

Analysis Report

calculate_difference(float*, float*, float*, float*, float, float*, float*)

Duration	154.543 ms (154,542,732 ns)
Grid Size	[64,1,1]
Block Size	[512,1,1]
Registers/Thread	30
Shared Memory/Block	4 KiB
Shared Memory Requested	112 KiB
Shared Memory Executed	112 KiB
Shared Memory Bank Size	4 B

[0] Tesla K80

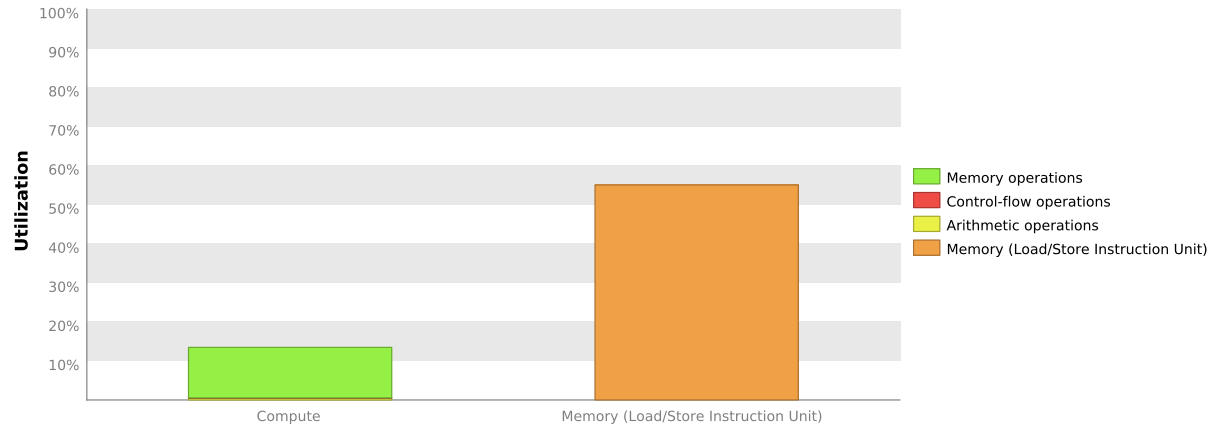
GPU UUID	GPU-094c6987-46f0-8163-ef91-00a3d25ad3d3
Compute Capability	3.7
Max. Threads per Block	1024
Max. Threads per Multiprocessor	2048
Max. Shared Memory per Block	48 KiB
Max. Shared Memory per Multiprocessor	112 KiB
Max. Registers per Block	65536
Max. Registers per Multiprocessor	131072
Max. Grid Dimensions	[2147483647, 65535, 65535]
Max. Block Dimensions	[1024, 1024, 64]
Max. Warps per Multiprocessor	64
Max. Blocks per Multiprocessor	16
Single Precision FLOP/s	4.111 TeraFLOP/s
Double Precision FLOP/s	1.37 TeraFLOP/s
Number of Multiprocessors	13
Multiprocessor Clock Rate	823.5 MHz
Concurrent Kernel	true
Max IPC	7
Threads per Warp	32
Global Memory Bandwidth	240.48 GB/s
Global Memory Size	11.172 GiB
Constant Memory Size	64 KiB
L2 Cache Size	1.5 MiB
Memcpy Engines	2
PCIe Generation	3
PCIe Link Rate	8 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 "calculate_difference" 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 "Tesla K80". 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 results below indicate that the GPU does not have enough work because instruction execution is stalling excessively.

2.1. Instruction Latencies May Be Limiting Performance

Instruction stall reasons indicate the condition that prevents warps from executing on any given cycle. The following chart shows the break-down of stalls reasons averaged over the entire execution of the kernel. The kernel has good theoretical and achieved occupancy indicating that there are likely sufficient warps executing on each SM. Since occupancy is not an issue it is likely that performance is limited by the instruction stall reasons described below.

Memory Dependency - A load/store cannot be made because the required resources are not available or are fully utilized, or too many requests of a given type are outstanding. Data request stalls can potentially be reduced by optimizing memory alignment and access patterns.

