

Analysis Report

```
horizontal_diffusion_gpu(int, int, int, int, int, int, int, int, int, int, int, int, int, int, int,
int, int, int, int, int, int, int, int, int, int, int, int, char, float*, float*, float*, float*,
float*, float*, float*, float*, float*, float*, float, float*, float, float)
```

Duration	11.996 ms (11,996,306 ns)
Grid Size	[28,39,1]
Block Size	[16,8,1]
Registers/Thread	63
Shared Memory/Block	3.516 KiB
Shared Memory Requested	16 KiB
Shared Memory Executed	16 KiB
Shared Memory Bank Size	4 B

[0] GeForce GT 730

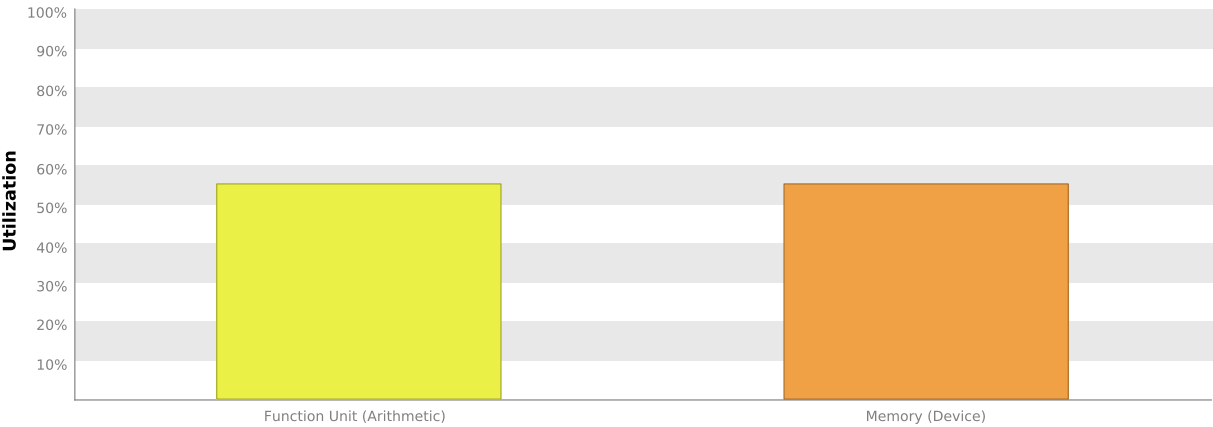
Compute Capability	2.1
Max. Threads per Block	1024
Max. Shared Memory per Block	48 KiB
Max. Registers per Block	32768
Max. Grid Dimensions	[65535, 65535, 65535]
Max. Block Dimensions	[1024, 1024, 64]
Max. Warps per Multiprocessor	48
Max. Blocks per Multiprocessor	8
Number of Multiprocessors	2
Multiprocessor Clock Rate	1.4 GHz
Concurrent Kernel	true
Max IPC	4
Threads per Warp	32
Global Memory Bandwidth	22.4 GB/s
Global Memory Size	2 GiB
Constant Memory Size	64 KiB
L2 Cache Size	128 KiB
Memcpy Engines	1
PCIe Generation	2
PCIe Link Rate	5 Gbit/s
PCIe Link Width	4

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 "horizontal_diffusion_gpu" 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 GT 730". 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. The results below indicate that occupancy can be improved by reducing the number of registers used by the kernel.

2.1. GPU Utilization Is Limited By Register Usage

The kernel uses 63 registers for each thread (8064 registers for each block). This register usage is likely preventing the kernel from fully utilizing the GPU. Device "GeForce GT 730" provides up to 32768 registers for each block. Because the kernel uses 8064 registers for each block each SM is limited to simultaneously executing 4 blocks (16 warps). Chart "Varying Register Count" below shows how changing register usage will change the number of blocks that can execute on each SM.

