

Current

376

Time: 28.70 usecond

Cycles: 41,629

Reqs: 24

GPU: GeForce GTX 1660 Ti

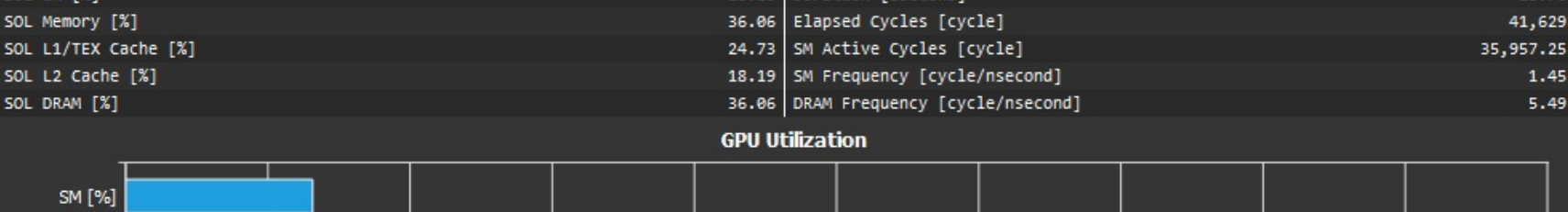
SM Frequency: 1.45 cycle/second

CC: 7.5

Process: [5972] memcached.exe

GPU Speed of Light

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.



Floating Point Operations Roofline



Recommendations

Bottleneck [Warning] 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 [Launch Statistics](#) and [Warp State Statistics](#) for potential reasons.

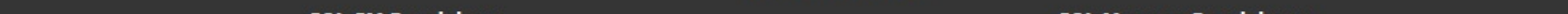
Roofline Analysis The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 0% of this device's fp32 peak performance and 0% of its fp64 peak performance.

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]	8.33	SM Busy [%]	9.54
Executed Ipc Active [inst/cycle]	8.38	Issue Slots Busy [%]	9.54
Issued Ipc Active [inst/cycle]	8.38		

Pipe Utilization



Recommendations

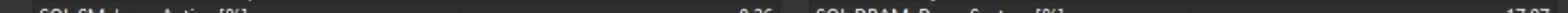
High Pipe Utilization [Warning] All pipelines are under-utilized. Either this kernel is very small or it doesn't issue enough warps per scheduler. Check the [Launch Statistics](#) and [Scheduler Statistics](#) sections for further details.

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 [gbyte/second]	95.18	Mem Busy [%]	18.19
L1/TEX Hit Rate [%]	40.18	Max Bandwidth [%]	36.86
L2 Hit Rate [%]	52.25	Mem Pipes Busy [%]	13.15

Memory Chart



Shared Memory

	Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	0	0	0	0	0
Shared Store	0	0	0	0	0
Shared Atomic	0	0	0	0	0
Other	-	-	8,448	0.85	0
Total	0	0	8,448	0.85	0

L1/TEX Cache

	Instructions	Requests	Wavefronts	% Peak	Sectors	Sectors/Req	Hit Rate
Local Load	0	0	0	0	0	0	0
Global Load	49,156	49,156	77,081	7.73	121,544	2,47	36.86
Surface Load	0	0	0	0	0	0	0
Texture Load	0	0	0	0	0	0	0
Global Store	8,192	8,192	8,192	0.82	8,192	1	87.5
Local Store	0	0	0	0	0	0	0
Surface Store	0	0	0	0	0	0	0
Global Reduction	0	0	0	0	0	0	0
Surface Reduction	0	0	0	0	0	0	0
Global Atomic ALU	0	0	0	0	0	0	0
Global Atomic CAS	0	0	0	0	0	0	0
Surface Atomic ALU	0	0	0	0	0	0	0
Surface Atomic CAS	0	0	0	0	0	0	0
Loads	49,156	49,156	77,081	7.73	121,544	2,47	36.86
Stores	8,192	8,192	8,192	0.82	8,192	1	87.5

