

Analysis Report

```
mult_kernel_compressed_data(unsigned char*, unsigned char*, unsigned char*, int, int, int, int)
```

| | |
|-------------------------|------------------------------|
| Duration | 8.34978 s (8,349,783,667 ns) |
| Grid Size | [750,750,1] |
| Block Size | [32,32,1] |
| Registers/Thread | 32 |
| Shared Memory/Block | 16.25 KiB |
| Shared Memory Executed | 32.5 KiB |
| Shared Memory Bank Size | 4 B |

[0] GeForce GTX 980 Ti

| | |
|---------------------------------------|------------------------------------------|
| GPU UUID | GPU-1efff26a-ba19-4be3-1414-34841cb941c3 |
| Compute Capability | 5.2 |
| Max. Threads per Block | 1024 |
| Max. Threads per Multiprocessor | 2048 |
| Max. Shared Memory per Block | 48 KiB |
| Max. Shared Memory per Multiprocessor | 96 KiB |
| Max. Registers per Block | 65536 |
| Max. Registers per Multiprocessor | 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 | 7.271 TeraFLOP/s |
| Double Precision FLOP/s | 227.216 GigaFLOP/s |
| Number of Multiprocessors | 22 |
| Multiprocessor Clock Rate | 1.291 GHz |
| Concurrent Kernel | true |
| Max IPC | 6 |
| Threads per Warp | 32 |
| Global Memory Bandwidth | 336.48 GB/s |
| Global Memory Size | 5.94 GiB |
| Constant Memory Size | 64 KiB |
| L2 Cache Size | 3 MiB |
| Memcpy Engines | 2 |
| PCIe Generation | 3 |
| PCIe Link Rate | 8 Gbit/s |
| PCIe Link Width | 16 |

1. Compute, Bandwidth, or Latency Bound

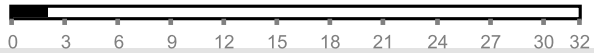
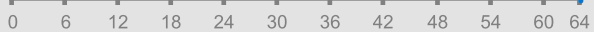
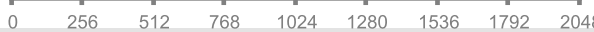
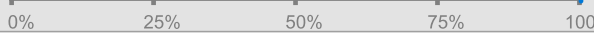








The first step in analyzing an individual kernel is to determine if the performance of the kernel is bounded by computation, memory bandwidth, or instruction/memory latency. Unfortunately, the device executing this kernel can not provide the profile data needed for this analysis.

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.