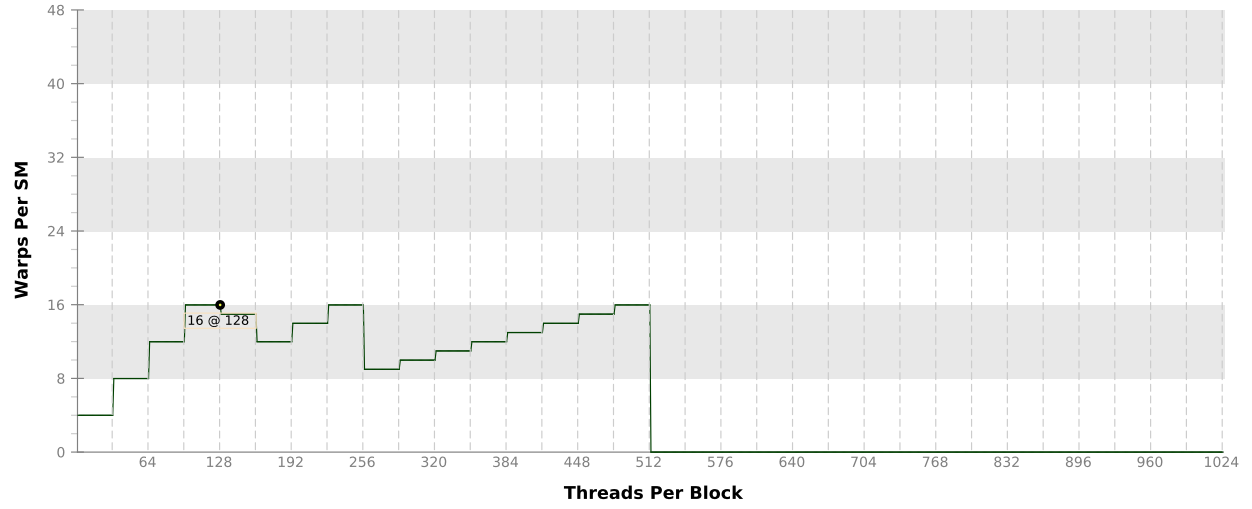
Optimization: Use the `-maxrregcount` flag or the `__launch_bounds__` qualifier to decrease the number of registers used by each thread. This will increase the number of blocks that can execute on each SM. On devices with Compute Capability 5.2 turning global cache off can increase the occupancy limited by register usage.

Variable	Achieved	Theoretical	Device Limit	Grid Size: [28,39,1] (1092 blocks) Block Size: [16,8,1] (128 threads)
Occupancy Per SM				
Active Blocks		4	8	
Active Warps	15.62	16	48	
Active Threads		512	1536	
Occupancy	32.5%	33.3%	100%	
Warps				
Threads/Block		128	1024	
Warps/Block		4	32	
Block Limit		12	8	
Registers				
Registers/Thread		63	63	
Registers/Block		8192	32768	
Block Limit		4	8	
Shared Memory				
Shared Memory/Block		3600	16384	
Block Limit		4	8	