Not Selected - Warp was ready to issue, but some other warp issued instead. You may be able to sacrifice occupancy without impacting latency hiding and doing so may help improve cache hit rates.

Instruction Fetch - The next assembly instruction has not yet been fetched.

Constant - A constant load is blocked due to a miss in the constants cache.

Synchronization - The warp is blocked at a `__syncthreads()` call.

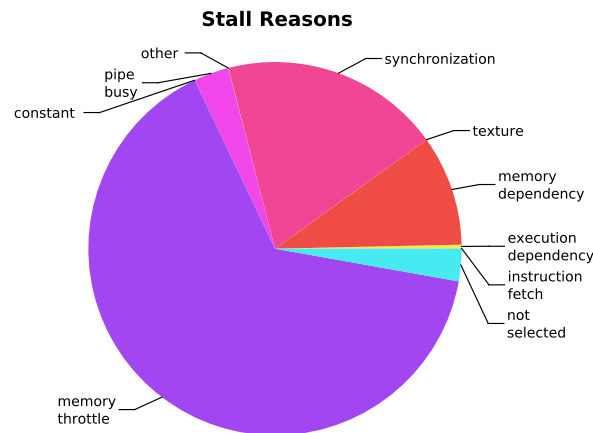
Pipeline Busy - The compute resource(s) required by the instruction is not yet available.

Execution Dependency - An input required by the instruction is not yet available. Execution dependency stalls can potentially be reduced by increasing instruction-level parallelism.

Texture - The texture sub-system is fully utilized or has too many outstanding requests.

Memory Throttle - Large number of pending memory operations prevent further forward progress. These can be reduced by combining several memory transactions into one.

Optimization: Resolve the primary stall issue; memory throttle.



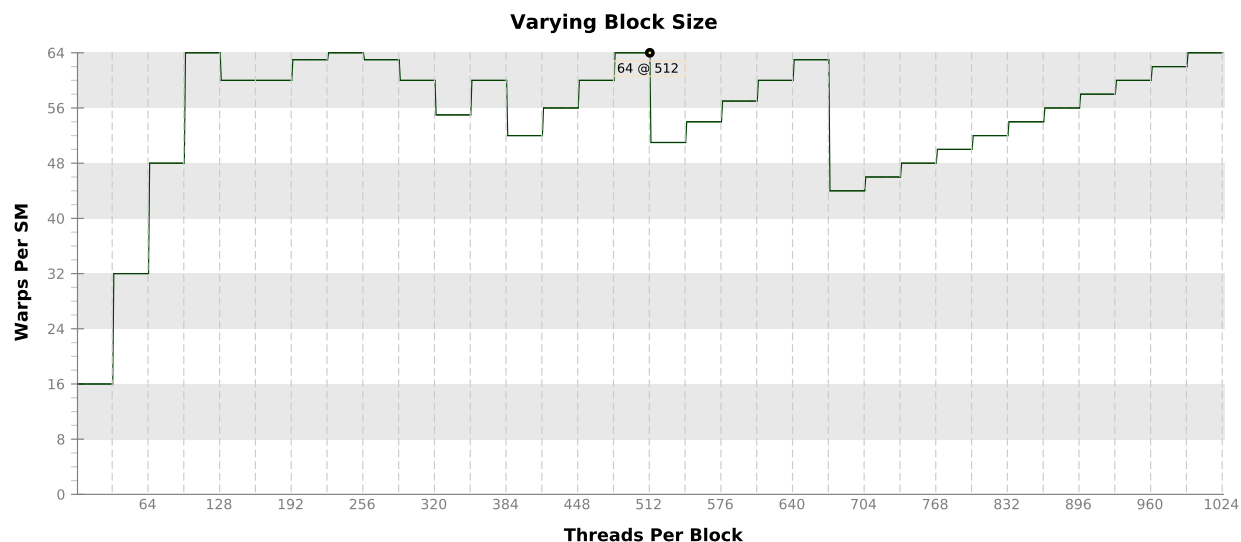
2.2. Occupancy Is Not Limiting Kernel Performance