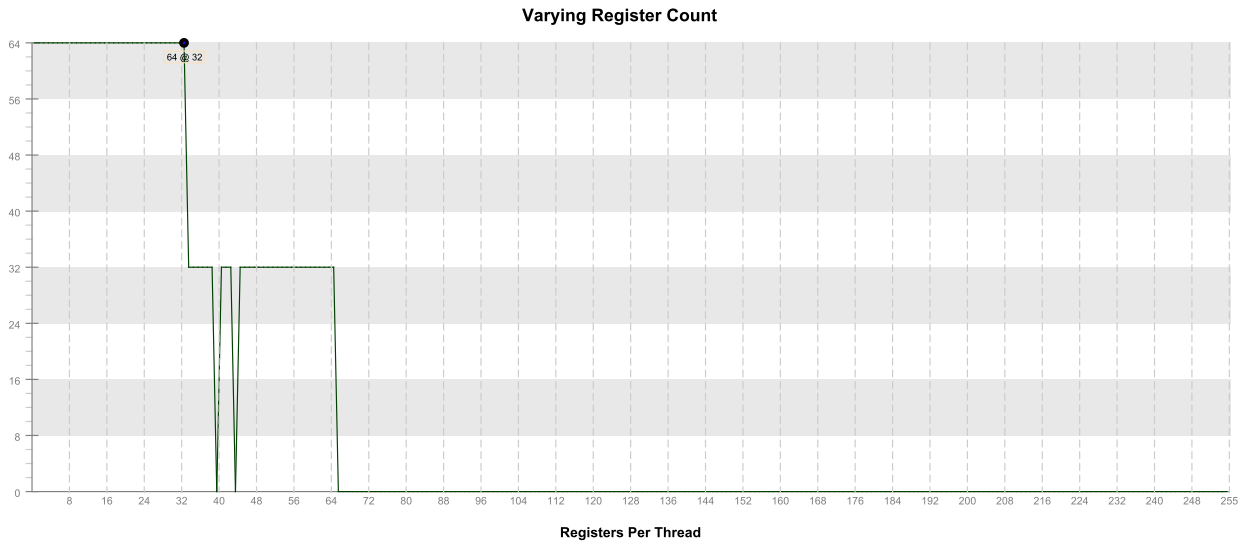
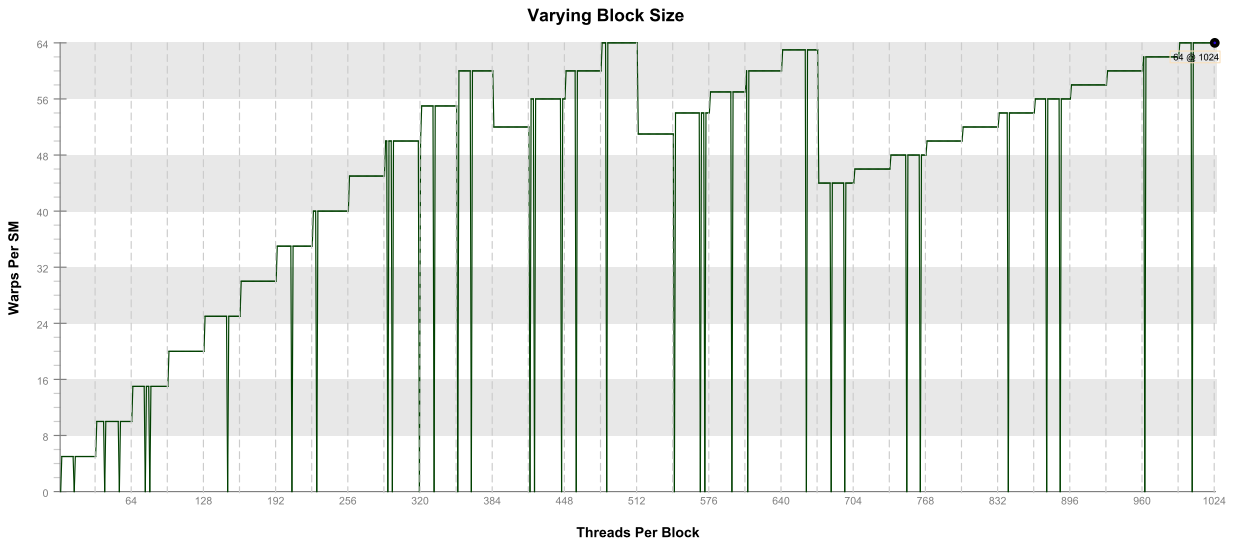
2.1. Occupancy Is Not Limiting Kernel Performance

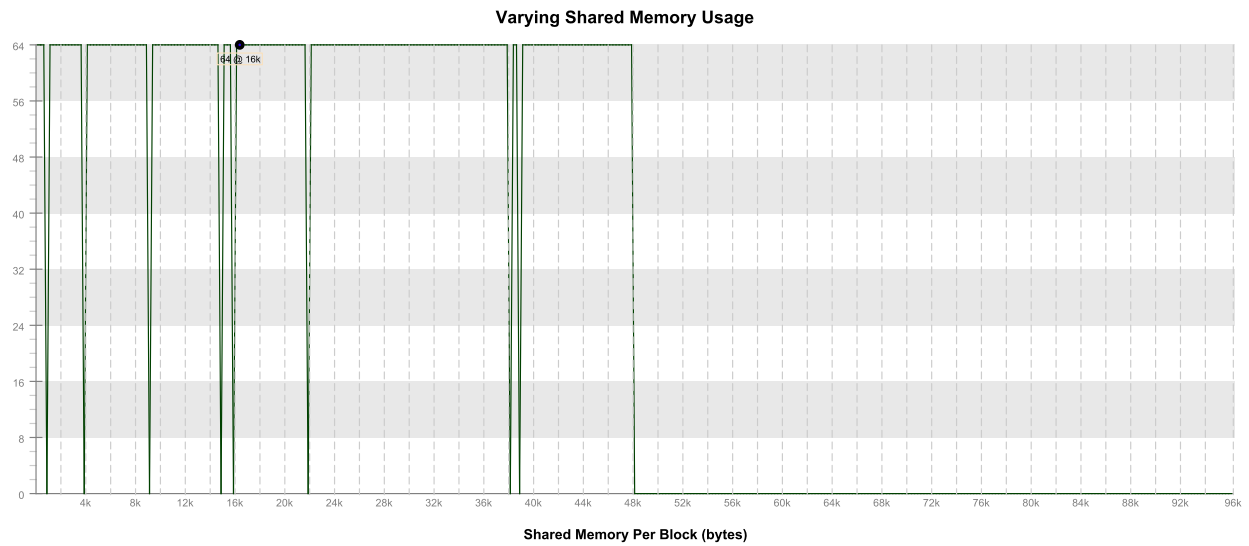
The kernel's block size, register usage, and shared memory usage allow it to fully utilize all warps on the GPU.

