

Analysis Report

void stencilCompute2D<unsigned char, int=4>(unsigned char*, unsigned char*, unsigned int, unsigned int, unsigned char const *)

Duration	608.811 μ s
Grid Size	[144,109,1]
Block Size	[16,16,1]
Registers/Thread	20
Shared Memory/Block	0 B
Shared Memory Requested	96 KiB
Shared Memory Executed	96 KiB
Shared Memory Bank Size	4 B

[0] GeForce GTX 960

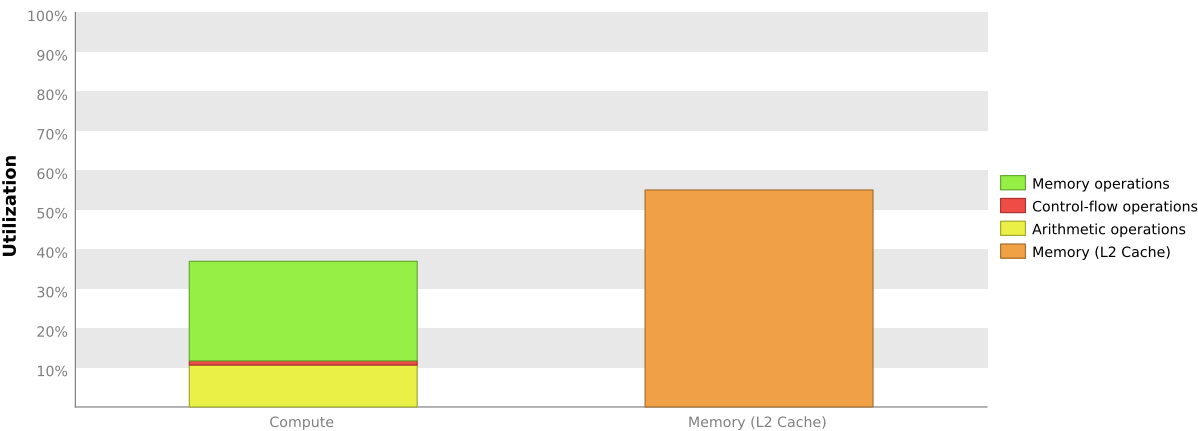
GPU UUID	GPU-0db32734-f94e-48a7-8b5d-4604317dc554
Compute Capability	5.2
Max. Threads per Block	1024
Max. Shared Memory per Block	48 KiB
Max. Registers per Block	65536
Max. Grid Dimensions	[2147483647, 65535, 65535]
Max. Block Dimensions	[1024, 1024, 64]
Max. Warps per Multiprocessor	64
Max. Blocks per Multiprocessor	32
Single Precision FLOP/s	2.644 TeraFLOP/s
Double Precision FLOP/s	82.624 GigaFLOP/s
Number of Multiprocessors	8
Multiprocessor Clock Rate	1.291 GHz
Concurrent Kernel	true
Max IPC	6
Threads per Warp	32
Global Memory Bandwidth	112.16 GB/s
Global Memory Size	4 GiB
Constant Memory Size	64 KiB
L2 Cache Size	1 MiB
Memcpy Engines	2
PCIe Generation	2
PCIe Link Rate	5 Gbit/s
PCIe Link Width	16

1. Compute, Bandwidth, or Latency Bound

The first step in analyzing an individual kernel is to determine if the performance of the kernel is bounded by computation, memory bandwidth, or instruction/memory latency. The results below indicate that the performance of kernel "void stencilCompute2D<unsig..." is most likely limited by instruction and memory latency. You should first examine the information in the "Instruction And Memory Latency" section to determine how it is limiting performance.

1.1. Kernel Performance Is Bound By Instruction And Memory Latency

This kernel exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of "GeForce GTX 960". These utilization levels indicate that the performance of the kernel is most likely limited by the latency of arithmetic or memory operations. Achieved compute throughput and/or memory bandwidth below 60% of peak typically indicates latency issues.

