Req

Hit Rate

Bytes

Sector Misses to L2

% Peak to L2

Returns to SM

Device Memory

Sectors

% Peak

Bytes

Throughput

Sector Misses to Device

Sector Misses to System

Sector Misses to Peer

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 warps are required to hide this latency. For each warp state, the chart shows the average number of cycles spent in that state per issued 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.

Warp Cycles Per Issued Instruction [cycle]12.27

Avg. Active Threads Per Warp32.00

Warp Cycles Per Executed Instruction [cycle]12.27

Avg. Not Predicted Off Threads Per Warp30.42

Mio Throttle Stalls

Est. Local Speedup: 38.90%

On average, each warp of this kernel spends 4.8 cycles being stalled waiting for the MIO (memory input/output) instruction queue to be not full. This stall reason is high in cases of extreme utilization of the MIO pipelines, which include special math instructions, dynamic branches, as well as shared memory instructions. When caused by shared memory accesses, trying to use fewer but wider loads can reduce pipeline pressure. This stall type represents about 38.9% of the total average of 12.3 cycles between issuing two instructions.

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.

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]88,368,128

Avg. Executed Instructions Per Scheduler [inst]552,300.80

Issued Instructions [inst]88,372,248

Avg. Issued Instructions Per Scheduler [inst]552,326.55

FP32 Non-Fused Instructions

Est. Speedup: 9.02%

This kernel executes 7045120 fused and 7045120 non-fused FP32 instructions. By converting pairs of 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.

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 Size1,024

Function Cache ConfigurationCache/PreferNone

Registers Per Thread [register/thread]	23	Static Shared Memory Per Block [byte/block]	0
Block Size	1,024	Dynamic Shared Memory Per Block [byte/block]	840
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 Branch 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 Occupancy [%]100

Block Limit Registers [block]32

Theoretical Active Warps per SM [warp]32

Block Limit Shared Mem [block]32

Achieved Occupancy [%]97.75

Block Limit Warps [block]31.28

Achieved Active Warps per SM [warp]16

Block Limit SM [blocks]

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 SM, SMP, SMSP, L1 & L2 caches, and DRAM

Average SM Active Cycles [cycle]	866,823.75	Average L1 Active Cycles [cycle]	866,823.75
Average L2 Active Cycles [cycle]	185,502.50	Average SMSP Active Cycles [cycle]	862,883.81
Average DRAM Active Cycles [cycle]	91,909	Total SM Elapsed Cycles [cycle]	35,457,416
Total L1 Elapsed Cycles [cycle]	35,457,416	Total L2 Elapsed Cycles [cycle]	41,434,528
Total SMSP Elapsed Cycles [cycle]	141,829,664	Total DRAM Elapsed Cycles [cycle]	60,613,632

Workload Distribution

Kernel

Min

Max

Sum

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]4,956,160

Branch Efficiency [%]

Branch Instructions Ratio [%]0.06

Avg. Divergent Branches6.40

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 [certain sections](#) to focus on selected performance aspects and make profiling faster.