

Details

Launch0 - 505-MM_V2

Time

Cycles

Regs

GPU

SM Frequency

OC

Process

505-MM_V2 (728, 1, 1)(16, 16, 1)

4.55 msec/cond

5,999,993

37

NVIDIA GeForce RTX 3060

1.32 cycle/msec/cond

8.6

[22484] MM_V2

Save as PDF

Current

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 roofline chart.

Compute (SM) Throughput [%]	20.23	Duration [msec/cond]	4.55
Memory Throughput [%]	90.68	Elapsed Cycles [cycle]	5,999,993
Compute (SM) Throughput [%]	99.21	SM Active Cycles [cycle]	5,483,872.43
L1/TEX Cache Throughput [%]	5.69	SM Frequency [cycle/msec/cond]	1.32
L2 Cache Throughput [%]	5.74	DRAM Frequency [cycle/msec/cond]	7.29

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 [Memory Workload](#) section.

Roofline Analysis

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

GPU Throughput

Compute (SM) [%]

Memory [%]

Speed of Light (SOL) [%]

Compute Throughput Breakdown

Memory Throughput Breakdown

SM: Inst Executed Pipe Lsu [%]	20.23	L1: Data Pipe Lsu Wavefronts [%]	90.68
SM: Mto Inst Executed [%] <td>6.74<td>L1: Lsu Requests [%]<td>20.23</td></td></td>	6.74 <td>L1: Lsu Requests [%]<td>20.23</td></td>	L1: Lsu Requests [%] <td>20.23</td>	20.23
SM: Mto Inst Executed Active [%] <td>5.25<td>L1: Lsu Wavefronts Active [%]<td>10.74</td></td></td>	5.25 <td>L1: Lsu Wavefronts Active [%]<td>10.74</td></td>	L1: Lsu Wavefronts Active [%] <td>10.74</td>	10.74
SM: Inst Executed Active [%] <td>4.56<td>DRAM Cycles Active [%]<td>5.74</td></td></td>	4.56 <td>DRAM Cycles Active [%]<td>5.74</td></td>	DRAM Cycles Active [%] <td>5.74</td>	5.74
SM: Inst Executed [%] <td>4.56<td>L2: Xbar2/Its Cycles Active [%]<td>5.29</td></td></td>	4.56 <td>L2: Xbar2/Its Cycles Active [%]<td>5.29</td></td>	L2: Xbar2/Its Cycles Active [%] <td>5.29</td>	5.29
SM: Pipe Fma Cycles Active [%] <td>1.96<td>L2: T Sectors [%]<td>4.61</td></td></td>	1.96 <td>L2: T Sectors [%]<td>4.61</td></td>	L2: T Sectors [%] <td>4.61</td>	4.61
SM: Pipe ALU Cycles Active [%] <td>1.23<td>DRAM Drain Sectors [%]<td>1.77</td></td></td>	1.23 <td>DRAM Drain Sectors [%]<td>1.77</td></td>	DRAM Drain Sectors [%] <td>1.77</td>	1.77
SM: Pipe Fma/Arithmetic Cycles Active [%] <td>0.24<td>L1: M L1toXbar Req Cycles Active [%]<td>3.20</td></td></td>	0.24 <td>L1: M L1toXbar Req Cycles Active [%]<td>3.20</td></td>	L1: M L1toXbar Req Cycles Active [%] <td>3.20</td>	3.20
SM: Mto Pq Read Cycles Active [%] <td>0.08<td>L2: T Tag Requests [%]<td>2.84</td></td></td>	0.08 <td>L2: T Tag Requests [%]<td>2.84</td></td>	L2: T Tag Requests [%] <td>2.84</td>	2.84
SM: Mto Pq Write Cycles Active [%] <td>0.08<td>L2: L1toXbar Cycles Active [%]<td>1.94</td></td></td>	0.08 <td>L2: L1toXbar Cycles Active [%]<td>1.94</td></td>	L2: L1toXbar Cycles Active [%] <td>1.94</td>	1.94
SM: Inst Executed Pipe Uniform [%] <td>0.00<td>L2: D Sectors Fill Device [%]<td>0.20</td></td></td>	0.00 <td>L2: D Sectors Fill Device [%]<td>0.20</td></td>	L2: D Sectors Fill Device [%] <td>0.20</td>	0.20
SM: Inst Executed Pipe ALU [%] <td>0.00<td>L2: D Sectors [%]<td>1.76</td></td></td>	0.00 <td>L2: D Sectors [%]<td>1.76</td></td>	L2: D Sectors [%] <td>1.76</td>	1.76
SM: Inst Executed Pipe Cpu Pred On Any [%] <td>0.00<td>L1: Data Bank Reads [%]<td>1.48</td></td></td>	0.00 <td>L1: Data Bank Reads [%]<td>1.48</td></td>	L1: Data Bank Reads [%] <td>1.48</td>	1.48
SM: Inst Executed Pipe Lsu [%] <td>0<td>L1: M Xbar2/Its Read Sectors [%]<td>1.14</td></td></td>	0 <td>L1: M Xbar2/Its Read Sectors [%]<td>1.14</td></td>	L1: M Xbar2/Its Read Sectors [%] <td>1.14</td>	1.14
SM: Inst Executed Pipe Tex [%] <td>0<td>L1: Data Bank Writes [%]<td>0.20</td></td></td>	0 <td>L1: Data Bank Writes [%]<td>0.20</td></td>	L1: Data Bank Writes [%] <td>0.20</td>	0.20
SM: Inst Executed Pipe Xu [%] <td>0<td>L2: D Sectors Fill System [%]<td>0.00</td></td></td>	0 <td>L2: D Sectors Fill System [%]<td>0.00</td></td>	L2: D Sectors Fill System [%] <td>0.00</td>	0.00
IDC: Request Cycles Active [%] <td>0<td>L1: Tensin Sm2tex Req Cycles Active [%]<td>0.00</td></td></td>	0 <td>L1: Tensin Sm2tex Req Cycles Active [%]<td>0.00</td></td>	L1: Tensin Sm2tex Req Cycles Active [%] <td>0.00</td>	0.00
SM: Pipe Fp64 Cycles Active [%] <td>0<td>L1: F Wavefronts [%]<td>0</td></td></td>	0 <td>L1: F Wavefronts [%]<td>0</td></td>	L1: F Wavefronts [%] <td>0</td>	0
SM: Pipe Tensor Cycles Active [%] <td>0<td>L1: Tex Writeback Active [%]<td>0</td></td></td>	0 <td>L1: Tex Writeback Active [%]<td>0</td></td>	L1: Tex Writeback Active [%] <td>0</td>	0
		L2: D Atomic Input Cycles Active [%] <td>0</td>	0

