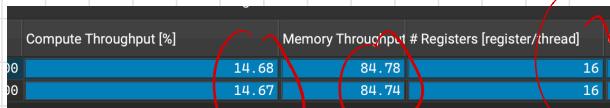


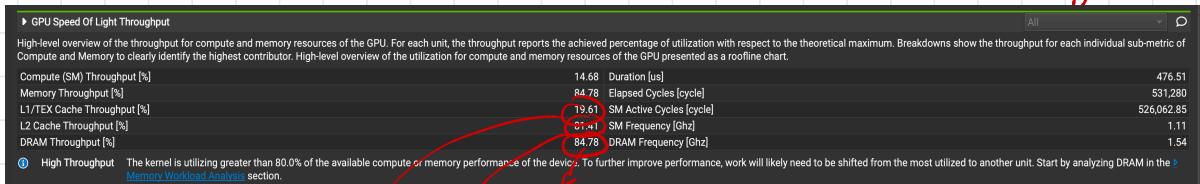
A100 profiling

VectorAdd



low reg pressure
" " high occupancy
" " more latency hiding
more latency hiding

(low util of compute units) Worse to map with of memory bus



maps
acts as pipeline
for all dram loads
has much lighter bandwidth load
than DRAM
so gets underutilized here

Compute Throughput Breakdown

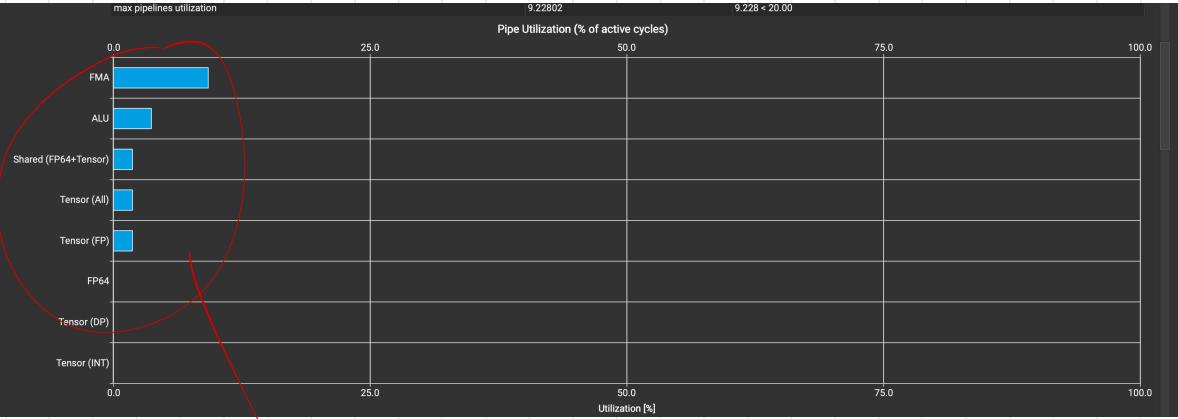
SM: Issue Active [%]	14.68
SM: Inst Executed [%]	14.67
SM: Inst Executed Pipe LSU [%]	9.17
SM: Mio Inst Issued [%]	9.17
SM: Pipe Fma Cycles Active [%]	7.34
SM: Mio2rf Writeback Active [%]	7.34
SM: Inst Executed Pipe Adu [%]	4.79
SM: Inst Executed Pipe Cbu Pred On Any [%]	3.67
SM: Mio Pq Write Cycles Active [%]	3.67
SM: Mio Pq Read Cycles Active [%]	1.83
SM: Pipe Alu Cycles Active [%]	1.83
SM: Pipe Tensor Cycles Active [%]	1.83
SM: Pipe Shared Cycles Active [%]	0
SM: Inst Executed Pipe Uniform [%]	0
SM: Inst Executed Pipe Tex [%]	0
SM: Inst Executed Pipe Ipa [%]	0
SM: Inst Executed Pipe Fp16 [%]	0
SM: Pipe Fp64 Cycles Active [%]	0
SM: Inst Executed Pipe Xu [%]	0
IDC: Request Cycles Active [%]	0

extremely fast at issuing load/store
 issued all quickly and then
 just waiting
 for memory
 to respond
 it hit its limit
 on how many
 in structure it
 will send very
 quickly