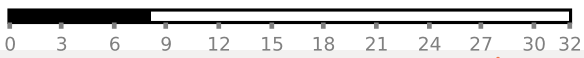
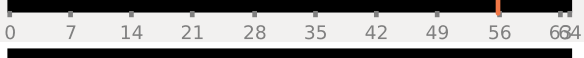



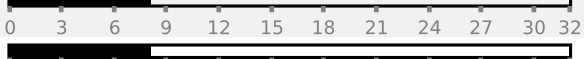



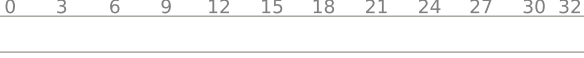
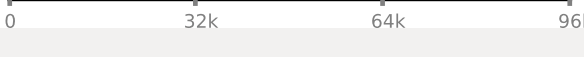


2. Instruction and Memory Latency

Instruction and memory latency limit the performance of a kernel when the GPU does not have enough work to keep busy. The performance of latency-limited kernels can often be improved by increasing occupancy. Occupancy is a measure of how many warps the kernel has active on the GPU, relative to the maximum number of warps supported by the GPU. Theoretical occupancy provides an upper bound while achieved occupancy indicates the kernel's actual occupancy.

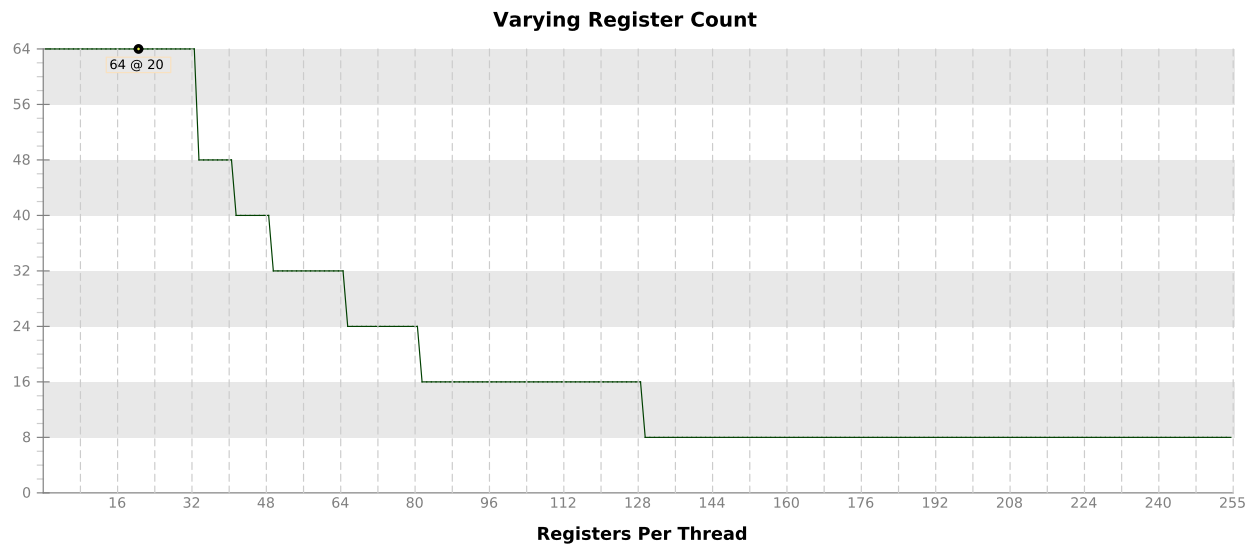
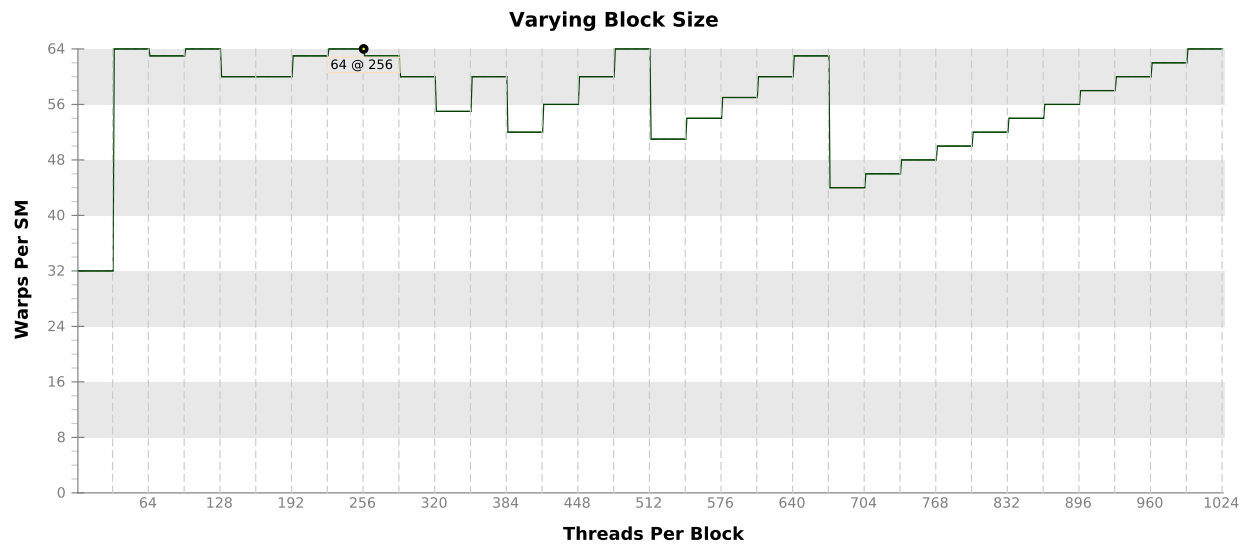
2.1. Occupancy Is Not Limiting Kernel Performance

