

Result

660 - brute\_force\_AL

Size

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

Time

2.10 ms

Cycles

1,228,788

GPU

0 - Tesla T4

SM Frequency

584.98 MHz

Process

[19994] exe

Attributes

0x00000000

Summary

Details

Source

Context

Command

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 roiline chart.

Compute (SM) Throughput [%]	90.10	Duration [ms]
Memory Throughput [%]	60.91	Elapsed Cycles [cycle]
L1/TEX Cache Throughput [%]	62.96	SM Active Cycles [cycle]
L2 Cache Throughput [%]	2.14	SM Frequency [MHz]
DRAM Throughput [%]	0.89	DRAM Frequency [GHz]

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](#) [Workload Analysis](#) section.

Refline 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 roline 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]

Maximum Buffer Size [bytes]

Pass Groups

Dropped Samples [sample]

SM

DRAM

L1 Cache

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]

Balanced

FMA is the highest-utilized pipeline (43.1%) based on active cycles, taking into account the rates of its different instructions. It executes 32-bit floating-point (FADD, FMUL, FMAD, ...) and integer (MUL, IMAD) operations. It is well-utilized, but should not be a bottleneck.

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 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 chart of the memory units. Detailed tables with data for each memory unit.

Memory Throughput [Gbyte/s]

L1/TEX Hit Rate [%]

L2 Hit Rate [%]

Memory L2 Compression

Memory L2 Compression

L1TEX Global Load Access Pattern

L1TEX Local Load Access Pattern

L1TEX Local Store Access Pattern

Memory Chart

Shared Memory

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]

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 warp parallelism is required to hide this latency. For each warp state, the chart shows the average number of cycles spent in this 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]

Warp Cycles Per Executed Instruction [cycle]

Mio Throttle Stalls

Warp Stall

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 few pipelines allows hiding latencies and enables parallel execution. Note that 'Instructions/Opcode' and 'Executed Instructions' are measured differently and can differ if cycles are spent in system calls.

Executed Instructions [Inst]

Achieved Occupancy [%]

Executed Instructions Per Scheduler [Inst]

Block Limit Warps [Block]

FP32 Non-Fused Instructions

Executed Instruction Mix

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

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 [%]

Follow the rules outlined to get guidance on how to navigate through the report and quickly discover performance bottlenecks in this kernel. You could also disable [Warp Stall Sampling](#) to focus on selected performance aspects and make profiling faster.