

	Result	Size	Time	Cycles	GPU	SM Frequency	Process	Attributes
Current	553 - vectorSum	(781250, 1, 1)x(256, 1, 1)	10.06 ms	18,406,507	0 - NVIDIA GeForce RTX 4060	1.83 GHz	[1703] test_vectorSum	@

Summary	Details	Source	Context	Comments	Raw	Session
---------	---------	--------	---------	----------	-----	---------

The report contains imported source files.

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 routine chart.

	11.92	Duration [ms]	10.06
Compute (SM) Throughput [%]	94.33	Elapsed Cycles [cycle]	18406507
Memory Throughput [%]	14.50	SM Active Cycles [cycle]	18402698.29
L1/TEX Cache Throughput [%]	34.35	SM Frequency [GHz]	1.83
L2 Cache Throughput [%]	94.33	DRAM Frequency [GHz]	8.24

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 DRAM in the [Memory Workload Analysis](#) section.

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 64:1. The kernel achieved close to 0% of this device's fp32 peak performance and 0% of its fp64 peak performance. See the [Kernel Profiling Guide](#) for more details on roofline analysis.

Floating Point Operations Roofline

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.

Maximum Sampling Interval [µs] 4 # Pass Groups 2

Maximum Buffer Size [Mbytes] 7.94 Dropped Samples [sample] 67

Overview

Average Active Warps Per Cycle	87.80 warp	
Total Active Warps Per Cycle	0	
Blocks Launched	515 block	
SM Active Cycles	0	
Executed Ipc Active	357m inst/cycle	

SM

SM Throughput	100 %	
SM ALU Pipe Throughput	0	
SM FMA Light Pipe Throughput	100 %	
SM FMA Heavy Pipe Throughput	0	
SM Tensor Pipe Throughput	100 %	

DRAM

DRAM Throughput	100 %	
DRAM Read Bandwidth	100 %	
DRAM Write Bandwidth	0	

L2 Cache

L2 Throughput	100 %	
L2 Hit Rate	100 %	

L1 Cache

L1 Throughput	100 %	
Writeback Throughput	511.54 cycle	
Hit Rate	100 %	
Wavefronts (Data)	739.71	
Workload Execution		vectorSum

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 (Peak Inst/cycle)	0.23	SM Busy [%]	5.66
Executed Ipc Active [Inst/cycle]	0.23	Issue Slots Busy [%]	5.66
Issued Ipc Active [Inst/cycle]	0.23		

Low Utilization

Est. Local Speedup: 97.17%

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.

Pipe Utilization (% of active cycles)

Pipe Utilization (% of peak instructions executed)

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/s]	248.81	Mem Busy [%]	34.35
L1/TEX Hit Rate [%]	0	Max Bandwidth [%]	94.33
L2 Hit Rate [%]	31.85	Mem Pipes Busy [%]	11.32
L2 Compression Success Rate [%]	0	L2 Compression Ratio	0

Memory Chart

Show As: Transfer Size

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]	10.23	No Eligible [%]	94.33
Eligible Warps Per Scheduler [warp]	0.06	One or More Eligible [%]	5.67
Issued Warp Per Scheduler	0.06		

Issue Slot Utilization

Est. Local Speedup: 5.67%

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 17.6 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 10.23 active warps per scheduler, but only an average of 0.06 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, avoid possible load imbalances due to highly different execution durations per warp. Reducing stalls indicated on the [Warp Stall Sampling](#) and [Warp State](#) sections can help, too.

Warps Per Scheduler

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]	180.38	Avg. Active Threads Per Warp	32
Warp Cycles Per Executed Instruction [cycle]	180.39	Avg. Not Predicted Off Threads Per Warp	30

Long Scoreboard Stalls

Est. Speedup: 5.67%

On average, each warp of this kernel spends 172.5 cycles being stalled waiting for a scoreboard dependency on a L1/TEX (local, global, surface, texture) operation. Find the instruction producing the data being waited upon to identify the culprit. To reduce the number of cycles waiting on L1/TEX data 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 95.6% of the total average of 180.4 cycles between issuing two instructions.

Warp Stall

Check the [Warp Stall Sampling \(All Samples\)](#) table for the top stall locations in your source based on sampling data. The [Kernel Profiling Guide](#) provides more details on each stall reason.