The kernel's block size, register usage, and shared memory usage allow it to fully utilize all warps on the GPU.

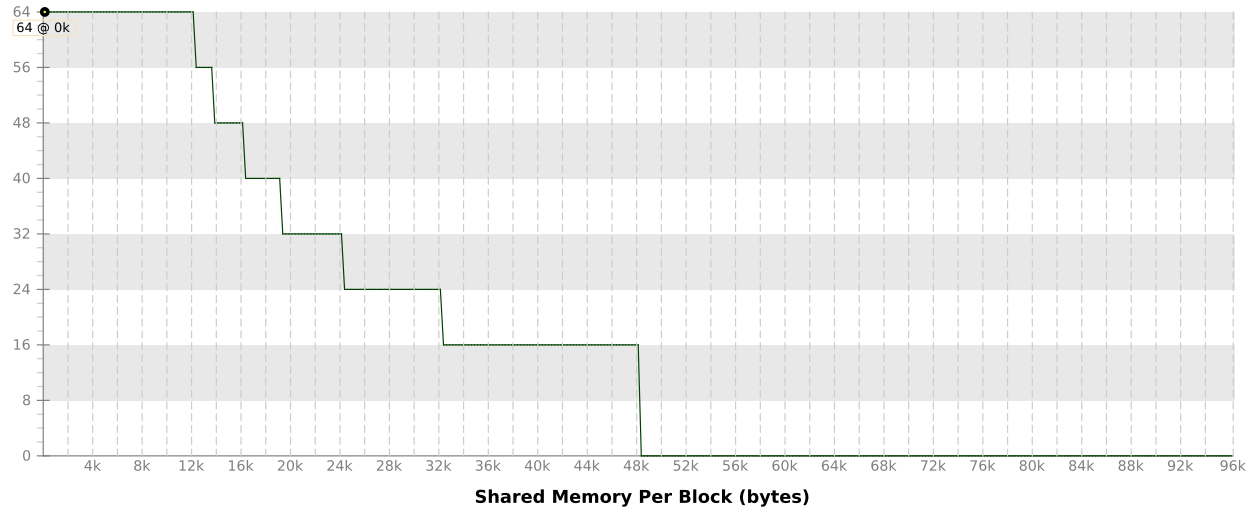
Variable	Achieved	Theoretical	Device Limit	Grid Size: [144,109,1] (15696 blocks) Block Size: [16,16,1] (256)
Occupancy Per SM				
Active Blocks		8	32	
Active Warps	55.61	64	64	
Active Threads		2048	2048	
Occupancy	86.9%	100%	100%	
Warps				
Threads/Block		256	1024	
Warps/Block		8	32	
Block Limit		8	32	
Registers				
Registers/Thread		20	255	
Registers/Block		6144	65536	
Block Limit		10	32	
Shared Memory				
Shared Memory/Block		0	98304	
Block Limit			32	

2.2. Occupancy Charts

The following charts show how varying different components of the kernel will impact theoretical occupancy.