L2 Cache

	Requests	Sectors	Sectors/Req	% Peak	Hit Rate	Bytes	Throug
L1/TEX Load	73,439	77,954	1.06	16.42	34.22	2,494,528	
L1/TEX Store	3,523	5,971	1.69	1.26	100	191,072	
L1/TEX Atomic ALU	0	0	0	0	0	0	
L1/TEX Atomic CAS	0	0	0	0	0	0	
L1/TEX Reduction	0	0	0	0	0	0	
L1/TEX Total	76,762	83,937	1.09	17.68	37.46	2,685,984	

Device Memory

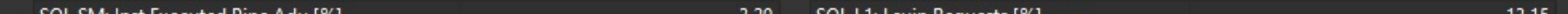
	Sectors	% Peak	Bytes	Throughput
Load	73,355	31.01	2,347,360	81,778,149,386.85
Store	11,952	5.05	382,464	13,324,414,715.72
Total	85,307	36.06	2,729,824	95,102,564,102.56

Scheduler Statistics

Summary of the activity of the schedulers issuing instructions. Each scheduler maintains a pool of warps that 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 issues 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]	7.48	No Eligible [%]	98.67
Eligible warps Per Scheduler [warp]	8.19	One or More Eligible [%]	9.33
Issued Warp Per Scheduler	8.89		

Warps Per Scheduler



Recommendations

Issue Slot Utilization [Warning] Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 10.7 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 7.48 active warps per scheduler, but only an average of 0.19 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 either increase the number of active warps or reduce the time the active warps are stalled.

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]	88.21	Avg. Active Threads Per Warp	32.88
Warp Cycles Per Executed Instruction [cycle]	88.58	Avg. Not Predicated Off Threads Per Warp	38.68

Warp State (All Cycles)



Recommendations

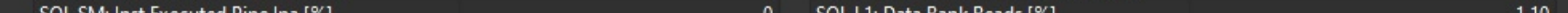
CPI Stall 'Long Scoreboard' [Warning] On average each warp of this kernel spends 71.0 cycles being stalled waiting for a scoreboard dependency on a L1TEX (local, global, surface, texture) operation. This represents about 88.6% of the total average of 80.2 cycles between issuing two instructions. To reduce the number of cycles waiting on L1TEX 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.

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]	327,696	Avg. Executed Instructions Per Scheduler [inst]	3,413.58
Issued Instructions [inst]	329,232	Avg. Issued Instructions Per Scheduler [inst]	3,429.58

Executed Instruction Mix



NVLink

High-level summary of NVLink connection status.

Physical Links	8	Logical Links	8
----------------	---	---------------	---

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.

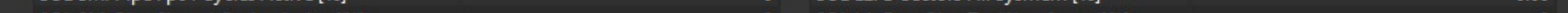
Grid Size	256	Registers Per Thread [register/thread]	24
Block Size	1,824	Static Shared Memory Per Block [byte/block]	8
Threads [thread]	262,144	Dynamic Shared Memory Per Block [byte/block]	8
Waves Per SM	18.67	Driver Shared Memory Per Block [byte/block]	8
Function Cache Configuration	cudaFuncCachePreferNone	Shared Memory Configuration Size [bbyte]	32.77

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 [%]	100	Block Limit Registers [block]	2
Theoretical Active Warps per SM [warp]	32	Block Limit Shared Mem [block]	16
Achieved Occupancy [%]	95.86	Block Limit Warps [block]	1
Achieved Active Warps per SM [warp]	38.42	Block Limit SM [block]	16

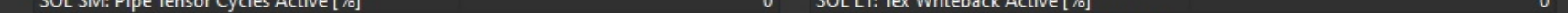
Impact of Varying Register Count Per Thread



Impact of Varying Block Size



Impact of Varying Shared Memory Usage Per Block

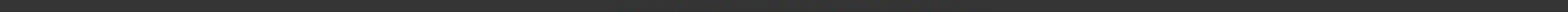


Source Counters

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]	8,194	Branch Efficiency [%]	8
Branch Instructions Ratio [%]	8.83	Avg. Divergent Branches	8

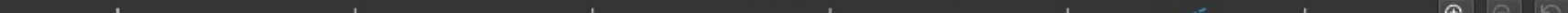
Sampling Data (All)



Sampling Data (Not Issued)



Most Instructions Executed



Recommendations

Uncoalesced Global Accesses [Warning] Uncoalesced global access, expected 8192 sectors, got 32766 (4.00x) at PC 0x700b90830