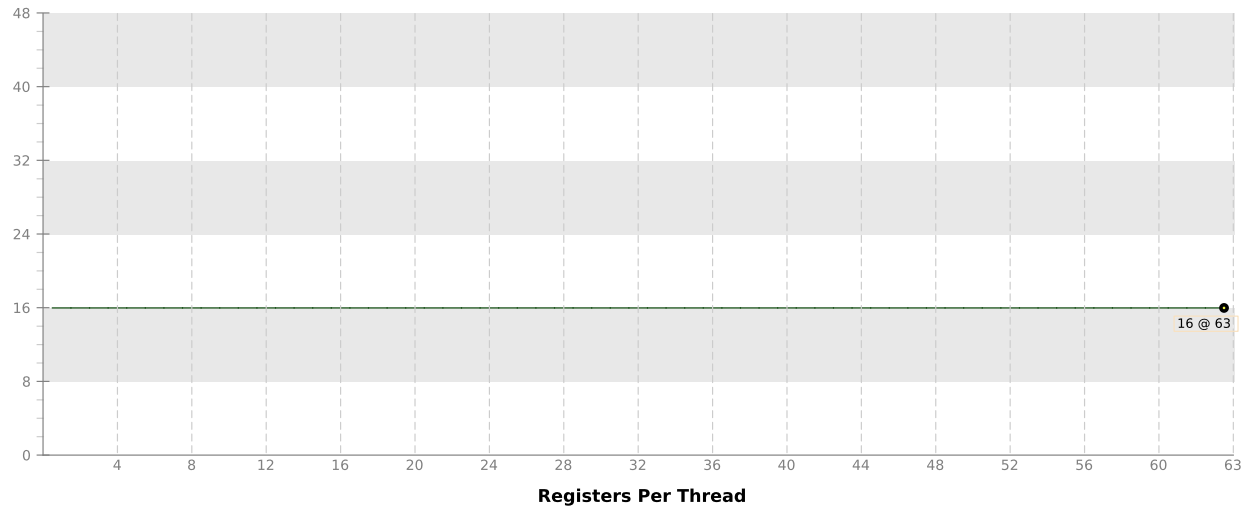
2.2. Occupancy Charts

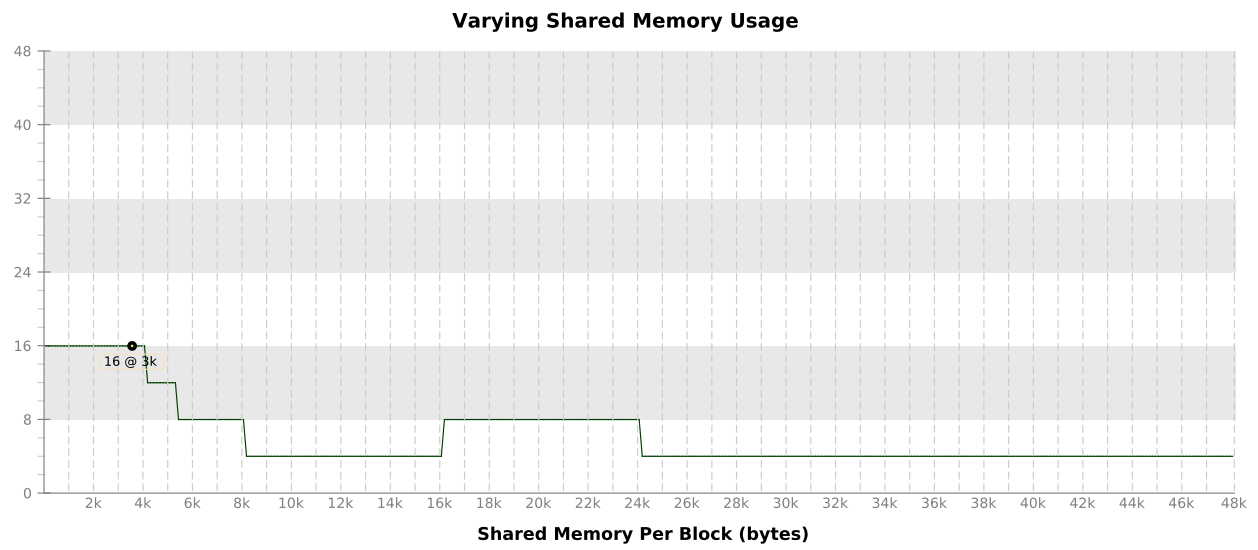
The following charts show how varying different components of the kernel will impact theoretical occupancy.

Varying Block Size



Varying Register Count





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

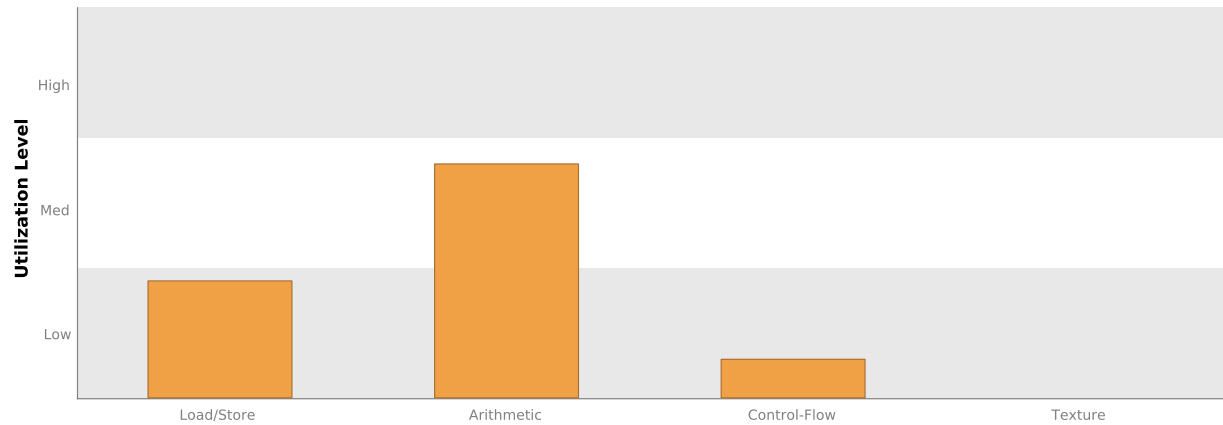
Different types of instructions are executed on different function units within each SM. Performance can be limited if a function unit is over-used by the instructions executed by the kernel. The following results show that the kernel's performance is not limited by overuse of any function unit.

Load/Store - Load and store instructions for local, shared, global, constant, etc. memory.

Arithmetic - All arithmetic instructions including integer and floating-point add and multiply, logical and binary operations, etc.

