

Result **Size** **Time** **Cycles** **GPU** **SM Frequency** **Process** **Attributes**

Current 563 - reduce_argmin (1024, 1, 1)(1024, 1, 1) 123.62 us 72,301 0 - Tesla T4 584.83 MHz [10345] exe

Primary **Details** **Source** **Context** **Comments** **Raw** **Session** **Compare** **Tools** **View** **Export**

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.

Metric	Value	Duration [s]
Compute (SM) Throughput [%]	24.41	123.62
Memory Throughput [%]	40.42	72.201
L1/TEX Cache Throughput [%]	33.50	66,030.25
L2 Cache Throughput [%]	11.73	584.83
DRAM Throughput [%]	40.42	5.00

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

Roofline Analysis The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 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.

GPU Throughput

Compute (SM) [%] Memory [%]

Speed of Light (SOL) [%]

Compute Throughput Breakdown

Metric	Value
SM: Inst Executed Pipe Lsu [%]	24.41
SM: Issue Active [%]	10.59
SM: Inst Executed [%]	10.63
SM: Memory Inst Executed [%]	9.76
SM: Pipe Alu Cycles Active [%]	6.22
SM: Mio Pq Write Cycles Active [%]	5.42
SM: Mio Pq Read Cycles Active [%]	5.35
SM: Inst Executed Pipe Adu [%]	4.83
SM: Inst Executed Pipe Fma [%]	4.55
SM: Inst Executed Pipe Tex [%]	3.87
SM: Inst Executed Pipe Cbu Pried On Any [%]	2.59
SM: Inst Executed Pipe Uniform [%]	1.14
SM: Memory Throughput Internal Activity [%]	0
SM: Pipe Tensor Cycles Active [%]	0
SM: Pipe Shared Cycles Active [%]	0
SM: Pipe Fp64 Cycles Active [%]	0
SM: Inst Executed Pipe Adu [%]	0
SM: Inst Executed Pipe Fma [%]	0
SM: Inst Executed Pipe Tex [%]	0
SM: Inst Executed Pipe Lsu [%]	0
SM: Inst Executed Pipe Fp16 [%]	0

Memory Throughput Breakdown

Metric	Value
DRAM: Cycles Active [%]	40.42
DRAM: Drim Sectors [%]	29.58
L1: Lsu Requests [%]	24.41
L1: Data Pipe Lsu Wavefronts [%]	16.78
L1: M Xbar-Write Read Sectors [%]	13.67
L2: T Sectors [%]	11.73
L2: Lsu2bar Cycles Active [%]	11.70
L2: D Sectors Fill Device [%]	11.65
L1: Lsu Writeback Active [%]	7.85
L2: D Sectors [%]	6.10
L1: Data Bank Writes [%]	3.79
L1: Data Bank Reads [%]	3.09
GPU Compute Memory Access Throughput Internal Activity [%]	3.00
L2: T Tag Requests [%]	1.93
L1: M L1tex2bar Req Cycles Active [%]	1.89
L2: Xbar2bar Cycles Active [%]	0.36
L1: Team Sm2bar Req Cycles Active [%]	0.00
L1: F Wavefronts [%]	0
L2: D Atomic Input Cycles Active [%]	0
L1: Tex Writeback Active [%]	0
L2: D Sectors Fill System [%]	0
L1: Data Pipe Tex Wavefronts [%]	0
GPU Compute Memory Request Throughput Internal Activity [%]	0

Floating Point Operations Roofline

Performance (FLOP/s) (1e1000000000)

Arithmetic Intensity (FLOP/byte)

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.

Metric	Value
Maximum Sampling Interval [cycle]	20,000
Maximum Buffer Size [kbytes]	512

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.

Metric	Value
Executed Ipc Elapsed [inst/cycle]	0.42
Executed Ipc Active [inst/cycle]	0.46
Issued Ipc Active [inst/cycle]	0.46

Low Utilization Est. Local Speedup: 93.21% 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.

Pipe Utilization (% of active cycles)

Pipe	Utilization [%]
ALU	~10
FMA	~10
FP64	~10
Shared (FP16+Tensor)	~10
Tensor (All)	~10
Tensor (FP)	~10
Tensor (INT)	~10

Pipe Utilization (% of peak instructions executed)

Pipe	Utilization [%]
LSU	~25
ALU	~10
ADU	~10
FMA	~10
Uniform	~10
CBU	~10
FP16	~10
FP64	~10
TEX	~10
Tensor (FP)	~10
Tensor (INT)	~10
XU	~10

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 (Mem Busy), 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.

Metric	Value
Memory Throughput [Gbyte/s]	129.44
L1/TEX Hit Rate [%]	0.15
L2 Hit Rate [%]	0.92

Memory L2 Compression The optional metric l2s_avg_compute_input_sector_success_rate_pct could not be found. Collecting it as an additional metric could enable the rule to provide more guidance.

DRAM Global Load Access Pattern Est. Speedup: 0.99% The memory access pattern for global loads from DRAM might not be optimal. On average, only 31.9 of the 32 bytes transmitted per sector are utilized by each thread. This applies to the 99.7% of sectors missed in L2. This could possibly be caused by a stride between threads. Check the [Global Load Access](#) section for uncoalesced global loads.

L1TEX Global Store Access Pattern Est. Speedup: 27.22% The memory access pattern for global stores to L1TEX might not be optimal. On average, only 6.0 of the 32 bytes transmitted per sector are utilized by each thread. This could possibly be caused by a stride between threads. Check the [Global Store Access](#) section for uncoalesced global stores.

Shared Load Bank Conflicts Est. Speedup: 4.27% The memory access pattern for shared loads might not be optimal and causes on average a 1.2 - way bank conflict across all 94060 shared load requests. This results in 14219 bank conflicts, which represent 12.75% of the overall 111509 wavefronts for shared loads. Check the [Shared Load Access](#) section for uncoalesced shared loads.

Memory Chart

Kernel

Global

Local

Texture

Surface

Shared

L1/TEX Cache

L2 Cache

System Memory

Device Memory

Peer Memory

Shared Memory

Metric	Value
Shared Load	94,060
Shared Store	102,992
Shared Atomic	0
Other	0
Total	197,052

L1/TEX Cache

Metric	Value
Global Load	66,560
Surface Load	0
Texture Load	0
Global Store	1,160
Local Store	0
Surface Store	0
Global Reduction	0
DSMEM Reduction	0
Surface Reduction	0
Global Atomic CAS	0
Global Atomic CAS	0
Surface Atomic CAS	0
Surface Atomic CAS	0
Loads	66,560
Stores	1,160
Atomic & Reductions	0
Total	67,720

L2 Cache

Metric	Value
Global Load	66,560
Surface Load	0
Texture Load	0
Global Store	1,160