

Launch0 - 508 - naive_kernel

TimeCyclesReqsGPU

508 - naive_kernel(256, 256, 1)(C32, 32, 1)6.12 mseccond8,379,59628NVIDIA GeForce GTX 16501.37 cycle/mseccond7.5 [1151] res

GPU Speed of Light Throughput

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 contributors. High-level overview of the utilization for compute and memory resources of the GPU presented as a roofline chart.

Compute (SM) Throughput [%]

Memory (SM) Throughput [%]

L1/TEX Cache Throughput [%]

L2 Cache Throughput [%]

DRAM Throughput [%]

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 [Stallbar Statistics](#) and [Issue 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 8% of this device's fp32 peak performance and 0% of its fp64 peak performance.

GPU Throughput

Compute (SM) [%]

Memory [%]

Speed of Light (SOL) [%]

Compute Throughput Breakdown

Memory Throughput Breakdown

Floating Point Operations Roofline

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 lpc Elapsed [inst/cycle]

Executed lpc Active [inst/cycle]

Issued lpc Active [inst/cycle]

SM Busy [%]

Issue Slots Busy [%]

Balanced

No pipeline is over-utilized.

Pipe 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 (Memory 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. Deprecated UI elements for deprecated UI elements compatibility.

Memory Throughput [Gbyte/second]

L1/TEX Hit Rate [%]

L2 Hit Rate [%]

Mem Busy [%]

Max Bandwidth [%]

Mem Pipes Busy [%]

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 [warp]

No Eligible [%]

One or More Eligible [%]

Issue Slot Utilization

On average, each warp of this kernel spends 16.3 cycles being stalled waiting for a scoreboard dependency on a L1/TEX (local, global, surface, texture) operation. This represents about 54.2% of the total average of 30.1 cycles between issuing two instructions. To reduce the number of cycles waiting on L1/TEX data accesses verify the memory access patterns are optimal for the target architecture, attempt to increase cache hit rates by increasing data locality or by changing the cache configuration, and consider moving frequently used data to shared memory.

Warp State (All Cycles)

Instruction Statistics

Statistics of the executed low-level assembly instructions (EASIS). 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]

Issued Instructions [inst]

Static Shared Memory Per Block [byte/block]

Dynamic Shared Memory Per Block [byte/block]

Global Shared Memory Per Block [byte/block]

Shared Memory Configuration Size [Kbyte]

NVLink

High-level summary of NVLink utilization. It shows the total received and transmitted (per) memory, as well as the overall link peak utilization. NVLink Topology: Detailed tables with properties for each NVLink.

Physical Links

NVLink Topology

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

Block Size

Threads Per Block

Warp Size

Function Cache Configuration

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 [%]

Theoretical Active Warps Per SM [warp]

Achieved Occupancy [%]

Achieved Active Warps Per SM [warp]

Occupancy Limiters

This kernel's theoretical occupancy is not impacted by any block limit. The difference between calculated theoretical (100.0%) and measured achieved occupancy (89.4%) 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.

Impact of Varying Block Size

Impact of Varying Shared Memory Usage Per Block

Source Counters

Source metrics, including branch efficiency and sampled warp stall reasons. Sampling Data 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 [%]

Uncoalesced Global Accesses

Uncoalesced Global Accesses

Uncoalesced Global Accesses

Uncoalesced Global Accesses

Uncoalesced Global Accesses

Uncoalesced Global Accesses

Sampling Data (All)

Sampling Data (Not Issued)

Most Instructions Executed