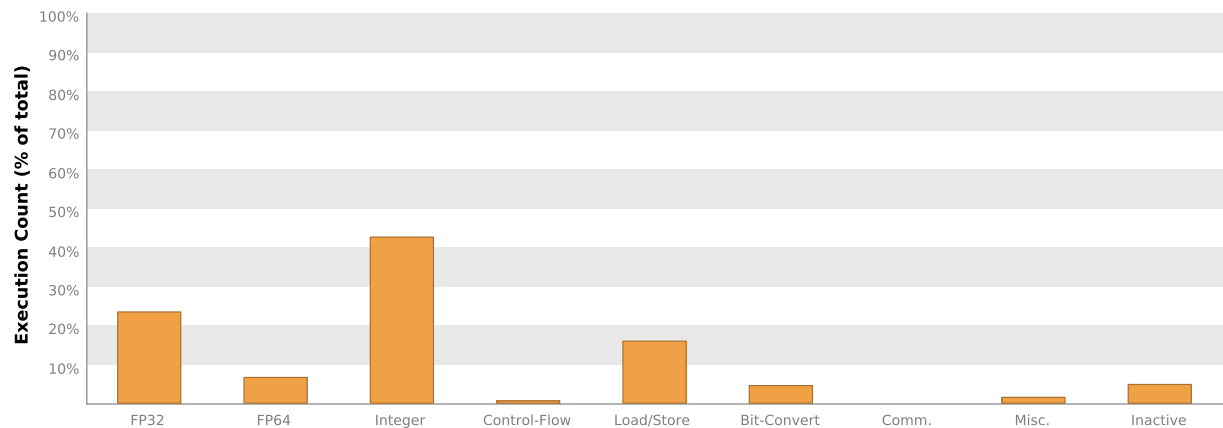
Control-Flow - Direct and indirect branches, jumps, and calls.

Texture - Texture operations.



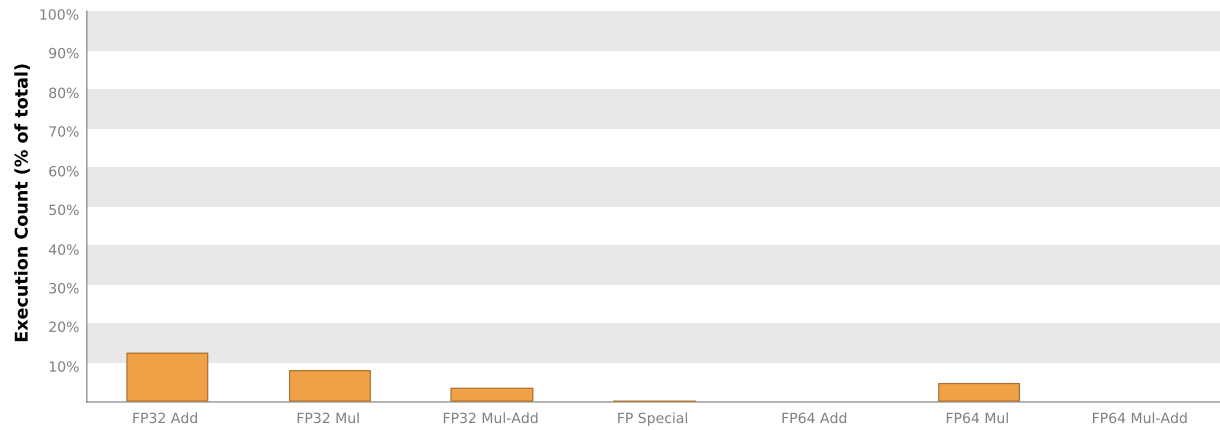
3.2. Instruction Execution Counts

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 device memory.

4.1. 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
L1/Shared Memory			
Local Loads	0	0 B/s	
Local Stores	0	0 B/s	
Shared Loads	123008	1.313 GB/s	
Shared Stores	153760	1.641 GB/s	
Global Loads	5025870	53.647 GB/s	
Global Stores	400276	2.061 GB/s	
Atomic	0	0 B/s	
L1/Shared Total	5702914	58.662 GB/s	
L2 Cache			
L1 Reads	5947208	15.87 GB/s	
L1 Writes	772504	2.061 GB/s	
Texture Reads	0	0 B/s	
Atomic	0	0 B/s	
Total	6719712	17.932 GB/s	
Texture Cache			
Reads	0	0 B/s	
Device Memory			
Reads	4171008	11.13 GB/s	
Writes	685102	1.828 GB/s	
Total	4856110	12.959 GB/s	
System Memory			
[PCIe configuration: Gen2 x4, 5 Gbit/s]			
Reads	0	0 B/s	
Writes	0	0 B/s	