

Result558 - brute\_forceSize(1024, 1, 1)(1024, 1, 1)1.61 ms943,4960 - Tests T4584.98 MHz162/74 exe @Attributes

SummaryDetailsSourceContextCommentsRawSession

GPU Speed of Light ThroughputAll

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 [%]	70.71	Duration [ms]	1.61
SM: Inst Executed [%]	70.71	L1: Data Pipe Lsu Wavefronts [%]	35.22
Memory Throughput [%]	57.53	Elapsed Cycles [cycle]	945,496
L1/TEX Cache Throughput [%]	70.52	SM Active Cycles [cycle]	919,148.23
L2 Cache Throughput [%]	19.32	SM Frequency [MHz]	584.98
DRAM Throughput [%]	3.82	DRAM Frequency [GHz]	5.01

High Compute Throughput

Compute is more heavily utilized than Memory. Look at the [Compute Workload Analysis](#) section to see what the compute pipelines are spending their time doing. Also, consider whether any computation is redundant and could be reduced or moved to look-up tables.

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 10% of this device's fp32 peak performance and 0% of its fp64 peak performance. See the [Roofline Profiling Guide](#) for more details on roofline analysis.

GPU Throughput

Compute (SM) [%]Memory [%]

Speed of Light (SOL) [%]

Compute Throughput BreakdownMemory Throughput Breakdown

SM: Issue Active [%]	70.71	L1: Lsu Requests [%]	57.53
SM: Inst Executed [%]	70.71	L1: Data Pipe Lsu Wavefronts [%]	35.22
SM: Inst Executed Pipe Lsu [%]	70.71	L1: Lsu Writeback Active [%]	22.91
SM: Pipe Alu Cycles Active [%]	55.97	L1: M1Tex2ubar Req Cycles Active [%]	22.70
SM: Inst Executed Pipe Adu [%]	42.70	L2: Xbar2iz Cycles Active [%]	19.32
IDC: Request Cycles Active [%]	41.71	L2: T Sectors [%]	16.44
SM: Pipe Tma Cycles Active [%]	37.89	L1: Data Bank Reads [%]	13.59
SM: Inst Executed Pipe Xu [%]	35.93	L2: D Sectors [%]	4.40
SM: Mio Inst Issued [%]	33.41	L2: T Tag Requests [%]	4.34
SM: Mio Pq Write Cycles Active [%]	26.03	GPU: Compute Memory Access Throughput Internal Activity [%]	4.34
SM: Mio Pq Read Cycles Active [%]	26.03	DRAM: Cycles Active [%]	3.82
SM: Minor Writeback Active [%]	21.92	DRAM: Dram Sectors [%]	2.80
SM: Inst Executed Pipe Cbu Pred On Any [%]	6.05	L1: Data Bank Writes [%]	2.65
SM: Inst Executed Pipe Uniform [%]	0.04	L1: M Xbar2l1tex Read Sectors [%]	0.98
SM: Pipe Tensor Cycles Active [%]	0	L2: L1s2ubar Cycles Active [%]	0.87
SM: Pipe Shared Cycles Active [%]	0	L2: D Sectors Fd device [%]	0.86
SM: Pipe Fp64 Cycles Active [%]	0	L1: Twin Sm2xor Req Cycles Active [%]	0.00
SM: Memory Throughput Internal Activity [%]	0	L1: F Wavefronts [%]	0.00
SM: Instruction Throughput Internal Activity [%]	0	L1: Tex Writeback Active [%]	0
SM: Inst Executed Pipe Tex [%]	0	L2: D Atomic Input Cycles Active [%]	0
SM: Inst Executed Pipe Isa [%]	0	GPU: Compute Memory Request Throughput Internal Activity [%]	0
SM: Inst Executed Pipe Fp16 [%]	0	L2: D Sectors Fill System [%]	0
		L1: Data Pipe Tex Wavefronts [%]	0

Floating Point Operations Roofline

Performance [FLOP/s] (1 = 10,000,000,000)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.

Maximum Sampling Interval [cycle]160,000# Pass Groups1

Maximum Buffer Size [Mbytes]2.56Dropped Samples [sample]0

SM

Average Active Warps Per Cycle59.20 warp

Total Active Warps Per Cycle1.18k warp

SM Active Cycles109k cycle

Executed Ipc Active2.99 inst/cycle

DRAM

DRAM Throughput100 %

DRAM Read Bandwidth100 %

DRAM Write Bandwidth100 %

L1 Cache

Writeback Throughput26.4k cycle

Hit Rate100 %

Wavefronts (Data)41.8k

Workload Executionbrute\_force

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	2.89	SM Busy [%]	72.36
Executed Ipc Active Inst/cycle	2.89	Issue Slots Busy [%]	72.36

Balanced

ALU is the highest-utilized pipeline (57.5%) based on active cycles, taking into account the rates of its different instructions. It executes integer and logic operations. It is well-utilized, but should not be a bottleneck.

