# Analysis Report

## mysgemm(int, int, int, float const *, float const *, float*)

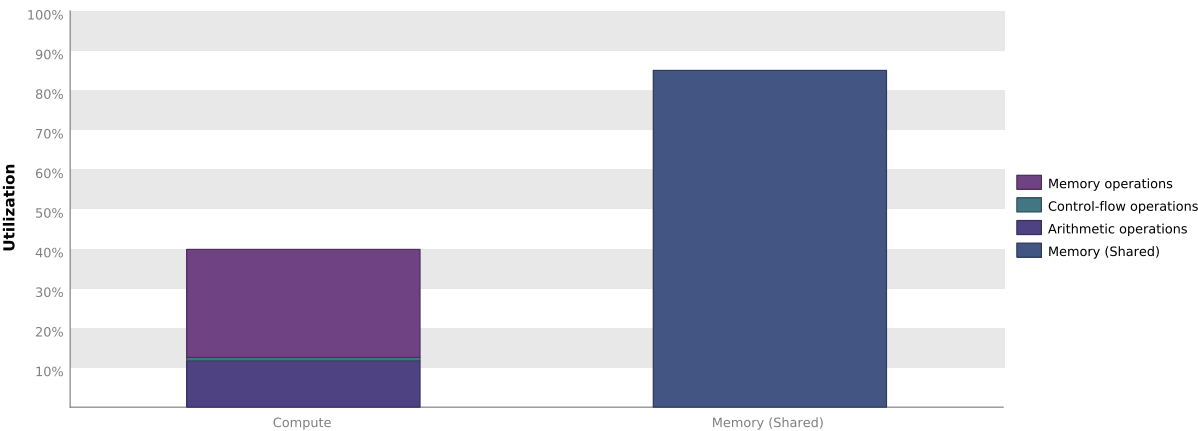| | |
|---|---|
| Duration | 7.31177 ms (7,311,771 ns) |
| Grid Size | [ 64,64,1 ] |
| Block Size | [ 16,16,1 ] |
| Registers/Thread | 30 |
| Shared  Memory/Block | 2 KiB |
| Shared Memory Requested | 96 KiB |
| Shared Memory Executed | 96 KiB |
| Shared Memory Bank Size | 4 B |

| [0] Tesla M60 | |
|---|---|
| GPU UUID | GPU-7dd05407-dba7-a874-62fe-55ac2820076c |
| 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 | 4.823 TeraFLOP/s |
| Double Precision FLOP/s | 150.72 GigaFLOP/s |
| Number of Multiprocessors | 16 |
| Multiprocessor Clock Rate | 1.177 GHz |
| Concurrent Kernel | true |
| Max IPC | 6 |
| Threads per Warp | 32 |
| Global Memory Bandwidth | 160.32 GB/s |
| Global Memory Size | 7.939 GiB |
| Constant Memory Size | 64 KiB |
| L2 Cache Size | 2 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 "mysgemm" is most likely limited by memory bandwidth. You should first examine the information in the "Memory Bandwidth" section to determine how it is limiting performance.

## 1.1. Kernel Performance Is Bound By Memory Bandwidth

For device "Tesla M60" the kernel's compute utilization is significantly lower than its memory utilization. These utilization levels indicate that the performance of the kernel is most likely being limited by the memory system. For this kernel the limiting factor in the memory system is the bandwidth of the Shared memory.

# 2. 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 shared memory.

## 2.1. GPU Utilization Is Limited By Memory Bandwidth

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. The results show that the kernel's performance is potentially limited by the bandwidth available from one or more of the memories on the device.