Varying Shared Memory Usage



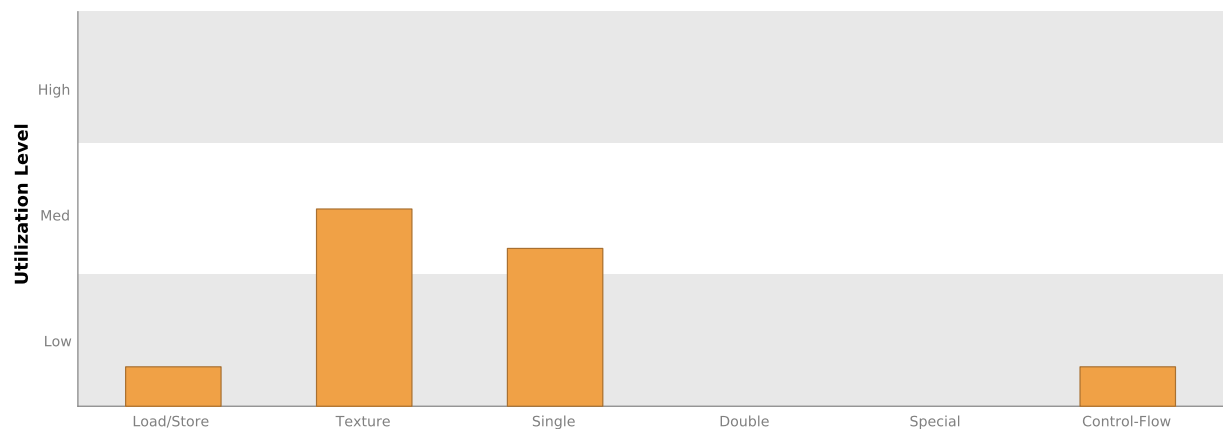
3. Compute Resources

GPU compute resources limit the performance of a kernel when those resources are insufficient or poorly utilized.

3.1. Function Unit Utilization

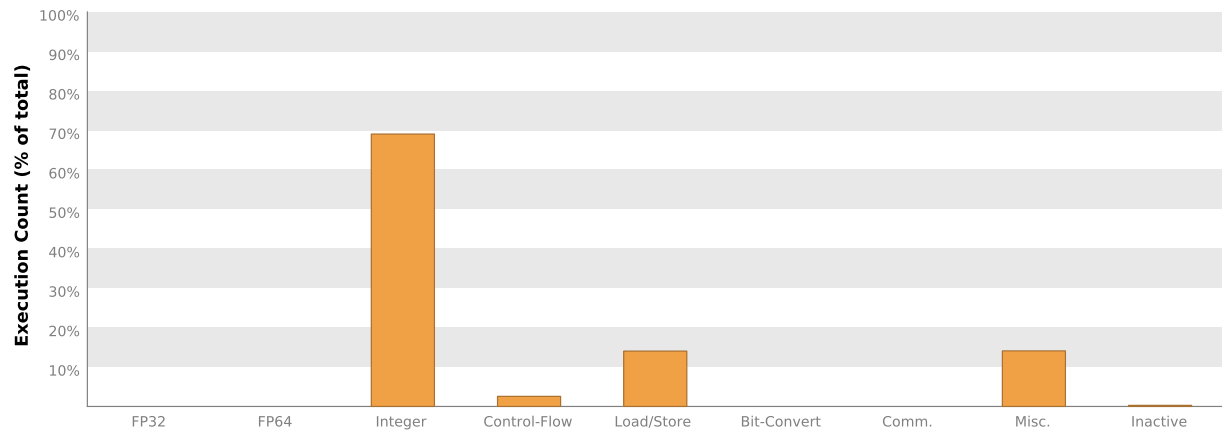
Different types of instructions are executed on different function units within each SM. Performance can be limited if a function unit is over-used by the instructions executed by the kernel. The following results show that the kernel's performance is not limited by overuse of any function unit.

- Load/Store - Load and store instructions for shared and constant memory.
- Texture - Load and store instructions for local, global, and texture memory.
- Single - Single-precision integer and floating-point arithmetic instructions.
- Double - Double-precision floating-point arithmetic instructions.
- Special - Special arithmetic instructions such as sin, cos, popc, etc.
- Control-Flow - Direct and indirect branches, jumps, and calls.



