

Result

550 - brute\_force

Size

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

Time

4.23 ms

CPU

2,473,819

GPU

0 - Tesla T4

SM Frequency

584.99 MHz

Process

[19994] exe

Attributes

Summary

Details

Source

Context

Command

Raw

Session

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 roffline chart.

Compute (SM) Throughput [%]	44.32	Duration [ms]	4.23
Memory Throughput [%]	38.88	Elapsed Cycles [cycle]	2,473,819
L1/TEX Cache Throughput [%]	50.78	SM Active Cycles [cycle]	2,408,916
L2 Cache Throughput [%]	18.32	SM Frequency [MHz]	584.99
DRAM Throughput [%]	30.42	DRAM Frequency [GHz]	5.01

Latency Issue

This kernel 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 the [Scheduler Statistics](#) and [Warp State Statistics](#) for potential reasons.

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32.1. The kernel achieved 4% of this device's fp32 peak performance and 9% of its fp64 peak performance. See the [Kernel Roofline Guide](#) for more details on roofline analysis.

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]	160,000	# Pass Groups	1
Maximum Buffer Size [bytes]	3.81	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]	1.22	SM Busy [%]	45.50
Executed IPC Active [inst/cycle]	1.25	Issue Slots Busy [%]	31.20
Issued IPC Active [inst/cycle]	1.25		

Balanced

FP64 is the highest-utilized pipeline (45.5%) based on active cycles, taking into account the rates of its different instructions. It is well-utilized, but should not be a bottleneck.

Pipeline Utilization (% of active cycles)

Pipeline 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 Bus), exhausting the available communication bandwidth between those units (Max Bandwidth), or by reaching the maximum throughput of issuing memory instructions (Mem Pipes Busy). Detailed tables with data for each memory unit.

Memory Throughput [GB/s]	97.61	Mem Busy [%]	25.39
L1/TEX Hit Rate [%]	84.90	Max Bandwidth [%]	38.88
L2 Hit Rate [%]	96.83	Mem Pipes Busy [%]	38.88

Memory L2 Compression

The optional metric l2s\_average\_gcomp\_input\_success\_rate\_pct could not be found. Collecting it as an additional metric could enable the rule to provide more guidance.

L1/TEX Global Load Access Pattern

Ext. Speedup: 45.24%

The memory access pattern for global loads from L1/TEX might not be optimal. On average, only 3.5 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 Columns](#) section for uncoalesced global loads.

L1/TEX Local Load Access Pattern

Ext. Speedup: 50.03%

The memory access pattern for local loads from L1/TEX might not be optimal. On average, only 0.5 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 Columns](#) section for uncoalesced local loads.

L1/TEX Global Load Access Pattern

Ext. Speedup: 49.93%

The memory access pattern for local stores to L1/TEX might not be optimal. On average, only 0.9 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 Columns](#) section for uncoalesced local stores.

Memory Chart

Show As: Transfer Size

Shared Memory

L1/TEX Caches

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. Leaving many skipped issue slots indicates pool latency hiding.

Active Warps Per Scheduler [warp]	7.79	No Eligible [%]	67.54
Eligible Warps Per Scheduler [warp]	1.07	One or More Eligible [%]	32.46
Issued Warps Per Scheduler	0.32		

Issue Slot Utilization

Ext. Load Speedup: 55.56%

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 3.1 cycles. This might leave hardware resources underutilized and may lead to less optimal performance. Out of the maximum of 8 warps per scheduler, this kernel allocates an average of 7.79 active warps per scheduler, but only an average of 1.07 warps were eligible per cycle. Eligible warps are the subset of active warps that are ready to issue their next instruction. Every cycle with no eligible warps results in no instruction being issued and the issue slot remains unused. To increase the number of eligible warps, avoid possible load imbalances due to highly different execution durations per warp. Reducing stalls indicated on the [Warp State Statistics](#) and [Source Columns](#) sections can help, too.

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 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 warps, the issue slot is skipped and no instruction is issued. Leaving many skipped issue slots indicates pool latency hiding.

23.99	Avg. Active Threads Per Warp	25.81
23.99	Avg. Not Predicted Off Threads Per Warp	24.61

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.

120,243,488	Avg. Executed Instructions Per Scheduler [inst]	751,521.80
120,249,157	Avg. Issued Instructions Per Scheduler [inst]	751,557.23

FP32 Non-Fused Instructions

Ext. Speedup: 5.20%

This kernel executes 660,480 fused and 660,480 non-fused FP32 instructions. By converting pairs of non-fused FP32 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.

FP64 Non-Fused Instructions

Ext. Speedup: 22.75%

This kernel executes 0 fused and 122,308 non-fused FP64 instructions. By converting pairs of non-fused FP64 instructions to their [fused](#), higher-throughput equivalent, the achieved FP64 performance could be increased by up to 50% (relative to its current performance). Check the Source page to identify where this kernel executes FP64 instructions.

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.

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 threads, 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/PreferNone
Registers Per Thread [register/thread]	33	Static Shared Memory Per Block [byte/block]	0
Block Size	1,024	Dynamic Shared Memory Per Block [byte/block]	0
Threads/Block	1,048,376	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 resulting number of possible active warps. Large discrepancies between the theoretical and the achieved occupancy during execution typically indicates highly imbalanced workloads.

100	Block Limit Registers [block]	1
32	Block Limit Shared Mem [block]	1
93.47	Block Limit Warps [block]	16
29.91	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, SMSP, L1 & L2 caches, and DRAM

Average SM Active Cycles [cycle]	2,408,916	Average L1 Active Cycles [cycle]	2,408,916
Average L2 Active Cycles [cycle]	2,345,925.22	Average DRAM Active Cycles [cycle]	2,315,699.29
Average DRAM Active Cycles [cycle]	6,449,875	Total SM Elapsed Cycles [cycle]	96,921,872
Total L1 Elapsed Cycles [cycle]	98,921,872	Total L2 Elapsed Cycles [cycle]	115,697,920
Total SMSP Elapsed Cycles [cycle]	395,687,488	Total DRAM Elapsed Cycles [cycle]	169,606,144

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.

8,464,056	Branch Efficiency [%]	97.34
0.07	Avg. Divergent Branches	1,237.40

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 individual sections to focus on selected performance aspects and make profiling faster.