| Variable | Achieved | Theoretical | Device Limit | Grid Size: [750,750,1] (562500 blocks) Block Size: [32,32,1] (1024 thread |
|---------------------|----------|-------------|--------------|--------------------------------------------------------------------------------------|
| Occupancy Per SM | | | | |
| Active Blocks | | 2 | 32 |  |
| Active Warps | 63.98 | 64 | 64 |  |
| Active Threads | | 2048 | 2048 |  |
| Occupancy | 100% | 100% | 100% |  |
| Warps | | | | |
| Threads/Block | | 1024 | 1024 |  |
| Warps/Block | | 32 | 32 |  |
| Block Limit | | 2 | 32 |  |
| Registers | | | | |
| Registers/Thread | | 32 | 255 |  |
| Registers/Block | | 32768 | 65536 |  |
| Block Limit | | 2 | 32 |  |
| Shared Memory | | | | |
| Shared Memory/Block | | 16640 | 98304 |  |
| Block Limit | | 5 | 32 |  |

2.2. Occupancy Charts

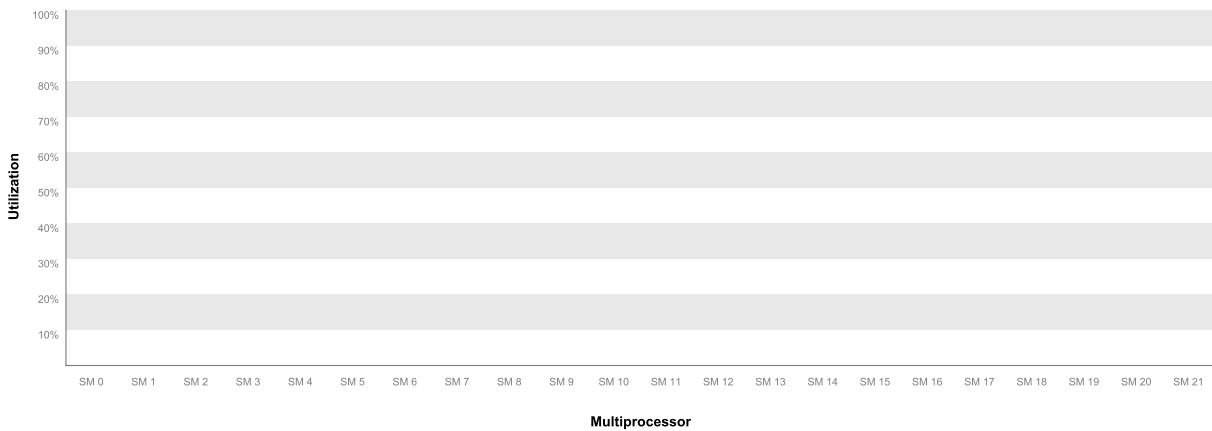
The following charts show how varying different components of the kernel will impact theoretical occupancy.





2.3. Multiprocessor Utilization

The kernel's blocks are distributed across the GPU's multiprocessors for execution. Depending on the number of blocks and the execution duration of each block some multiprocessors may be more highly utilized than others during execution of the kernel. The following chart shows the utilization of each multiprocessor during execution of the kernel.