Memory Throughput Breakdown

DRAM: Cycles Active [%]	84.78
DRAM: Dram Sectors [%]	61.69
L2: D Sectors [%]	44.05
L2: D Sectors Fill Device [%]	41.63
L2: T Sectors [%]	41.32
L2: T Tag Requests [%]	31.10
L2: Xbar2lts Cycles Active [%]	29.26
L2: Lts2xbar Cycles Active [%]	20.83
L1: M L1tex2xbar Req Cycles Active [%]	19.49
L1: Data Pipe LSU Wavefronts [%]	16.75
L1: M Xbar2l1tex Read Sectors [%]	14.67
L1: Lsuin Requests [%]	14.67
L1: Lsu Writeback Active [%]	11.00
L1: Data Bank Reads [%]	5.50
L1: Data Bank Writes [%]	5.50
L1: F Wavefronts [%]	0.00
L1: Texin Sm2tex Req Cycles Active [%]	0.00
L2: D Atomic Input Cycles Active [%]	0
L2: D Sectors Fill Sysmem [%]	0
L1: Tex Writeback Active [%]	0
L1: Data Pipe Tex Wavefronts [%]	0

every transfer DRAM does, 60% full
 higher cap



✓ can clearly see 14-1. compute util here
 (9-1. FMA +)
 3-1. ALU)

Compute Workload Analysis

Detailed analysis of the compute resources of the streaming multiprocessors (SM), including the achieved instructions per clock (IPC) and

Executed Ipc Elapsed [inst/cycle]

Executed Ipc Active [inst/cycle]

Issued Ipc Active [inst/cycle]

⚠ Low Utilization All compute pipelines are under-utilized. Either this kernel is very small or it doesn't issue enough warps per scheduled warp slot(s).

Handwritten note: 4 (1 per warp slot(s))

0.59
0.59
0.59

) and the utilization of each available pipeline. Pipelines with very high utilization might limit the overall performance.

9 SM Busy [%] 14.77

9 Issue Slots Busy [%] 14.77

9

heudler. Check the [Launch Statistics](#) and [Scheduler Statistics](#) sections for further details.

► Key Performance Indicators

▼ 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 (Bandwidth), or by reaching the maximum throughput of issuing memory instructions (Mem Pipes Busy). Detailed chart of the memory unit

Memory Throughput [Tbyte/s]

L1/TEX Hit Rate [%]

L2 Hit Rate [%]

L2 Compression Success Rate [%]

Local Memory Usage Detects local memory usage and register spilling and estimates its impact on performance.