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

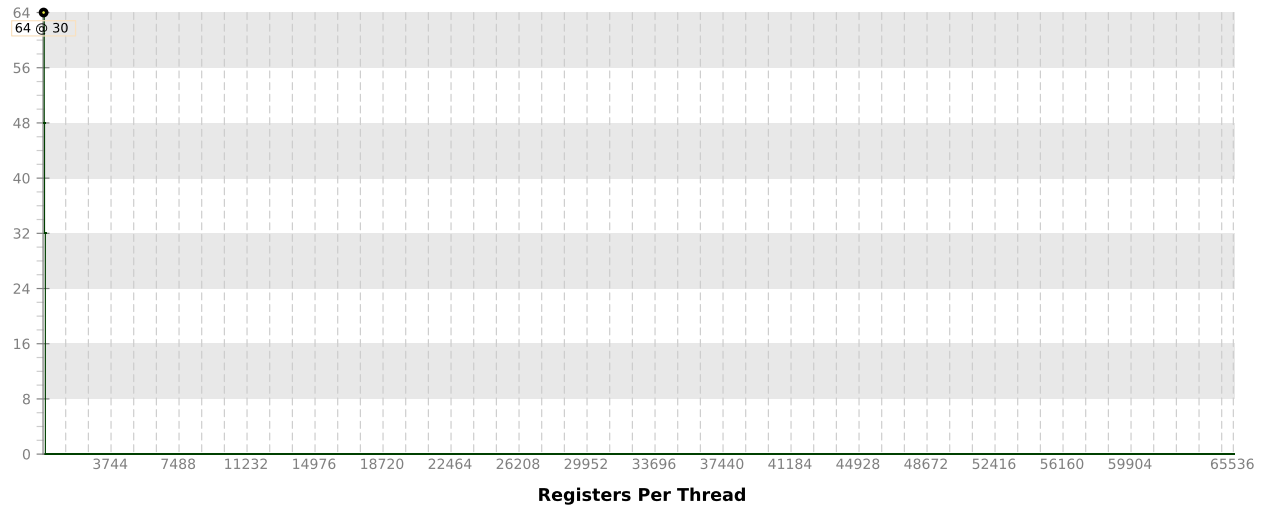
Variable	Achieved	Theoretical	Device Limit	Grid Size: [64,1,1] (64 blocks) Block Size: [512,1,1] (512 threads)
Occupancy Per SM				
Active Blocks		4	16	
Active Warps	62.61	64	64	
Active Threads		2048	2048	
Occupancy	97.8%	100%	100%	
Warps				
Threads/Block		512	1024	
Warps/Block		16	32	
Block Limit		4	16	
Registers				
Registers/Thread		30	65536	
Registers/Block		16384	131072	
Block Limit		8	16	
Shared Memory				
Shared Memory/Block		4096	114688	
Block Limit		28	16	

