

Page: DetailsLaunch: 0 - 117 - computeAdd BaselineApply RulesSave as PDF

Current117 - compute (250, 1000, 1)x(64, 16, 1)Time: 197.72 msecondCycles: 284,616,594Regs: 27GPU: GeForce RTX 3080SM Frequency: 1.44 cycle/nsecondCC: 8.6Process: [14152] no DP.exe

GPU Speed Of Light

All

High-level overview of the utilization for compute and memory resources of the GPU. For each unit, the Speed Of Light (SOL) reports the achieved percentage of utilization with respect to the theoretical maximum. High-level overview of the utilization for compute and memory resources of the GPU presented as a roofline chart.

SOL SM [%]	94.58	Duration [msecond]	197.72
SOL Memory [%]	1.63	Elapsed Cycles [cycle]	284,616,594
SOL L1/TEX Cache [%]	1.65	SM Active Cycles [cycle]	283,752,197.91
SOL L2 Cache [%]	1.63	SM Frequency [cycle/nsecond]	1.44
SOL DRAM [%]	1.88	DRAM Frequency [cycle/nsecond]	9.24

GPU Utilization

SM [%]

Memory [%]

Speed Of Light [%]

SOL SM Breakdown

SOL Memory Breakdown

SOL SM: Issue Active [%]	94.58	SOL L2: T Sectors [%]	1.63
SOL SM: Inst Executed [%]	94.58	SOL L2: Xbar2lts Cycles Active [%]	1.52
SOL SM: Pipe Fma Cycles Active [%]	66.15	SOL GPU: Dram Throughput [%]	1.08
SOL SM: Pipe Fmaheavy Cycles Active [%]	57.14	SOL L1: Lsuin Requests [%]	0.91
SOL SM: Pipe Alu Cycles Active [%]	37.52	SOL L1: M L1tex2xbar Req Cycles Active [%]	0.83
SOL SM: Inst Executed Pipe Lsu [%]	0.91	SOL L1: Data Pipe Lsu Wavefronts [%]	0.55
SOL SM: Inst Executed Pipe Cbu Pred On Any [%]	0.67	SOL L2: D Sectors [%]	0.45
SOL SM: Mio Pq Read Cycles Active [%]	0.46	SOL L2: T Tag Requests [%]	0.41
SOL SM: Mio Pq Write Cycles Active [%]	0.46	SOL L1: Lsu Writeback Active [%]	0.25
SOL SM: Mio Inst Issued [%]	0.36	SOL L1: Data Bank Reads [%]	0.19
SOL SM: Inst Executed Pipe Adu [%]	0.33	SOL L2: Lts2xbar Cycles Active [%]	0.14
SOL SM: Inst Executed Pipe Xu [%]	0.33	SOL L2: D Sectors Fill Device [%]	0.14
SOL SM: Mio2rf Writeback Active [%]	0.17	SOL L1: Data Bank Writes [%]	0.10
SOL SM: Inst Executed Pipe Uniform [%]	0.02	SOL L1: Texin Sm2tex Req Cycles Active [%]	0.00
SOL IDC: Request Cycles Active [%]	0	SOL L1: M Xbar2l1tex Read Sectors [%]	0.00
SOL SM: Inst Executed Pipe lpa [%]	0	SOL L1: Data Pipe Tex Wavefronts [%]	0
SOL SM: Inst Executed Pipe Tex [%]	0	SOL L1: F Wavefronts [%]	0
SOL SM: Pipe Fp64 Cycles Active [%]	0	SOL L1: Tex Writeback Active [%]	0
SOL SM: Pipe Tensor Cycles Active [%]	0	SOL L2: D Atomic Input Cycles Active [%]	0
		SOL L2: D Sectors Fill Sysmem [%]	0

Floating Point Operations Roofline

Performance [FLOP/s]
(1 = 1e+12)

Arithmetic Intensity [FLOP/byte]

Recommendations

Bottleneck

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.

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 64:1. The kernel achieved 36% of this device's fp32 peak performance and 0% of its fp64 peak performance.

Compute Workload Analysis

All

Summary of the activity of the compute resources of the streaming multiprocessors (SM), including the achieved instructions per dock (IPC) and the utilization of each available pipeline. Pipelines with very high utilization might limit the overall performance.

Executed Ipc Elapsed [inst/cycle]	3.78	SM Busy [%]	94.87
Executed Ipc Active [inst/cycle]	3.79	Issue Slots Busy [%]	94.87
Issued Ipc Active [inst/cycle]	3.79		

Pipe Utilization

FMA

ALU

LSU

ADU

XU

CBU

Uniform

FMA (FP16)

FP64

TEX

Tensor (FP)

Tensor (INT)

Utilization [%]

Memory Workload Analysis

All

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 [gbyte/second]	8.81	Mem Busy [%]	1.63
L1/TEX Hit Rate [%]	88.87	Max Bandwidth [%]	1.52
L2 Hit Rate [%]	91.58	Mem Pipes Busy [%]	8.91
L2 Compression Success Rate [%]	8	L2 Compression Ratio	8

Scheduler Statistics

All

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]	7.78	No Eligible [%]	3.52
Eligible Warps Per Scheduler [warp]	4.88	One or More Eligible [%]	96.48
Issued Warp Per Scheduler	8.96		

Warp State Statistics

All

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]	8.86	Avg. Active Threads Per Warp	31.31
Warp Cycles Per Executed Instruction [cycle]	8.86	Avg. Not Predicated Off Threads Per Warp	28.57

Instruction Statistics

All

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]	73,217,888,243	Avg. Executed Instructions Per Scheduler [inst]	269,183,412.66
Issued Instructions [inst]	73,218,066,038	Avg. Issued Instructions Per Scheduler [inst]	269,184,066.32

Launch Statistics

All

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	250,000	Registers Per Thread [register/thread]	27
Block Size	1,024	Static Shared Memory Per Block [byte/block]	8
Threads [thread]	256,000,000	Dynamic Shared Memory Per Block [byte/block]	8
Waves Per SM	3,676.47	Driver Shared Memory Per Block [kbyte/block]	1.82
		Shared Memory Configuration Size [kbyte]	8.19

Occupancy

All

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 [%]	66.67	Block Limit Registers [block]	2
Theoretical Active Warps per SM [warp]	32	Block Limit Shared Mem [block]	100
Achieved Occupancy [%]	63.73	Block Limit Warps [block]	1
Achieved Active Warps Per SM [warp]	38.59	Block Limit SM [block]	16

Source Counters

All

Source metrics, including 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.