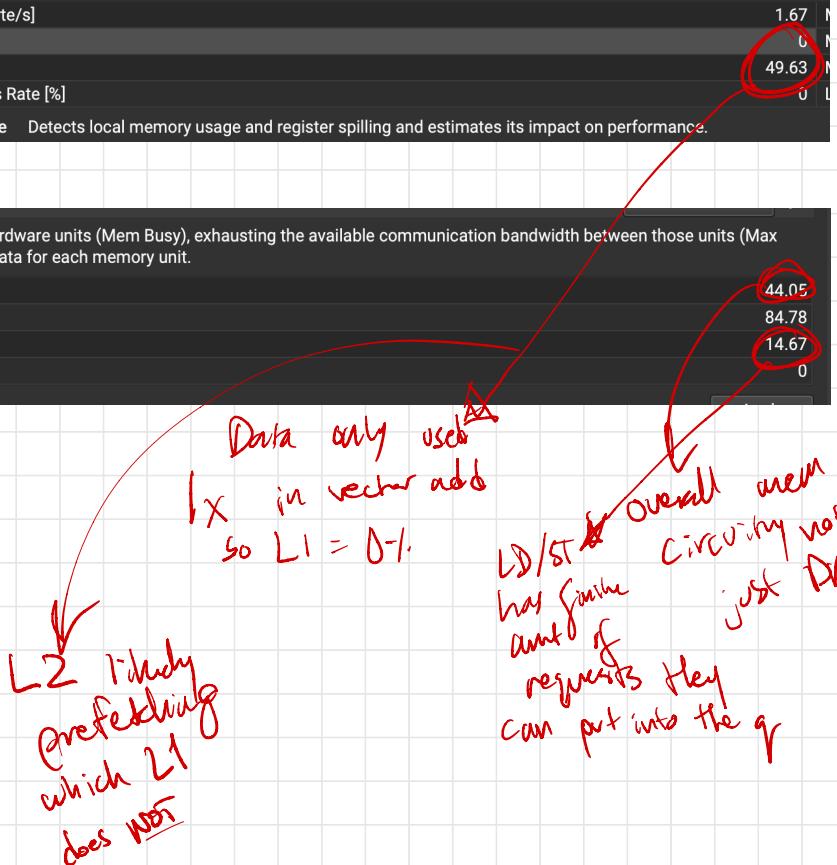
en fully utilizing the involved hardware units (Mem Busy), exhausting the available communication bandwidth between those units (Max Bandwidth), or by reaching the maximum throughput of issuing memory instructions (Mem Pipes Busy). Detailed tables with data for each memory unit.

.67 Mem Busy [%]

0 Max Bandwidth [%]

.63 Mem Pipes Busy [%]

0 L2 Compression Ratio



		L1/TEX Cache										
		Instructions	Requests	Wavefronts	% Peak	Sectors	Hit Rate	Bytes	Sector Misses to L2	% Peak to L2	Returns to SM	% Peak to SM
Local Load		0	4,194,304	0	4,194,304	0	16,777,216	4	0	536,870,912	14.67	6,291,456
Global Load To Shared Store (access)		0	0	4,194,304	7.34	0	0	0	0	0	-	-
Global Load To Shared Store (bypass)		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		2,097,152	2,097,152	3.67	8,388,608	4	0	268,495,456	0	8,388,608	14.67	-
Local Store		0	0	0	0	0	0	0	0	0	-	-
Surface Store		0	0	0	0	0	0	0	0	0	-	-
Global Reduction		0	0	0	0	0	0	0	0	0	-	-
Surface Reduction		0	0	0	0	0	0	0	0	0	see above	see above
Global Atomic ALU		0	0	0	0	0	0	0	0	0	0	0
Global Atomic CAS		0	0	0	0	0	0	0	0	0	0	0
Surface Atomic ALU		0	0	0	0	0	0	0	0	0	0	0
Surface Atomic CAS		0	4,194,304	4,194,304	7.34	16,777,216	4	0	536,870,912	14.67	6,291,456	11.00
Loads		4,194,304	2,097,152	2,097,152	3.67	8,388,608	4	0	268,495,456	8,388,608	14.67	-
Stores		0	0	0	0	0	0	0	0	0	0	-
Atomics & Handouts		0	6,291,456	6,291,456	11.00	25,165,824	4	0	805,306,368	25,165,824	29.34	6,291,456
Total		6,291,456	6,291,456	6,291,456	0	0	0	0	0	0	0	11.00

perfect
coherency

L1 to L2
bandwidth 301.
1:1, utilised
(STAT1)

