

Result

728 -brute_force_AL_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

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 (Est. Throughput) [%]

Memory Throughput [%]

L1/TEX Cache Throughput [%]

L2 Cache Throughput [%]

DRAM Throughput [%]

81.86

Duration [ms]

2.01

Elapsed Cycles [cycle]

1,264,490

2.16

1.55

SM Active Cycles [cycle]

1,256,878.15

0.67

SM Frequency [MHz]

584.99

0.00

2.01

DRAM Frequency [GHz]

4.99

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.

FP64/32 Utilization

Est. Speedup: 53.66%

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved close to 1% of this device's fp32 peak performance and 61% of its fp64 peak performance. If [Compute Workload Analysis](#) determines that this kernel is fp64 bound, consider using 32-bit precision floating point operations to improve its performance. See the [Kernel Profiling Guide](#) for more details on roffline analysis.

GPU Throughput

Compute (SM) [%]

Memory [%]

Speed of Light (SOL) [%]

Compute Throughput Breakdown

Memory Throughput Breakdown

Floating Point Operations Roffline

Performance [FLOP/s]
(N = 10,000,000,000)

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]

3.06

Dropped Samples [sample]

1

0

SM

Average Active Warps Per Cycle

31.84 warp

0

Total Active Warps Per Cycle

636.89 warp

0

SM Active Cycles

27.4k cycle

0

Executed lpc Active

1.50 inst/cycle

0

DRAM

DRAM Throughput

100 %

0

DRAM Read Bandwidth

100 %

0

DRAM Write Bandwidth

100 %

0

L1 Cache

Writeback Throughput

20.32 cycle

0

Hit Rate

100 %

0

Wavefronts (Data)

82.7%

0

Workload Execution

brute_force_AL_coarsening

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 Active [inst/cycle]

40,000

SM Busy [%]

1.37

Issue Slots Busy [%]

1.37

Issue Slots Busy [%]

34.36

Very High Utilization

FP64 is the highest-utilized pipeline (92.4%) based on active cycles, taking into account the rates of its different instructions. It executes 64-bit floating point operations. The pipeline is over-utilized and likely a performance bottleneck. Based on the number of executed instructions, the highest utilized pipeline (92.4%) is FP64. It executes 64-bit floating point 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 roffline analysis. 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 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 [Gbyte/s]

6.41

Mem Busy [%]

0.67

L1/TEX Hit Rate [%]

0

Max Bandwidth [%]

2.01

L2 Hit Rate [%]

99.41

Mem Pipes Busy [%]

26.89

Memory L2 Compression

The optional metric lts_average_gomp_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

Shared Memory

Instructions

Requests

Wavefronts

% Peak

Bank Conflicts

L1/TEX Cache

Instructions

Requests

Wavefronts

% Peak

Sectors

Sectors/Req

Hit Rate

Bytes

Sector Misses to L2

% Peak to L2

Returns to SM

L2 Cache

Requests

Sectors

Sectors/Req

% Peak

Hit Rate

Bytes

Throughput

Sector Misses to Device

Sector Misses to System

Sector Misses to Peer

Device Memory

Sectors

% Peak

Bytes

Throughput

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]

3.96

No Eligible [%]

65.61

Eligible Warps Per Scheduler [warp]

0.48

One or More Eligible [%]

34.39

Issued Warp Per Scheduler

0.34

Issue Slot Utilization: 18.14%

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 2.9 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 3.96 active warps per scheduler, but only an average of 0.48 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 warp 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 Stall Sampling \(All Samples\)](#) and [Warp Stall Sampling \(Per Warp\)](#) sections can help too.

Warps Per Scheduler

GPU Maximum Warps Per Scheduler

Theoretical Warps Per Scheduler

Active Warps Per Scheduler

Eligible Warps Per Scheduler

Issued Warp 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]

11.50

Avg. Active Threads Per Warp

24.77

Warp Cycles Per Executed Instruction [cycle]

11.51

Avg. Not Predicted Off Threads Per Warp

22.48

Long Executed Stalls

Est. Speedup: 18.14%

On average, each warp of this kernel spends 5.7 cycles being stalled waiting for a scoreboard dependency on a L1/TEX local/global/surface/textured operation. Find the instruction producing the data being waited upon to identify the culprit. 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 (coalescing), or by changing the cache configuration. Consider moving frequently used data to shared memory. This stall type represents about 49.4% of the total average of 11.5 cycles between issuing two instructions.

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.

Thread Divergence

Est. Speedup: 24.36%

Instructions are executed in warps, which are groups of 32 threads. Optimal instruction throughput is achieved if all 32 threads of a warp execute the same instruction. The chosen launch configuration, early thread completion, and divergent flow control can significantly lower the number of active threads in a warp per cycle. This kernel achieves an average of 24.8 threads being active per cycle. This is further reduced to 22.5 threads per warp due to predication. The compiler may use predication to avoid an actual branch. Instead, all instructions are scheduled, but a per-thread condition code or predicate controls which threads execute the instructions. Try to avoid different execution paths within a warp when possible.