Pipe Utilization (% of active cycles)

Utilization [%]

ALU57.5%FMA22.7%FP640.0%

Pipe Utilization (% of peak instructions executed)

Utilization [%]

LSU35.2%ALU22.9%ADU22.7%FMA19.3%XU13.6%CBU4.4%

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 tables with data for each memory unit.

Memory Throughput [byte/s]	12.24	Mem Busy [%]	35.22
L1/TEX Hit Rate [%]	98.49	Max Bandwidth [%]	57.53
L2 Hit Rate [%]	99.32	Mem Pipes Busy [%]	57.53

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.

L1/TEX Local Load Access Pattern

Est. Speedup: 69.95%

The memory access pattern for local loads from L1/TEX might not be optimal. On average, only 0.3 of the 32 bytes transmitted per sector are utilized by each thread. This could possibly be caused by a stride between threads. Check the [Source Analysis](#) section for uncoalesced local loads.

L1/TEX Local Store Access Pattern

Est. Speedup: 69.84%

The memory access pattern for local stores to L1/TEX might not be optimal. On average, only 0.3 of the 32 bytes transmitted per sector are utilized by each thread. This could possibly be caused by a stride between threads. Check the [Source Analysis](#) section for uncoalesced local stores.

Memory Chart

Show As:Transfer Size

Kernel

Global32.77 K Inst

Local10.60 M Inst

Texture0.00 Inst

Surface0.00 Inst

Shared0.00 Inst

L1/TEX CacheHit Rate: 98.49 %

L2 CacheHit Rate: 99.32 %

System Memory369.48u

Device Memory1.19 MB

Peer Memory18.55 MB

Shared Memory

	Instructions	Requests	Wavefronts	% Peak	Bank Conflicts
Shared Load	0	0	0	0	0
Shared Load Matrix	0	0	0	0	0
Shared Store	0	0	0	0	0
Shared Atomic	0	0	0	0	0
Other	-	-	186,368	0.50	0
Total	0	0	186,368	0.50	0

L1/TEX Cache

	Instructions	Requests	Wavefronts	% Peak	Sectors	Sectors/Req	Hit Rate	Bytes	Sector Misses to L2	% Peak to L2	Returns to SM
Global Load	0	0	0	0	0	0	0	0	0	0.98	8,587,264
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	32,768	32,768	32,768	0.09	131,072	4	0	4,194,304	6,858,978	18.23	-
Local Store	2,015,232	2,113,936	2,113,936	5.62	8,060,928	3.81	98.54	257,949,696	0	0	-
Surface Store	0	0	0	0	0	0	0	0	0	0	-
Global Reduction	0	0	0	0	0	0	0	0	0	0	-
DSMEM Reduction	0	0	0	0	0	0	0	0	0	0	-
Surface Reduction	0	0	0	0	0	0	0	0	0	0	-
Global Atomic ALU	0	0	0	0	0	0	0	0	0	0	-
Global Atomic CAS	0	0	0	0	0	0	0	0	0	0	see above
Surface Atomic ALU	0	0	0	0	0	0	0	0	0	0	see above
Surface Atomic CAS	0	0	0	0	0	0	0	0	0	0	-
Loads	8,587,264	8,587,264	8,587,264	22.82	32,951,360	3.84	98.85	1,054,443,520	369,480	0.98	8,587,264
Stores	2,048,000	2,146,304	2,146,304	5.70	8,192,000	3.82	96.96	262,144,000	6,858,978	18.23	-
Atomic & Reductions	0	0	0	0	0	0	0	0	0	0	-
Total	10,635,264	10,733,568	10,733,568	28.53	41,143,360	3.83	98.48	1,316,587,520	7,228,458	19.21	8,587,264

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 Load	95,585	371,886	3.89	0.84	100	11,900,352	7,787,397,682.63	0	0	0
L1/TEX Store	1,813,824	6,852,286	3.78	15.53	99.64	219,273,152	135,352,660,008.71	24,196	0	0
L1/TEX Atomic ALU	0	0	0	0	0	0	0	0	0	0
L1/TEX Atomic CAS	0	0	0	0	0	0	0	0	0	0
L1/TEX Reduction	0	0	0	0	0	0	0	0	0	0
L1/TEX Total	1,908,031	7,228,380	3.79	16.40	99.66	231,564,160	143,573,270,001.95	24,252	0	0
Eco Total	24,276	0.01	0.01	0.01	0.01	775,832	481,647,553.67	24,276	0	0
GPU Total	1,915,736	7,254,541	3.79	16.44	99.66	232,145,312	143,933,593,905.00	24,870	0	0

