

	Result	Size	Time	Cycles	GPU	SM Frequency	Process	Attributes
Current	578 mvt_kernel1	(149,1)x(704,1,1)	9.67 ms	16,215,639	0 - NVIDIA GeForce RTX 3050 Laptop GPU	1.68 GHz	[9319] mvt_KL_shared_loading_4_parallel_compute_warps.exe	
Summary	Details	Source	Context	Comments	Raw	Session		

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 rooftop chart.

Compute (SM) Throughput [%]	42.35	Duration [ms]	9.67
Memory Throughput [%]	42.35	Elapsed Cycles [cycle]	16215639
L1/TEX Cache Throughput [%]	42.94	SM Active Cycles [cycle]	15991646.80
L2 Cache Throughput [%]	13.43	SM Frequency [GHz]	1.68
DRAM Throughput [%]	29.00	DRAM Frequency [GHz]	11.99

Latency Issue This workload 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 [Pipeline Statistics](#) and [Warp State Statistics](#) for potential reasons.

Key Performance Indicators

Roofline Analysis The ratio of peak float (FP32) to double (FP64) performance on this device is 64:1. The workload achieved 2% of this device's FP32 peak performance and 0% of its FP64 peak performance. See the [Roofline Graph](#) for more details on rooftop analysis.

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.

Minimum Sampling Interval [us]	3	# Pass Groups	1
Maximum Buffer Size [Mbyte]	32.57		

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]	1.42	SM Busy [%]	36.08
Executed ipc: Active [Inst/cycle]	1.44	Issue Slots Busy [%]	36.08
Issued ipc: Active [Inst/cycle]	1.46		

Low Utilization Est. Local Speedup: 89.54% All compute pipelines are under-utilized. Either this workload is very small or it doesn't issue enough warps per scheduler. 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 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 [Obyte/s]	111.28	Mem Busy [%]	23.11
L1/TEX Hit Rate [%]	5.03	Max Bandwidth [%]	42.35
L2 Hit Rate [%]	21.44	Mem Pipes Busy [%]	42.35
L2 Compression Ratio	0	Local Memory Spilling Requests	0
L2 Compression Success Rate [%]	0	Local Memory Spilling Request Overhead [%]	0
	0	L2 Persisting Size [Mbyte]	6.29

Scheduler Statistics

Summary of the actions 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]	11.01	No Eligible [%]	63.42
Eligible Warps Per Scheduler [warp]	1.63	One or More Eligible [%]	36.58
Issued Warp Per Scheduler [warp]	0.37		

Issue Slot Utilization Est. Local Speedup: 57.65% Every scheduler is capable of issuing one instruction per cycle, but for this workload each scheduler only issues an instruction every 2.7 cycles. This might leave hardware resources underutilized and may lead to less optimal performance. Out of the maximum of 12 warps per scheduler, this workload 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 37.4% of the total average of 30.1 cycles between issuing two instructions.

Key Performance Indicators

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 warps are spent in system calls.

Executed Instructions [Inst]	45964942	Exec. Executed Instructions Per Scheduler [Inst]	5745618.65
Issued Instructions [Inst]	468038100	Issued Instructions Per Scheduler [Inst]	5850476.25

FP32 Non-Fused Instructions Est. Speedup: 2.43% This kernel executes 6388608 fused and 141943040 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 42% (relative to its current performance).

Key Performance Indicators

NVLink Topology

NVLink Topology diagram shows logical NVLink connections with transmit/receive throughput.

NVLink Tables

Detailed tables with properties for each NVLink.

Logical NVLink Properties	
The system does not have any NVLink connections.	

NUMA Affinity

Non-uniform memory access (NUMA) affinities based on compute and memory distances for all GPUs.

NUMA ID Table	
NUMA information is not available on the target system.	

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	1490	Function Cache Configuration	Cache/PreferNone
Cluster Size	0	Preferred Cluster Size	0
Registers Per Thread [register/thread]	35	Cluster Scheduling Policy	PolicySpread
Static Shared Memory Per Block [byte/block]	256	Block Size	704
Dynamic Shared Memory Per Block [byte/block]	0	Threads [thread]	1048960
Driver Shared Memory Per Block [byte/block]	1.02	Waves Per SM	37.25
Shared Memory Configuration Size [kbyte]	8.19	Uses Green Context	0
Stack Size	1024	# SMs [SM]	all
# TPCs	10	Enabled TPC IDs	20

World/C/LC Info				
World/C/LC Requests	Average	Min	Max	Sum
World/C/LC Requests Granted	0	0	0	0
World/C/LC Requests Granted as CTAs	0	0	0	0

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 [%]	91.67	Block Limit Registers [block]	2
Achieved Occupancy [%]	61	Block Limit Shared Mem [block]	6
Achieved Active Warps Per SM [warp]	43.96	Block Limit Warps [block]	2
Cluster Occupancy [%]	0	Block Limit Barriers [block]	24
Max Active Clusters [cluster]	0	Max Cluster Size [block]	8
Overall GPU Occupancy [%]	0		

GPU and Memory Workload Distribution

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

Average SM Active Cycles [cycle]	15991646.80	Average L1 Active Cycles [cycle]	15991646.80
Average L2 Active Cycles [cycle]	15992011.54	Average SMSP Active Cycles [cycle]	15992011.54
Average DRAM Active Cycles [cycle]	15991646.80	Total SM Elapsed Cycles [cycle]	324263540
Total L1 Elapsed Cycles [cycle]	15171094.56	Total L2 Elapsed Cycles [cycle]	250196784
Total SMSP Elapsed Cycles [cycle]	1297050160	Total DRAM Elapsed Cycles [cycle]	463917056

Workload Distribution				
SM Active Cycles	Average	Min	Max	Sum
SMSP Active Cycles	15991646.80	15805121	16158058	319823236
L1 Active Cycles	15992011.54	15736705	16149548	327966023
L2 Active Cycles	15991646.80	15805121	16158058	319823236
DRAM Active Cycles	15171094.56	15058463	15309116	242737513
	33636768	33636224	33637888	134547072

Source Counters

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 reasons. Only focus on stalls if the schedulers fail to issue every cycle.

Branch Instructions [Inst]	135791628	Branch Efficiency [%]	89.00
Branch Instructions Rate [%]	0.30	Avg Divergent Branches [branches]	104897.60

Warp Stall Sampling (All Samples)		Most Instructions Executed			
Location	Value	Value (%)	Location	Value	Value (%)
0x001160d0 in mvt_kernel1	57,837	37	0x00116090 in mvt_kernel1	16,780,288	6
0x001160a0 in mvt_kernel1	19,297	19	0x00116090 in mvt_kernel1	8,391,680	4
0x00116090 in mvt_kernel1	12,855	8	0x00116090 in mvt_kernel1	8,391,680	4
0x00116070 in mvt_kernel1	7,231	5	0x00116070 in mvt_kernel1	8,391,680	4
0x001160a0 in mvt_kernel1	4,427	3	0x00116090 in mvt_kernel1	8,391,680	4

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