

Result

577 - brute_force_AL_coarsening

Size

Time

Cycles

GPU

SM Frequency

Process

Attributes

Summary

Details

Source

Context

Comm

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 [%]	81.94	Duration [ms]	2.15
Memory Throughput [%]	2.01	Elapsed Cycles [cycle]	1,263,892
L1/TEX Cache Throughput [%]	1.56	SM Active Cycles [cycle]	1,257,481.05
L2 Cache Throughput [%]	0.67	SM Frequency [MHz]	584.98
DRAM Throughput [%]	2.01	DRAM Frequency [GHz]	5.00

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.

FP64/2 Utilization

Est. Speedup: 53.64%

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 [Expert Profiles Guide](#) for more details on roiline analysis.

GPU Throughput

Compute (SM) [%]

Memory [%]

Speed of Light (SOL) [%]

Compute Throughput Breakdown

Memory Throughput Breakdown

SM: Pipe Fp64 Cycles Active [%]	81.94	DRAM: Cycles Active [%]	2.01
SM: Pipe ALU Cycles Active [%]	43.51	DRAM: Dram Streams [%]	1.47
SM: Issue Active [%]	34.17	L1: M L1to2Dir Req Cycles Active [%]	0.78
SM: Inst Executed [%]	34.17	L2: T Sectors [%]	0.67
SM: Mio Pq Write Cycles Active [%]	26.92	L2: Xbar2Its Cycles Active [%]	0.67
SM: Inst Executed Pipe ALU [%]	19.71	L1: Lcam Requests [%]	0.39
EC: Request Cycles Active [%]	11.92	L2: D Sectors [%]	0.52
SM: Pipe Fma Cycles Active [%]	10.23	L1: Data Pipe Lsu Wavefronts [%]	0.28
SM: Mio Inst Issued [%]	10.12	L2: T Tag Requests [%]	0.17
SM: Mio Pq Read Cycles Active [%]	5.92	GPU: Compute Memory Access Throughput Internal Activity [%]	0.17
SM: Inst Executed Pipe Cpu Pced On Any [%]	4.60	L1: Data Bank Reads [%]	0.10
SM: Mio2W Writeback Active [%]	2.89	L1: Data Bank Writes [%]	0.10
SM: Inst Executed Pipe Lsu [%]	0.39	L1: Lsu Writeback Active [%]	0.06
SM: Inst Executed Pipe Xu [%]	0.13	L2: D Sectors Fill Device [%]	0.00
SM: Pipe Tensor Cycles Active [%]	0	L2: L2c2Dir Cycles Active [%]	0.00
SM: Inst Executed Pipe Fp16 [%]	0	L1: Team SndDir Req Cycles Active [%]	0.00
SM: Pipe Shared Cycles Active [%]	0	L1: F Wavefronts [%]	0.00
SM: Memory Throughput Internal Activity [%]	0	L1: M Xbar2Tltx Read Sectors [%]	0
SM: Instruction Throughput Internal Activity [%]	0	L1: Tex Writeback Active [%]	0
SM: Inst Executed Pipe Uniform [%]	0	L2: D Atomic Input Cycles Active [%]	0
SM: Inst Executed Pipe Tex [%]	0	L2: D Sectors Fill Sysmem [%]	0
SM: Inst Executed Pipe Ipa [%]	0	L1: Data Pipe Tex Wavefronts [%]	0
GPU: Compute Memory Request Throughput Internal Activity [%]	0		0

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 [cycles]: 40,000 # Pass Groups: 1

Maximum Buffer Size [bytes]: 3,06 Dropped Samples [sample]: 0

SM

0ms +0.2ms +0.4ms +0.6ms +0.8ms +1ms +1.2ms +1.4ms +1.6ms +1.8ms +2ms +2.2ms

Average Active Warps Per Cycle: 31.84 warp

Total Active Warps Per Cycle: 636.83 warp

SM Active Cycles: 27.4k cycle

Executed Ipc Active: 1.54 inst/cycle

DRAM

DRAM Throughput: 100 %

DRAM Read Bandwidth: 100 %

DRAM Write Bandwidth: 100 %

L1 Cache

Writeback Throughput: 20.63 cycle

Hit Rate: 100 %

Wavefronts (Data): 82.55

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 Ipc Elapsed [inst/cycle]	1.37	SM Busy [%]	82.34
Executed Ipc Active [inst/cycle]	1.37	Issue Slots Busy [%]	34.34
Issued Ipc Active [inst/cycle]	1.37		

Very High Utilization

FP64 is the highest-utilized pipeline (82.3%) 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 (82.3%) 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 [Memory Workload Analysis](#) section for more details on pipeline utilization. The [Performance Summary](#) section shows the mix of executed instructions in this kernel. Check the [Pipe Utilization](#) section for which 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 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 chart of the memory units. Detailed tables with data for each memory unit.