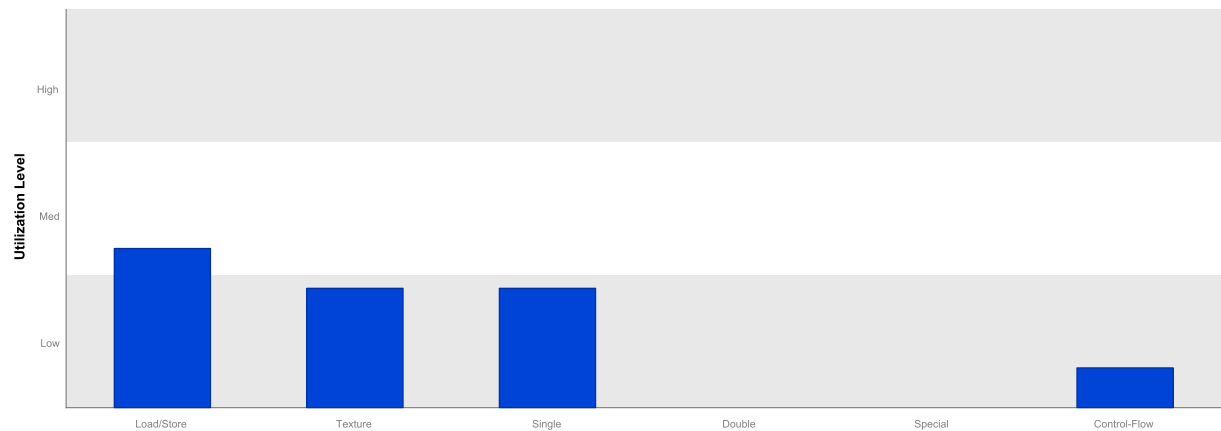
3. Compute Resources

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

3.1. Function Unit Utilization

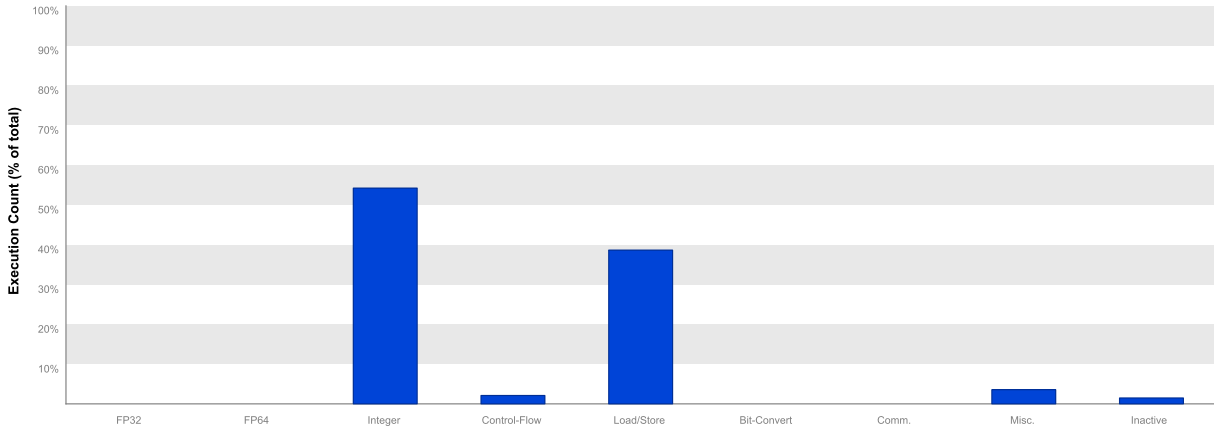
Different types of instructions are executed on different function units within each SM. Performance can be limited if a function unit is over-used by the instructions executed by the kernel. The following results show that the kernel's performance is not limited by overuse of any function unit.

- 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.



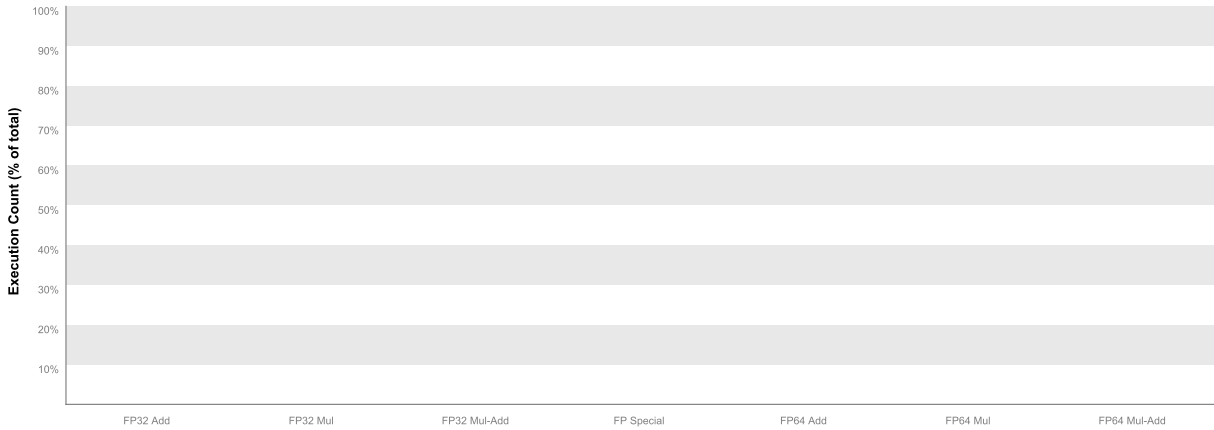
3.2. Instruction Execution Counts

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. Unfortunately, the device executing this kernel can not provide the profile data needed for this analysis.