2.3. Occupancy Charts

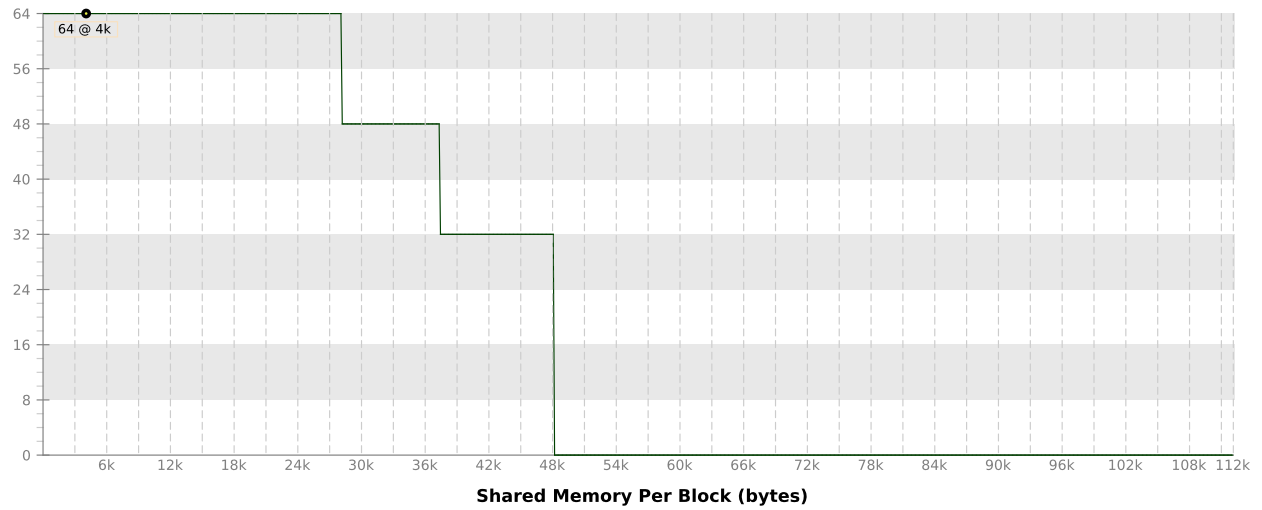
The following charts show how varying different components of the kernel will impact theoretical occupancy.



Varying Register Count



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

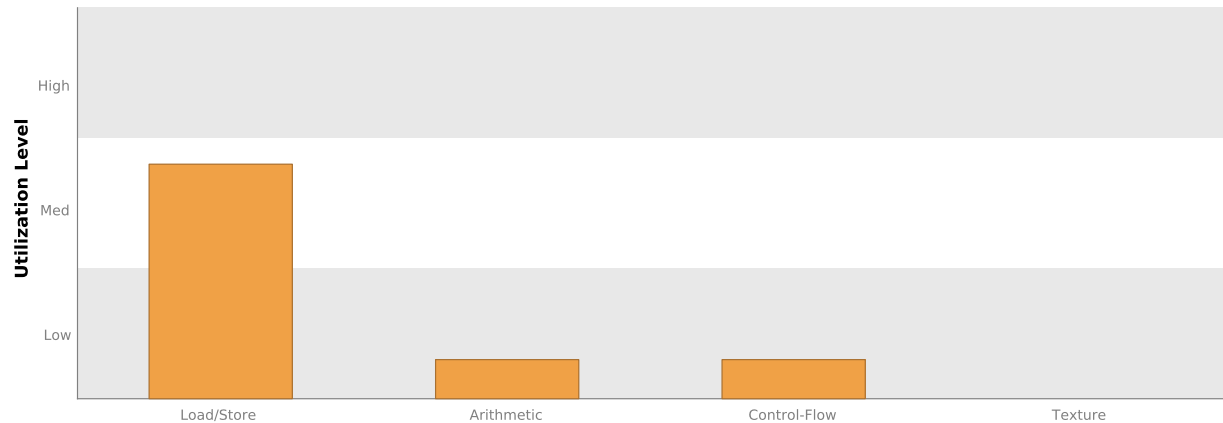
Different types of instructions are executed on different function units within each SM. Performance can be limited if a function unit is over-used by the instructions executed by the kernel. The following results show that the kernel's performance is not limited by overuse of any function unit.