3.2. Instruction Execution Counts

The following chart shows the mix of instructions executed by the kernel. The instructions are grouped into classes and for each class the chart shows the percentage of thread execution cycles that were devoted to executing instructions in that class. The "Inactive" result shows the thread executions that did not execute any instruction because the thread was predicated or inactive due to divergence.



3.3. Floating-Point Operation Counts

The following chart shows the mix of floating-point operations executed by the kernel. The operations are grouped into classes and for each class the chart shows the percentage of thread execution cycles that were devoted to executing operations in that class. The results do not sum to 100% because non-floating-point operations executed by the kernel are not shown in this chart.



4. Memory Bandwidth

Memory bandwidth limits the performance of a kernel when one or more memories in the GPU cannot provide data at the rate requested by the kernel. The results below indicate that the kernel is limited by the bandwidth available to the L2 cache.

4.1. Global Memory Alignment and Access Pattern

Memory bandwidth is used most efficiently when each global memory load and store has proper alignment and access pattern.



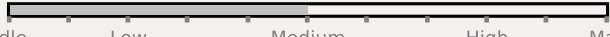


Optimization: Each entry below points to a global load or store within the kernel with an inefficient alignment or access pattern. For each load or store improve the alignment and access pattern of the memory access.

`/home/adas/cuda-workspace/CudaVisionSysDeploy/Release/./src/init/./device/LBPHist/LBPcompute.cuh`

Line 62	Global Load L2 Transactions/Access = 3, Ideal Transactions/Access = 1 [372810 L2 transactions for 124992 total executions]
Line 62	Global Load L2 Transactions/Access = 5.6, Ideal Transactions/Access = 1 [702045 L2 transactions for 124992 total executions]
Line 62	Global Load L2 Transactions/Access = 2, Ideal Transactions/Access = 1 [249696 L2 transactions for 124992 total executions]
Line 62	Global Load L2 Transactions/Access = 3, Ideal Transactions/Access = 1 [372810 L2 transactions for 124992 total executions]
Line 62	Global Load L2 Transactions/Access = 3, Ideal Transactions/Access = 1 [372810 L2 transactions for 124992 total executions]
Line 62	Global Load L2 Transactions/Access = 3, Ideal Transactions/Access = 1 [372810 L2 transactions for 124992 total executions]
Line 62	Global Load L2 Transactions/Access = 2, Ideal Transactions/Access = 1 [249696 L2 transactions for 124992 total executions]
Line 62	Global Load L2 Transactions/Access = 3, Ideal Transactions/Access = 1 [372810 L2 transactions for 124992 total executions]
Line 62	Global Load L2 Transactions/Access = 2, Ideal Transactions/Access = 1 [249696 L2 transactions for 124992 total executions]
Line 62	Global Load L2 Transactions/Access = 3, Ideal Transactions/Access = 1 [372810 L2 transactions for 124992 total executions]
Line 62	Global Store L2 Transactions/Access = 2, Ideal Transactions/Access = 1 [249696 L2 transactions for 124992 total executions]
Line 64	Global Store L2 Transactions/Access = 1.9, Ideal Transactions/Access = 1 [3772 L2 transactions for 2028 total executions]

4.2. Memory Bandwidth And Utilization

The following table shows the memory bandwidth used by this kernel for the various types of memory on the device. The table also shows the utilization of each memory type relative to the maximum throughput supported by the memory.

Transactions	Bandwidth	Utilization	
Shared Memory			
Shared Loads	0	0 B/s	
Shared Stores	0	0 B/s	
Shared Total	0	0 B/s	
L2 Cache			
Reads	3051636	162.104 GB/s	
Writes	253474	13.465 GB/s	
Total	3305110	175.569 GB/s	
Unified Cache			
Local Loads	0	0 B/s	
Local Stores	0	0 B/s	
Global Loads	6883283	246.267 GB/s	
Global Stores	253468	13.464 GB/s	
Texture Reads	4993920	265.279 GB/s	
Unified Total	12130671	525.01 GB/s	
Device Memory			
Reads	125023	6.641 GB/s	
Writes	126532	6.721 GB/s	
Total	251555	13.363 GB/s	
System Memory			
[PCIe configuration: Gen2 x16, 5 Gbit/s]			
Reads	0	0 B/s	
Writes	5	265.601 kB/s	