Device Memory

	Sectors	% Peak	Bytes	Throughput
Load	37,280	0.23	1,192,960	739,653,188.37
Store	579,591	3.59	18,546,912	11,495,365,104.56
Total	616,871	3.82	19,739,872	12,239,018,292.92

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]	7.61	No Eligible [%]	25.23
Eligible Warps Per Scheduler [warp]	2.49	One or More Eligible [%]	74.77
Issued Warp Per Scheduler	0.75		

Warps Per Scheduler

GPU Maximum Warps Per Scheduler919,148

Theoretical Warps Per Scheduler919,148

Active Warps Per Scheduler919,148

Eligible Warps Per Scheduler919,148

Issued Warp Per Scheduler919,148

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]	10.18	Avg. Active Threads Per Warp	25.94
Warp Cycles Per Executed Instruction [cycle]	10.18	Avg. Not Predicted Off Threads Per Warp	24.51

Warp State (All Cycles)

Stall Wait2.56

Stall Not Selected2.56

Stall Long Scoreboard1.12

Stall Math Pipe Throttle1.12

Selected0.00

Stall LG Throttle1.12

Stall Short Scoreboard1.12

Stall MIO Throttle1.12

Stall Branch Resolving1.12

Stall Dispatch Stall1.12

Stall No Instruction1.12

Stall Drain0.00

Stall IMC Miss0.00

Stall Misc0.00

Stall Barrier0.00

Stall Membar0.00

Stall Sleeping0.00

Stall Tex Throttle0.00

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]	106,416,992	Avg. Executed Instructions Per Scheduler [inst]	665,106.20
Issued Instructions [inst]	106,420,249	Avg. Issued Instructions Per Scheduler [inst]	665,126.56

FP32 Non-Fused Instructions

Est. Speedup: 9.69%

This kernel executes 660,480 fused and 660,480 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 25% (relative to its current performance). Check the Source page to identify where this kernel executes FP32 instructions.

Executed Instruction Mix

IMAD15,000,000

LDP315,000,000

IAD0315,000,000

LDL15,000,000

LDC15,000,000

BRA15,000,000

ISETP15,000,000

I2F15,000,000

FMUL15,000,000

FFMA15,000,000

PRMT15,000,000

STL15,000,000

SHL15,000,000

BSYNC15,000,000

SEL15,000,000

BSYS15,000,000

BMOV15,000,000

S2R15,000,000

LEA15,000,000

EXIT15,000,000

ULDC15,000,000

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	1,024	Function Cache Configuration	CachePreference
Registers Per Thread [register/thread]	22	Static Shared Memory Per Block [byte/block]	0
Block Size	1,024	Dynamic Shared Memory Per Block [byte/block]	0
Threads Thread	1,048,576	Driver Shared Memory Per Block [byte/block]	0
Waves Per SM	25.60	Shared Memory Configuration Size [kbyte]	32.77
Uses Green Context	0	# SMs [SM]	40

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]	2
Average L2 Active Cycles [cycle]	308,435.50	Total SM Elapsed Cycles [cycle]	16
Achieved Occupancy [%]	92.22	Block Limit Shared Mem [block]	1
Achieved Active Warps Per SM [warp]	29.51	Block Limit SM [block]	16

Impact of Varying Register Count Per Thread

Warp Occupancy32

Impact of Varying Block Size

Warp Occupancy32

Impact of Varying Shared Memory Usage Per Block

Warp Occupancy32

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]	919,148.28	Average L1 Active Cycles [cycle]	919,148.28
Average L2 Active Cycles [cycle]	966,711.59	Average SMSP Active Cycles [cycle]	889,537.35
Average DRAM Active Cycles [cycle]	37,622,416	Total SM Elapsed Cycles [cycle]	37,622,416
Total L1 Elapsed Cycles [cycle]	37,622,416	Total L2 Elapsed Cycles [cycle]	44,128,824
Total SMSP Elapsed Cycles [cycle]	150,493,664	Total DRAM Elapsed Cycles [cycle]	64,606,208

L2 Slices Workload Imbalance

Est. Speedup: 51.9%

One or more L2 Slices have a much lower number of active cycles than the average number of active cycles. Maximum instance value is 7.32% above the average, while the minimum instance value is 10.07% below the average.

Workload Distribution

	Average	Min	Max	Sum
SM Active Cycles	919,148.28	901,652	938,605	36,765,931
SMSP Active Cycles	889,537.35	837,699	928,020	142,325,976
L1 Active Cycles	919,148.28	901,652	938,605	36,765,931
L2 Active Cycles	966,711.59	869,344	1,042,069	30,934,771
DRAM Active Cycles	308,435.50	279,936	340,596	2,467,484

Source Comments

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]	8,768,352	Branch Efficiency [%]	97.72
Branch Instructions Ratio [%]	0.08	Avg. Divergent Branches	11.19

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.