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

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

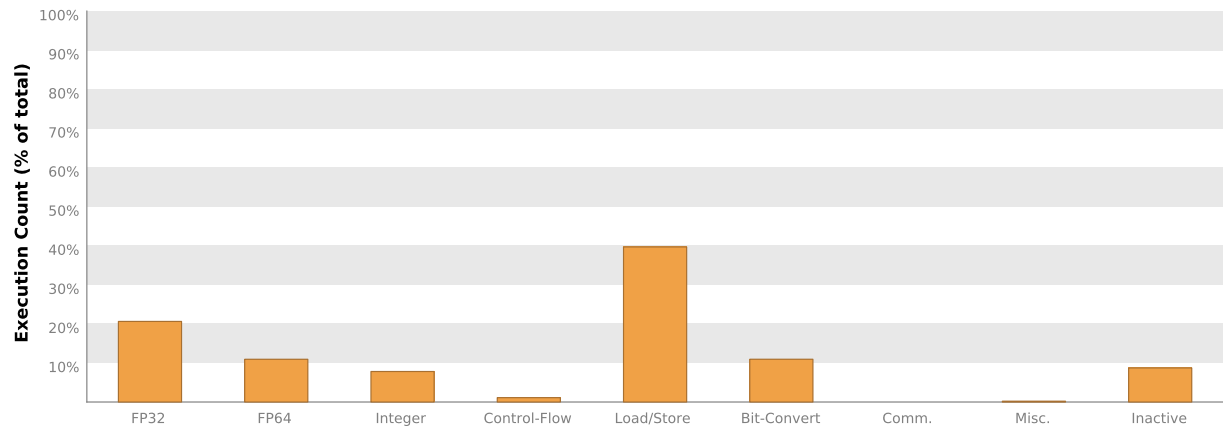
Control-Flow - Direct and indirect branches, jumps, and calls.

Texture - Texture operations.



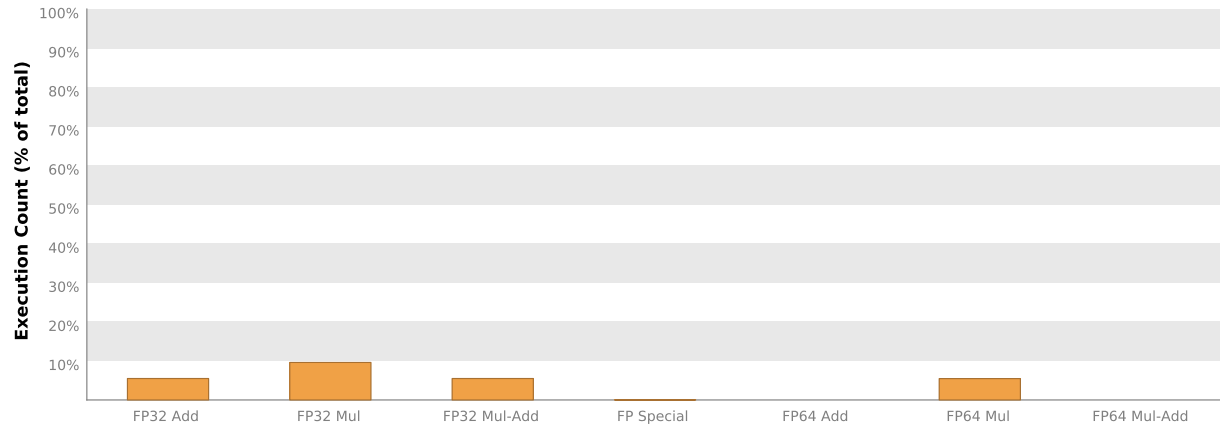
3.2. Instruction Execution Counts

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.

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	524288	868.483 MB/s	
Shared Stores	2182	3.614 MB/s	
Global Loads	135268352	28.662 GB/s	
Global Stores	67109888	13.897 GB/s	
Atomic	0	0 B/s	
L1/Shared Total	202904710	43.43 GB/s	
L2 Cache			
L1 Reads	138420224	28.662 GB/s	
L1 Writes	67112960	13.897 GB/s	
Texture Reads	0	0 B/s	
Noncoherent Reads	0	0 B/s	
Atomic	0	0 B/s	
Total	205533184	42.558 GB/s	
Texture Cache			
Reads	0	0 B/s	
Device Memory			
Reads	215492163	44.62 GB/s	
Writes	100606000	20.832 GB/s	
Total	316098163	65.452 GB/s	
ECC Overhead	155244570	32.145 GB/s	
System Memory			
[PCIe configuration: Gen3 x16, 8 Gbit/s]			
Reads	0	0 B/s	
Writes	4	828 B/s	