Warp State (All Cycles)

Stall Long Scoreboard

Stall Wait

Stall Short Scoreboard

Selected

Stall Not Selected

Stall Math Pipe Throttle

Stall Branch Throttle

Stall MIO Throttle

Stall Tex Throttle

Stall No Instruction

Stall Drain

Stall Dispatch Stall

Stall IMC Miss

Stall Misc

Stall Barrier

Stall LG Throttle

Stall Membar

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 implies a 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]

69,091,228

Avg. Executed Instructions Per Scheduler [inst]

431,930.00

Executed Instructions [inst]

69,093,419

Avg. Issued Instructions Per Scheduler [inst]

431,833.87

FP32 Non-Fused Instructions

Est. Speedup: 5.14%

This kernel executes 0 fused and 491,320 non-fused FP32 instructions. By converting pairs of non-fused instructions to their [fused](#), higher-throughput equivalent, the achieved FP32 performance could be increased by up to 50% (relative to its current performance). Check the Source page to identify where this kernel executes FP32 instructions.

FP64 Non-Fused Instructions

Est. Speedup: 19.26%

This kernel executes 134,348 fused and 1,179,648 non-fused FP64 instructions. By converting pairs of non-fused instructions to their [fused](#), higher-throughput equivalent, the achieved FP64 performance could be increased by up to 23% (relative to its current performance). Check the Source page to identify where this kernel executes FP64 instructions.

Executed Instruction Mix

LDP

STP

SHF

IADD3

IMAD

LDC

FADD

BRA

ISETP

PRMT

DFMA

DADD

BSYNC

BSSV

BMOV

DMUL

LEA

SEL

FSEL

STG

S2R

FS2TP

CS2R

RET

MUJ

MOV

I2F

F2I

EXIT

CALL

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 Configuration

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

32,768

Function Cache Configuration

CachePreference

Registers Per Thread [register/thread]

26

Static Shared Memory Per Block [byte/block]

0

Block Size

82

Dynamic Shared Memory Per Block [byte/block]

0

Threads Per Thread

1,048,576

Driver Shared Memory Per Block [byte/block]

0

Waves Per SM

51.20

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 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 [%]

60

Block Limit Registers [block]

64

Theoretical Active Warps per SM [warp]

216.561

Block Limit Shared Mem [block]

32

Achieved Occupancy [%]

49.42

Block Limit Warps [block]

16

Achieved Active Warps Per Scheduler [warp]

15.81

Block Limit SM [block]

16

Theoretical Occupancy

Est. Speedup: 18.14%

The 4.00 theoretical warps per scheduler this kernel can issue according to its occupancy are below the hardware maximum of 8. This kernel's theoretical occupancy (50.0%) is limited by the number of blocks that can fit on the SM. This kernel's theoretical occupancy (50.0%) is limited by the required amount of shared memory.

Impact of Varying Register Count Per Thread

Warp Occupancy

Registers Per Thread

Impact of Varying Block Size

Warp Occupancy

Block Size

Impact of Varying Shared Memory Usage Per Block

Warp Occupancy

Shared Memory 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]

1,256,878.15

Average L1 Active Cycles [cycle]

109,970

Average SMSP Active Cycles [cycle]

1,256,878.15

Average L2 Active Cycles [cycle]

216.561

Total SM Elapsed Cycles [cycle]

50,599,792

Total L1 Elapsed Cycles [cycle]

90,138,976

Total SMSP Elapsed Cycles [cycle]

86,233,136

Total L2 Elapsed Cycles [cycle]

50,599,792

Workload Distribution

SM Active Cycles

1,256,878.15

Min

1,252,440

Max

1,260,697

Sum

50,275,126

SMSP Active Cycles

1,256,878.15

Min

1,250,094

Max

1,268,335

Sum

200,913,925

L1 Active Cycles

1,256,878.15

Min

1,252,440

Max

1,260,697

Sum

50,275,126

L2 Active Cycles

109,970

Min

109,285

Max

110,757

Sum

3,519,040

DRAM Active Cycles

216,561

Min

216,004

Max

217,504

Sum

1,732,488

Source Comments

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]

5,619,712

Branch Efficiency [%]

96.54

Branch Instructions Ratio [%]

0.08

Avg. Divergent Branches

1,024

Warp Stall Sampling (All Samples)

Location

Value

Location

Value

Most Instructions Executed

Location

Value

Location

Value

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 sections](#) to focus on selected performance aspects and make profiling faster.