Floating Point Operations Roofline

Performance [FLOP/s] (1+1e12)

Arithmetic Intensity [FLOP/byte]

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]	0.18	SM Busy [%]	7.38
Executed ipc Active [inst/cycle] <td>0.20<td>Issue Slots Busy [%]<td>4.99</td></td></td>	0.20 <td>Issue Slots Busy [%]<td>4.99</td></td>	Issue Slots Busy [%] <td>4.99</td>	4.99
Issued ipc Active [inst/cycle] <td>0.20<td></td><td></td></td>	0.20 <td></td> <td></td>		

Balanced

No pipeline is over-utilized.

Pipe Utilization

LSU

FMA

ALU

Uniform

ADU

CBU

FMA (FP16)

FP64

TEX

Tensor (FP)

Tensor (INT)

XU

Utilization [%]

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 unit.

Memory Throughput [GB/sec/cond]	20.08	Mem Busy [%]	90.68
L1/TEX Hit Rate [%] <td>96.44<th>Max Bandwidth [%]</th><td>20.23</td></td>	96.44 <th>Max Bandwidth [%]</th> <td>20.23</td>	Max Bandwidth [%]	20.23
L2 Hit Rate [%] <td>60.57<th>Mem Pipes Busy [%]</th><td>20.23</td></td>	60.57 <th>Mem Pipes Busy [%]</th> <td>20.23</td>	Mem Pipes Busy [%]	20.23
L2 Compression Success Rate [%] <td>0<th>L2 Compression Ratio</th><td>0</td></td>	0 <th>L2 Compression Ratio</th> <td>0</td>	L2 Compression Ratio	0

Memory Chart

Kernel

Global

Local

Texture

Surface

Load Global Store Shared

Shared

L1/TEX Cache

L2 Cache

System Memory

Device Memory

Peer Memory

L2 Compression

% Peak

Shared Memory

Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	0	0	0	0
Shared Load Matrix	0	0	0	0
Shared Store	0	0	0	0
Shared Store From Global Load	0	0	0	0
Shared Atomic	0	0	0	0
Other	-	-	2,048	0.00
Total	0	0	2,048	0.00