Sum issues Show
L2 is write-back, success

	Requests	Sectors	Sectors/Req	% Peak	Hit Rate	Bytes	Throughput	Sector Misses to Device	Sector Misses to System	Sector Misses to Peer
L1/TEX Load	4,194,304	16,777,216	4	20.82	0	536,870,912	1,126,661,188,838,90	16,777,216	0	0
L1/TEX Store	2,097,152	8,388,608	4	13.88	100	208,435,456	563,334,094,419,45	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	6,291,456	25,165,824	4	31.23	33.33	805,306,468	1,690,002,283,258,34	16,777,216	0	0
GCC Total	0	0	0	0	0	0	0	0	0	0
ECC Total	-	3	-	0.00	-	96	201,443.97	3	-	-
L2 Fabric Total	6,229,877	24,777,236	3.98	61.39	66.13	792,871,552	1,663,906,789,335,84	8,401,984	0	0
GPU Total	12,529,594	49,555,691	3.99	41.32	49.63	1,588,582,112	3,334,757,303,068,97	25,161,179	0	0
	First	Hit Rate	Last	Hit Rate	Normal	Bytes	Throughput	L2 Cache Eviction Policies	Device Memory	
L1/TEX Load	0	0	0	0	0	16,777,216	0	Normal	0	Hit Rate 0
L1/TEX Store	0	0	0	0	0	8,388,608	100	Normal	0	Hit Rate 0
L1/TEX Atomic	0	0	0	0	0	25,165,824	0	Normal	0	Hit Rate -
L1/TEX Total	0	0	0	0	0	12,581,36	33.33	0	0	0
L2 Fabric Total	0	0	0	0	0	37,740,305	33.35	0	0	0
GPU Total	2,573	99.73	0	0	0	0	0	0	0	0
	Sectors	% Peak	Bytes	Throughput						
Load	16,777,228	57.57	536,871,296	1,126,661,934,644,78						
Store	8,083,400	27.57	258,668,800	542,837,955,812,24						
Total	24,860,628	84.78	795,540,996	1,669,506,950,507,92						

L2 much faster

loads more esp through stores

(plus vector add has 2x loads from stores)

10 just waiting

Scheduler Statistics

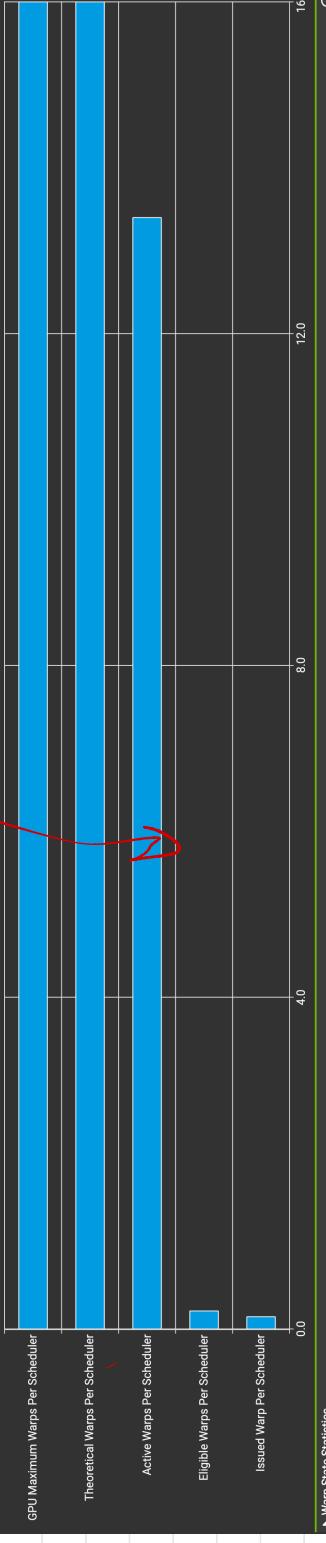
Summary of the activity of the scheduler issuing instructions. Each scheduler maintains a pool of warps in the pool (Active Warps). Active warps that are not stalled (Eligible Warps) are ready to issue their next instruction. From the set of eligible warps the scheduler selects a single warp from which to issue one or more instructions (Issued Warps). In practice, with no available 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

▲ Issue Slot Utilization: Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler can issue an instruction every 6.8 cycles. For this scheduler, for any given average of 0.22 warps per cycle, there are 1.6 warps per cycle. Every warp is 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 deadlocks, and to highly different execution durations per warp, Reducing stalls indicated on the [Warp State Statistics](#) and [Warp State Details](#).



Give in your lowest
Offering first and
I'll be glad to take it.

No divergence