Warp State (All Cycles)

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/Opcodes' and 'Executed Instructions' are measured differently and can diverge if cycles are spent in system calls.

Executed Instructions [Inst]	100000000	Avg. Executed Instructions Per Scheduler [Inst]	1041666.67
Theoretical Occupancy [%]	100	Block Limit Registers [block]	1041717.67
Executed Instructions [Inst]	100004896	Avg. Issued Instructions Per Scheduler [Inst]	

FP32 Non-Fused Instructions

Est. Speedup: 1.42%

This kernel executes 0 fused and 6250000 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 50% (relative to its current performance). Check the [Source](#) page to identify where this kernel executes FP32 instructions.

Executed Instruction Mix

NVLink Topology

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

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.

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	781250	Function Cache Configuration	CachePreferNone
Registers Per Thread [register/thread]	16	Static Shared Memory Per Block [byte/block]	0
Block Size	256	Dynamic Shared Memory Per Block [byte/block]	0
Theoretical Occupancy [%]	78205020	Driver Shared Memory Per Block [byte/block]	1.02
Waves Per SM	5425.35	Shared Memory Configuration Size [kbyte]	16.38
Uses Green Context	0	# SMs [SM]	24

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 typically indicates highly imbalanced workloads.

Achieved Occupancy [%]	100	Block Limit Registers [block]	16
Theoretical Active Warps per SM [warp]	48	Block Limit Shared Mem [block]	16
Achieved Occupancy [%]	85.56	Block Limit Warps [block]	6
Achieved Active Warps Per SM [warp]	41.07	Block Limit SM [block]	24

Achieved Occupancy

Est. Speedup: 5.67%

The difference between calculated theoretical (100.0%) and measured achieved occupancy (85.6%) 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. See the [CUDA Best Practices Guide](#) for more details on optimizing occupancy.

Impact of Varying Register Count Per Thread

Impact of Varying Block Size

Impact of Varying Shared Memory Usage Per Block

GPU and Memory Workload Distribution

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

Average SM Active Cycles [cycle]	18402698.29	Average L1 Active Cycles [cycle]	18402698.29
Average L2 Active Cycles [cycle]	16603652.33	Average SMSP Active Cycles [cycle]	18375450.41
Average DRAM Active Cycles [cycle]	78205020	Total SM Elapsed Cycles [cycle]	441756128
Total L1 Elapsed Cycles [cycle]	441756128	Total L2 Elapsed Cycles [cycle]	1767024512
Total SMSP Elapsed Cycles [cycle]	1767024512	Total DRAM Elapsed Cycles [cycle]	331628544

Workload Distribution

	Average	Min	Max	Sum
SM Active Cycles	18402698.29	18401343	18404434	441664759
SMSP Active Cycles	18375450.41	18372279	1764043229	441644759
L1 Active Cycles	18402698.29	18401343	18404434	441664759
L2 Active Cycles	16603652.33	16603213	16604392	199243828
DRAM Active Cycles	78205020	78188928	78223344	312820080

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

Branch Instructions [Inst]	12500000	Branch Efficiency [%]	0
Branch Instructions Ratio [%]	0.12	Warp Divergent Branches	0

Warp Stall Sampling (All Samples)

Location	Value	Value (%)	Location	Value	Value (%)
vectorSum.cu:18 (0x500e72240 in vectorSum)	4117038	59	vectorSum.cu:18 (0x500e722f0 in vectorSum)	6250000	8
vectorSum.cu:18 (0x500e72240 in vectorSum)	5866	1	vectorSum.cu:16 (0x500e72260 in vectorSum)	6250000	8
vectorSum.cu:14 (0x500e72240 in vectorSum)	3634	1	vectorSum.cu:16 (0x500e72240 in vectorSum)	6250000	8
vectorSum.cu:15 (0x500e72280 in vectorSum)	1851	0	vectorSum.cu:16 (0x500e722d0 in vectorSum)	6250000	8
vectorSum.cu:16 (0x500e72290 in vectorSum)	1268	0	vectorSum.cu:18 (0x500e72240 in vectorSum)	6250000	8

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.