Memory Throughput [GB/s]	6.42	Mem Busy [%]	0.67
L1/TEX Hit Rate [%]	0	Max Bandwidth [%]	2.01
L2 Hit Rate [%]	99.29	Mem Pipes Busy [%]	25.92

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
Shared Load	0	0			0
Shared Load Matrix	0	0	0	0	0
Shared Store	0	0	0	0	0
Shared Atomic	0	0	0	0	0
Other	0	0	32,768	0.06	0
Total	0	0	32,768	0.06	0

L1/TEX Cache

	Instructions	Requests	Wavefronts	% Peak	Sectors	Sectors/Req	Hit Rate	Bytes	Sector Misses to L2	% Peak to L2	Returns to SM
Local Load	0	0	0	0	0	0	0	0	0	0	0
Global Load	0	0	0	0	0	0	0	0	0	0	0
Surface Load	0	0	0	0	0	0	0	0	0	0	0
Texture Load	0	0	0	0	0	0	0	0	0	0	0
Global Store	65,536	65,536	65,536	0.13	393,216	6	0	12,582,912	393,216	0.78	-
Local Store	0	0	0	0	0	0	0	0	0	0	-
Surface Store	0	0	0	0	0	0	0	0	0	0	-
Global Reduction	0	0	0	0	0	0	0	0	0	0	-
DGMEM Reduction	0	0	0	0	0	0	0	-	-	-	-
Surface Reduction	0	0	0	0	0	0	0	0	0	0	-
Global Atomic ALU	0	0	0	0	0	0	0	0	0	0	see above
Global Atomic CAS	0	0	0	0	0	0	0	0	0	0	see above
Surface Atomic ALU	0	0	0	0	0	0	0	0	0	0	-
Surface Atomic CAS	0	0	0	0	0	0	0	0	0	0	-
Loads	0	0	0	0	0	0	0	0	0	0	0
Stores	65,536	65,536	65,536	0.13	393,216	6	0	12,582,912	393,216	0.78	-
Atomics & Reductions	0	0	0	0	0	0	0	0	0	0	-
Total	65,536	65,536	65,536	0.13	393,216	6	0	12,582,912	393,216	0.78	-

L2 Cache

	Requests	Sectors	Sectors/Req	% Peak	Hit Rate	Bytes	Throughput	Sector Misses to Device	Sector Misses to System	Sector Misses to Peer
L1/TEX Load	0	0	0	0	0	0	0	0	0	0
L1/TEX Store	98,304	393,216	4	0.67	100	12,582,912	5,828,530,771.97	0	0	0
L1/TEX Atomic ALU	0	0	0	0	0	0	0	0	0	0
L1/TEX Atomic CAS	0	0	0	0	0	0	0	0	0	0
L1/TEX Reduction	0	0	0	0	0	0	0	0	0	0
L1/TEX Total	98,304	393,216	4	0.67	100	12,582,912	5,828,530,771.97	0	0	0
ECO Total	-	8	-	0.00	-	256	118,581.76	8	-	-
GPU Total	99,910	397,476	3.98	0.67	99.99	12,719,232	5,891,675,560.30	2,492	0	0

Device Memory

	Sectors	% Peak	Bytes	Throughput
Load	358	0.00	11,456	5,305,533.86
Store	433,030	2.01	13,856,960	6,418,682,556.62
Total	433,388	2.01	13,868,416	6,423,989,090.48

Scheduler Statistics

Summary of the activities 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]	3.95	No Eligible [%]	65.62
Eligible Warps Per Scheduler [warp]	0.48	One or More Eligible [%]	34.38
Issued Warp Per Scheduler	0.34		

Issue Slot Utilization

Est. Load Speedup: 18.06%

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 3.95 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 \(Per SM\)](#) sections can help too.

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]	11.50	Avg. Active Threads Per Warp	24.77
Warp Cycles Per Executed Instruction [cycle]	11.50	Avg. Not Predicted Off Threads Per Warp	22.48

Long Scoreboard Stalls

Est. Speedup: 18.06%

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, texture) 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 49.5% 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: 24.93%

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)

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]	69,091,328	Avg. Executed Instructions Per Scheduler [inst]	431,820.80
Issued Instructions [inst]	69,093,413	Avg. Issued Instructions Per Scheduler [inst]	431,833.83

FP32 Non-Fused Instructions

Est. Speedup: 5.14%

This kernel executes 0 fused and 491,520 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: 19.25%

This kernel executes 134,348 fused and 117,964 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

GPU and Memory Workload Distribution

Analysis of workload distribution in active cycles of SM, BMP, SMP, L1 & L2 caches, and DRAM

Average SM Active Cycles [cycle]	1,257,481.05	Average L1 Active Cycles [cycle]	1,257,481.05
Average L2 Active Cycles [cycle]	110,155.81	Average SMPSP Active Cycles [cycle]	1,256,013.51
Average DRAM Active Cycles [cycle]	216,694	Total SM Elapsed Cycles [cycle]	50,548,464
Total L1 Elapsed Cycles [cycle]	50,548,464	Total L2 Elapsed Cycles [cycle]	50,064,192
Total SM Elapsed Cycles [cycle]	202,193.55	Total DRAM Elapsed Cycles [cycle]	86,269,592

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 [%]	98.54
Branch Instructions Rate [%]	0.08	Avg. Divergent Branches	1,024

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