L1/TEX Cache

Instructions	Requests	Wavefronts	% Peak	Hit Rate	Sectors	Sectors/Req	Hit Rate	Bytes	Sector Misses to L2	% Peak to L2	Returns to SM
Local Load	65,536	65,536	65,536	0.04	261,924	4.00	0.81	8,381,568	1,915,924	1.14	18,04
Global Load	16,777,216	16,777,216	0	24.97	1,50,994,679	9.00	98.90	4,831,829,728	0	0	0
Global Load To Shared Store (access)	0	0	41,943,068	0	0	0	0	0	0	0	0
Global Load To Shared Store (bypass)	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	262,144	0.16	2,097,152	32	90.63	67,108,864	0	0	0
Local Store	81,920	131,072	131,072	0.08	524,200	4.00	44.08	16,774,400	2,621,440	1.56	0
Surface Store	0	0	0	0	0	0	0	0	0	0	0
Global Reduction	0	0	0	0	0	0	0	0	0	0	0
Surface Reduction	0	0	0	0	0	0	0	0	0	0	0
Global Atomic ALU	0	0	0	0	0	0	0	0	0	0	see a
Global Atomic CAS	0	0	0	0	0	0	0	0	0	0	see a
Surface Atomic ALU	0	0	0	0	0	0	0	0	0	0	0
Surface Atomic CAS	0	0	0	0	0	0	0	0	0	0	0
Loads	16,842,752	16,842,752	42,008,604	25.01	1,251,256,603	8.98	98.73	4,840,211,296	1,915,934	1.14	18,04
Stores	147,456	196,608	393,216	0.23	2,671,352	13.33	81.32	83,883,264	2,621,440	1.56	0

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 Store	1,719,437	1,915,543	1.11	1.87	8.24	61,297,376	13,471,303,395.50	1,760,778	0
L1/TEX Load	1,166,381	2,621,440	2.25	2.56	100	83,886,080	18,436,436,267.48	0	0
L1/TEX Atomic ALU	0	0	0	0	0	0	0	0	0
L1/TEX Atomic CAS	0	0	0	0	0	0	0	0	0
L1/TEX Reduction	0	0	0	0	0	0	0	0	0
L1/TEX Total	2,895,306	4,536,662	1.57	4.43	60.77	145,237,184	31,920,147,978.11	1,761,778	0
GPU Total	2,891,743	4,625,449	1.60	4.51	61.45	148,014,368	32,530,515,936.65	2,006,056	433

Device Memory

Sectors	% Peak	Bytes	Throughput
Load	1,978,136	3.98	63,300,352
Store	877,472	1.76	28,079,104
Total	2,855,608	5.74	91,379,456

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]	5.83	No Eligible [%]	95.01
Eligible Warps Per Scheduler [warp] <th>0.17</th> <th>One or More Eligible [%]</th> <th>4.99</th>	0.17	One or More Eligible [%]	4.99
Issued Warp Per Scheduler <th>0.05</th> <th></th> <th></th>	0.05		

Issue Slot Utilization

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 20.0 cycles. This might leave hardware resources underutilized and may lead to less optimal performance. Out of the maximum of 12 warps per scheduler, this kernel allocates an average of 5.83 active warps per scheduler, but only an average of 0.17 warps are 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, reduce the time the active warps are stalled by inspecting the top stall reasons on the [Warp State](#) and [Warp Stall](#) sections.

Warp State (All Cycles)

GPU Maximum Warps Per Scheduler

Theoretical Warps Per Scheduler

Active Warps Per Scheduler

Eligible Warps Per Scheduler

Issued Warp Per Scheduler

Cycles per Instruction

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 cycles per instruction the kernel spends in that 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]	116.76	Avg. Active Threads Per Warp	32
Warp Cycles Per Executed Instruction [cycle] <th>116.79</th> <th>Avg. Not Predicted Off Threads Per Warp</th> <th>31.93</th>	116.79	Avg. Not Predicted Off Threads Per Warp	31.93

Warp Stall

On average, each warp of this kernel spends 99.0 cycles being stalled waiting for the local/global instruction queue to be not full. This represents about 84.8% of the total average of 116.8 cycles between issuing two instructions. Typically this stall occurs only when executing local or global memory instructions extremely frequently. If applicable, consider combining multiple lower-level memory operations into fewer higher-level memory operations and try interleaving memory operations and math instructions.

