

Result

Size

Time

Cycles

GPU

SM Frequency

Process

Attributes

Current

6507 - square\_kernel

(512, 1, 1)x(128, 1, 1)

14.24 us

8,283

0 - Tesla T4

581.14 Mhz

[25497] python3.11

Summary

Details

Source

Context

Comments

Raw

Session

Compare

Tools

View

Export

GPU Speed Of Light Throughput

GPU Throughput Chart

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 [%]	7.78	Duration [us]	14.24
Memory Throughput [%]	93.10	Elapsed Cycles [cycle]	8283
L1/TEX Cache Throughput [%]	55.95	SM Active Cycles [cycle]	5111.12
L2 Cache Throughput [%]	34.22	SM Frequency [Mhz]	581.14
DRAM Throughput [%]	93.10	DRAM Frequency [Ghz]	4.03

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 32:1. The kernel achieved 3% 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.

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 [cycle]	20000	# Pass Groups	1
Maximum Buffer Size [Kbytes]	64	Dropped Samples [sample]	0

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]	0.29	SM Busy [%]	10.10
Executed Ipc Active [inst/cycle]	0.37	Issue Slots Busy [%]	10.10
Issued Ipc Active [inst/cycle]	0.40		

Low Utilization

Est. Local Speedup: 93.49%

All compute 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

Memory Chart

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]	239.94	Mem Busy [%]	34.22
L1/TEX Hit Rate [%]	0	Max Bandwidth [%]	93.10
L2 Hit Rate [%]	50.47	Mem Pipes Busy [%]	7.71

Memory L2 Compression

The optional metric lts\_\_average\_gcomp\_input\_sector\_success\_rate.pct could not be found. Collecting it as an additional metric could enable the rule to provide more guidance.

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]	5.99	No Eligible [%]	89.81
Eligible Warps Per Scheduler [warp]	0.18	One or More Eligible [%]	10.19
Issued Warp Per Scheduler	0.10		

Issue Slot Utilization

Est. Local Speedup: 6.90%

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 9.8 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 5.99 active warps per scheduler, but only an average of 0.18 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, reduce the time the active warps are stalled by inspecting the top stall reasons on the [Warp State Statistics](#) and [Source Counters](#) sections.

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]	58.85	Avg. Active Threads Per Warp	32
Warp Cycles Per Executed Instruction [cycle]	64.17	Avg. Not Predicated Off Threads Per Warp	32

Long Scoreboard Stalls

Est. Speedup: 6.90%

On average, each warp of this kernel spends 47.7 cycles being stalled waiting for a scoreboard dependency on a L1TEX (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 L1TEX 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 81.0% of the total average of 58.9 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.

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]	75776	Avg. Executed Instructions Per Scheduler [inst]	473.60
Issued Instructions [inst]	82625	Avg. Issued Instructions Per Scheduler [inst]	516.41

NVLink Topology

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

NVLink Tables

Detailed tables with properties for each NVLink.

NUMA Affinity

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

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	512	Function Cache Configuration	CachePreferNone
Registers Per Thread [register/thread]	18	Static Shared Memory Per Block [byte/block]	0
Block Size	128	Dynamic Shared Memory Per Block [byte/block]	0
Threads [thread]	65536	Driver Shared Memory Per Block [byte/block]	0
Waves Per SM	1.60	Shared Memory Configuration Size [Kbyte]	32.77
Uses Green Context	0	# SMs [SM]	40

Tail Effect

Est. Speedup: 50.00%

A wave of thread blocks is defined as the maximum number of blocks that can be executed in parallel on the target GPU. The number of blocks in a wave depends on the number of multiprocessors and the theoretical occupancy of the kernel. This kernel launch results in 1 full waves and a partial wave of 192 thread blocks. Under the assumption of a uniform execution duration of all thread blocks, the partial wave may account for up to 50.0% of the total kernel runtime with a lower occupancy of 24.3%. Try launching a grid with no partial wave. The overall impact of this tail effect also lessens with the number of full waves executed for a grid. See the [Hardware Model](#) description for more details on launch configurations.

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]	21
Theoretical Active Warps per SM [warp]	32	Block Limit Shared Mem [block]	16
Achieved Occupancy [%]	75.66	Block Limit Warps [block]	8
Achieved Active Warps Per SM [warp]	24.21	Block Limit SM [block]	16

Achieved Occupancy

Est. Speedup: 6.90%

The difference between calculated theoretical (100.0%) and measured achieved occupancy (75.7%) 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.

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]	5111.12	Average L1 Active Cycles [cycle]	5111.12
Average L2 Active Cycles [cycle]	7318.31	Average SMSP Active Cycles [cycle]	5070
Average DRAM Active Cycles [cycle]	53386.50	Total SM Elapsed Cycles [cycle]	265536
Total L1 Elapsed Cycles [cycle]	265536	Total L2 Elapsed Cycles [cycle]	387232
Total SMSP Elapsed Cycles [cycle]	1062144	Total DRAM Elapsed Cycles [cycle]	458752

SMSPs Workload Imbalance

Est. Speedup: 5.10%

One or more SMSPs have a much higher number of active cycles than the average number of active cycles. Additionally, other SMSPs have a much lower number of active cycles than the average number of active cycles. Maximum instance value is 6.68% above the average, while the minimum instance value is 8.22% below the average.

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]	2048	Branch Efficiency [%]	0
Branch Instructions Ratio [%]	0.03	Avg. Divergent Branches	0

To customize your report even further, you might want to learn about [custom sections](#) and [writing your own rules](#). You might also want to consider [adding individual metrics](#).