

Result

681 - brute_force_AL

Size

(4, 1, 1)x(1024, 1, 1)

Time

9.39 ms

Cycles

5,491,446

GPU

0 - Tesla T4

SM Frequency

584.99 Mhz

Process

[12774] exe

Attributes

Summary

Details

Source

Context

Comm

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 rooftop chart.

Compute (SM) Throughput [%]	8.81	Duration [ms]	9.39
Memory (SM) Throughput [%]	0.00	Elapsed Cycles [cycle]	5,491,446
L1/TEX Cache Throughput [%]	0.01	SM Active Cycles [cycle]	541,649.50
L2 Cache Throughput [%]	0.00	SM Frequency [Mhz]	584.99
DRAM Throughput [%]	0.00	DRAM Frequency [Ghz]	5.00

Small Grid

This kernel grid is too small to fill the available resources on this device, resulting in only 0.1 full waves across all SMs. Look at [Launch Strategies](#) for more details.

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved close to 0% of this device's fp32 peak performance and close to 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

SM: Pipe Als Cycles Active [%]

SM: Issue Active [%]

SM: Inst Executed [%]

SM: Inst Executed Pipe Adu [%]

IDC: Request Cycles Active [%]

SM: Pipe Fma Cycles Active [%]

SM: Mio Pq Write Cycles Active [%]

SM: Mio Inst Issued [%]

SM: Mio Pq Read Cycles Active [%]

SM: Inst Executed Pipe Cbu Pred On Any [%]

SM: Mio2r Writeback Active [%]

SM: Pipe Fp64 Cycles Active [%]

SM: Inst Executed Pipe Lsu [%]

SM: Inst Executed Pipe Xu [%]

SM: Memory Throughput Internal Activity [%]

SM: Pipe Tensor Cycles Active [%]

SM: Pipe Shared Cycles Active [%]

SM: Instruction Throughput Internal Activity [%]

SM: Inst Executed Pipe Uniform [%]

SM: Inst Executed Pipe Tex [%]

SM: Inst Executed Pipe Isa [%]

SM: Inst Executed Pipe Fp16 [%]

L2: T Sectors [%]

L2: Xbar2ls Cycles Active [%]

L2: L4x2ls Cycles Active [%]

L1: M1l1as2lsar Req Cycles Active [%]

GPU: Compute Memory Access Throughput Internal Activity [%]

L2: D Sectors Fill Memory Access [%]

L2: D Sectors [%]

L2: T Tag Requests [%]

DRAM: Cycles Active [%]

L1: Lsuin Requests [%]

DRAM: Dram Sectors [%]

L1: Data Pipe Lsu Wavefronts [%]

L1: Data Bank Writes [%]

L1: Data Bank Reads [%]

L1: Lsu Writeback Active [%]

L1: Train Sectors Req Cycles Active [%]

L1: F Wavefronts [%]

L1: Tex Writeback Active [%]

L2: D Atomic Input Cycles Active [%]

L1: M4 Xbar2l1as Read Sectors [%]

L2: D Sectors Fill System [%]

GPU: Compute Memory Request Throughput Internal Activity [%]

L1: Data Pipe Tex Wavefronts [%]

Floating Point Operations Roofline

Performance [FLOP/s] (FLOP/s = 1000000000000)

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]

Maximum Buffer Size [byte]

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]

Executed Ipc Active [Inst/cycle]

Issued Ipc Active [Inst/cycle]

ALU is the highest-utilized pipeline (89.3%) 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 (89.3%) 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](#) for more details on pipeline utilization. The [Performance Analysis](#) section shows the mix of executed instructions in this kernel. Check the [Launch Strategies](#) 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 [Mbytes/s]

L1/TEX Hit Rate [%]

L2 Hit Rate [%]

Memory L2 Compression

The optional metric l2_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

Show As: Transfer Size

Kernel

Global

Local

Texture

Surface

Shared

L1/TEX Cache

L2 Cache

Peer Memory

Peer Memory

Peer Memory

% Peak

Shared Memory

Instructions

Requests

Wavefronts

% Peak

Bank Conflicts

Shared Load

Shared Load Matrix

Shared Store

Shared Atomic

Other

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

Local Load

Global Load

Surface Load

Texture Load

Global Store

Local Store

Surface Store

Global Atomic CAS

Surface Atomic CAS

Global Atomic CAS

Surface Atomic CAS

Loads

Stores

Atomic & Reductions

L2 Cache

Requests

Sectors

Sectors/Req

% Peak

Hit Rate

Bytes

Throughput

Sector Misses to Device

Sector Misses to System

Sector Misses to Peer

L1/TEX Load

L1/TEX Store

L1/TEX Atomic ALU

L1/TEX Atomic CAS

L1/TEX Reduction

L1/TEX Total

Loc Total

GPU Total

Device Memory

Sectors

% Peak

Bytes

Throughput

Load

Store

Total

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]

Eligible Warps Per Scheduler [warp]

Issued Warp Per Scheduler

Warps 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 warps are stalled in that state. 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 34.3% of the total average of 11.3 cycles between issuing two instructions.

Warp Cycles Per Issued Instruction [cycle]

Warp Cycles Per Executed Instruction [cycle]

Not Selected: 34.31%

Math Pipe Throttle Stalls: Est. Local Speedup: 30.06%

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.

Warp State (All Cycles)

Stall Not Selected

Stall Math Pipe Throttle

Stall War

Selected

Stall Short Scoreboard

Stall Branch Resolving

Stall No Instruction

Stall Tex Throttle

Stall Long Scoreboard

Stall Dispatch Stall

Stall MID Throttle

Stall IMC Miss

Stall Drain

Stall Misc

Stall Barrier

Stall LG Throttle

Stall Member

Stall Sleeping

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 indicates 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]

Issued Instructions [Inst]

FP32 Non-Fused Instructions

Est. Speedup: 9.78%

FP64 Non-Fused Instructions

Est. Speedup: 0.17%

Executed Instruction Mix

NVLink Topology

NVLink Topology diagram shows logical NVLink connections with transmit/receive throughput.

NVLink Tables

Detailed tables with properties for each NVLink.

NJMA Affinity

Non-uniform memory access (NJMA) 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

Registers Per Thread [register/thread]

Block Size

Threads Per Block

Waves Per SM

Uses Group Context

Small Grid

The grid for this launch is configured to execute only 4 blocks, which is less than the GPU's 40 multiprocessors. This can underutilize some multiprocessors if you do not intend to execute this kernel concurrently with other workloads, consider reducing the block size to have at least one block per multiprocessor or increase the size of the grid to fully utilize the available hardware resources. See the [Launch Strategies](#) section for more details on launch configurations.

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 [%]

SMSP Active Cycles [%]

Achieved Occupancy [%]

Achieved Active Warps Per SM [warp]

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, SMSP, L1 & L2 caches, and DRAM.

SM Active Cycles

SMSP Active Cycles [cycle]

Average DRAM Active Cycles [cycle]

Total L1 Elapsed Cycles [cycle]

Total L2 Elapsed Cycles [cycle]

SMs Workload Imbalance

Est. Speedup: 8.89%

SMSPs Workload Imbalance

Est. Speedup: 8.87%

L1 Slices Workload Imbalance

Est. Speedup: 8.89%

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]

Branch Instructions Ratio [%]