Check the [Warp State](#) section for the top stall reasons in your source based on sampling data.

Warp State (All Cycles)

Stall L0 Throttle

Stall Long Scoreboard

Stall Not Selected

Stall Wait

Selected

Stall Branch Resolving

Stall No Instruction

Stall IMC Miss

Stall MIO Throttle

Stall Drain

Stall Dispatch Stall

Stall Math Pipe Throttle

Stall Misc

Stall Barrier

Stall Member

Stall Short Scoreboard

Stall Sleeping

Stall Tex Throttle

Cycles per Instruction

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]	30,641,152	Avg. Executed Instructions Per Scheduler [inst]	273,581.71
Issued Instructions [inst] <th>30,648,809</th> <th>Avg. Issued Instructions Per Scheduler [inst]</th> <th>273,650.08</th>	30,648,809	Avg. Issued Instructions Per Scheduler [inst]	273,650.08

Executed Instruction Mix

LDG

FFMA

IADD3

ISETP

BRA

IMAD

MOV

SHF

LDP3

STL

LEA

STO

LDL

UAD3D

ULDC

S2R

USHF

UMOV

EXT

Executed Instructions/Opcode

NVLink Topology

NVLink Topology

The system does not have any NVLink connections.

NVLink Tables

Detailed tables with properties for each NVLink.

Logical NVLink Properties

The system does not have any NVLink connections.

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	128	Registers Per Thread [register/thread]	37
Threads/Thread <th>256</th> <th>Static Shared Memory Per Block [byte/block]</th> <th>100</th>	256	Static Shared Memory Per Block [byte/block]	100
Waves Per SM <th>32,768</th> <th>Dynamic Shared Memory Per Block [byte/block]</th> <th>1.02</th>	32,768	Dynamic Shared Memory Per Block [byte/block]	1.02
Function Cache Configuration <th>0.76</th> <th>Driver Shared Memory Per Block [byte/block]</th> <th>8.19</th>	0.76	Driver Shared Memory Per Block [byte/block]	8.19

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]	6
Theoretical Active Warps per SM [warp] <th>48</th> <th>Block Limit Shared Mem [block]</th> <th>6</th>	48	Block Limit Shared Mem [block]	6
Achieved Occupancy [%] <th>48.57</th> <th>Block Limit Warps [block]</th> <th>16</th>	48.57	Block Limit Warps [block]	16
Achieved Active Warps per SM [warp] <th>23.31</th> <th>Block Limit SM [block]</th> <th></th>	23.31	Block Limit SM [block]	

Occupancy Limits

This kernel's theoretical occupancy is not impacted by any block limit. The difference between calculated theoretical (100.0%) and measured achieved occupancy (48.57%) can be the result of warp scheduling overheads or workload imbalances during the kernel execution. Load imbalances can occur between warps within a block as well as across blocks of the same kernel.

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

Source Counts

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]	590,848	Branch Efficiency [%]	100
Branch Instructions Ratio <th>0.02</th> <th>Avg. Divergent Branches</th> <th>0</th>	0.02	Avg. Divergent Branches	0

Uncoalesced Global Accesses

Uncoalesced global access, expected 2068229 sectors, got 8243993 (3.99%) at PC [0x7f7221af60](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 2068199 sectors, got 8243843 (3.99%) at PC [0x7f7221af60](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 2068133 sectors, got 8243513 (3.99%) at PC [0x7f7221af60](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 2068109 sectors, got 8243393 (3.99%) at PC [0x7f7221af60](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 2068082 sectors, got 8243258 (3.99%) at PC [0x7f7221af60](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 2067827 sectors, got 8241983 (3.99%) at PC [0x7f7221af60](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 2067794 sectors, got 8241818 (3.99%) at PC [0x7f7221af60](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 2067698 sectors, got 8241338 (3.99%) at PC [0x7f7221af60](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 2067668 sectors, got 8241188 (3.99%) at PC [0x7f7221af60](#)

Uncoalesced Global Accesses

Uncoalesced global access, expected 2067656 sectors, got 8241128 (3.99%) at PC [0x7f7221af60](#)

Sampling Data (All)

Location

Value

Value (%)

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

Sampling Data (Not Issued)

Location

Value

Value (%)

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

Most Instructions Executed

Location

Value

Value (%)

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8

0x7f7221af60 in MM_V2

8,290

8