

Result

559 - brute_force_coarsening

Size

Time

Cycles

GPU

SM Frequency

Process

Attributes

Summary

Details

Source

Context

Comments

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 rooiline chart.

Compute (SM) [%]	79.94	Duration [ms]	2.22
Memory Throughput [%]	11.00	Elapsed Cycle [cycle]	1,297,249
L1/TEX Cache Throughput [%]	11.85	SM Active Cycles [cycle]	1,292,228.20
L2 Cache Throughput [%]	0.65	SM Frequency [MHz]	584.99
DRAM Throughput [%]	1.97	DRAM Frequency [GHz]	4.99

Compute is more heavily utilized than Memory. Look at the [Compute Workload Analysis](#) section to see what the compute pipelines are spending their time doing. Also, consider whether any computation is redundant and could be reduced or moved to look-up tables.

GPU Throughput

Compute Throughput Breakdown

Memory Throughput Breakdown

Floating Point Operations Rooiline

Performance (FLOP/cycle) (1e+10/100,000,000.000)

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 [Mbytes]

Pass Groups

Dropped Samples [sample]

Timeline

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]

SM Busy [%]

Issue Slots Busy [%]

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 tables with data for each memory unit.

Memory Throughput [Gbytes/s]

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 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 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 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 States

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]

Issued Instructions [Inst]

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 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

Waves Per SM

Uses Green Context

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 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 [%]

Achieved Occupancy [%]

Theoretical Active Warps per SM [warp]

Achieved Active Warps per SM [warp]

Theoretical Occupancy

Local Speedup: 50.00%

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]

Average L2 Active Cycles [cycle]

Average DRAM Active Cycles [cycle]

Total L1 Elapsed Cycles [cycle]

Total SMSP Elapsed Cycles [cycle]

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 [%]

Branch Efficiency [%]

Avg. Divergent Branches

Warp Stall Sampling (All Samples)

Most Instructions Executed

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.