*Optimization: Try the following optimizations for the memory with high bandwidth utilization.*
*Shared Memory - If possible use 64-bit accesses to shared memory and 8-byte bank mode to achieved 2x throughput.*
*L2 Cache - Align and block kernel data to maximize L2 cache efficiency.*
*Unified Cache - Reallocate texture data to shared or global memory. Resolve alignment and access pattern issues for global loads and stores.*
*Device Memory - Resolve alignment and access pattern issues for global loads and stores.*
*System Memory (via PCIe) - Make sure performance critical data is placed in device or shared memory.*

| Transactions | Bandwidth | Utilization | |
|---|---|---|---|
| **Shared Memory** | | | |
| Shared Loads | 49545216 | 1,666.742 GB/s | |
| Shared Stores | 4128768 | 138.895 GB/s | |
| Shared Total | 53673984 | 1,805.637 GB/s | Idle · Low · Medium · High · Max |
| **L2 Cache** | | | |
| Reads | 16000071 | 134.564 GB/s | |
| Writes | 187006 | 1.573 GB/s | |
| Total | 16187077 | 136.137 GB/s | Idle · Low · Medium · High · Max |
| **Unified Cache** | | | |
| Local Loads | 0 | 0 B/s | |
| Local Stores | 0 | 0 B/s | |
| Global Loads | 32000000 | 134.563 GB/s | |
| Global Stores | 187000 | 1.573 GB/s | |
| Texture Reads | 16128000 | 135.64 GB/s | |
| Unified Total | 48315000 | 271.776 GB/s | Idle · Low · Medium · High · Max |
| **Device Memory** | | | |
| Reads | 7132832 | 59.989 GB/s | |
| Writes | 156755 | 1.318 GB/s | |
| Total | 7289587 | 61.307 GB/s | Idle · Low · Medium · High · Max |
| **System Memory** | | | |
| [ PCIe configuration: Gen3 x16, 8 Gbit/s ] | | | |
| Reads | 0 | 0 B/s | Idle · Low · Medium · High · Max |
| Writes | 5 | 42.051 kB/s | Idle · Low · Medium · High · Max |

3

# 3. 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.
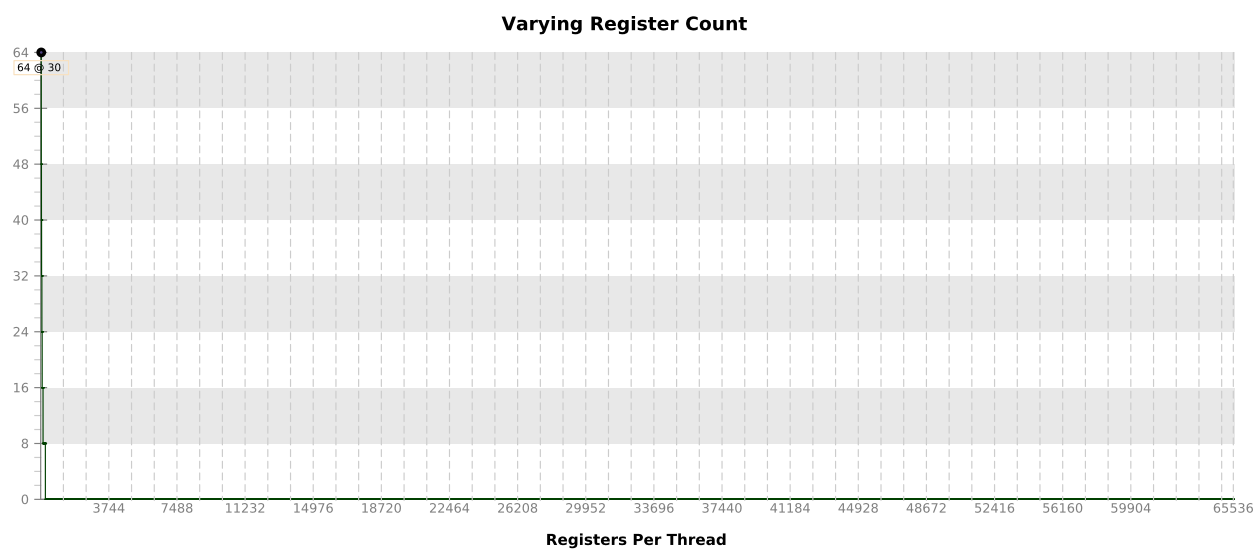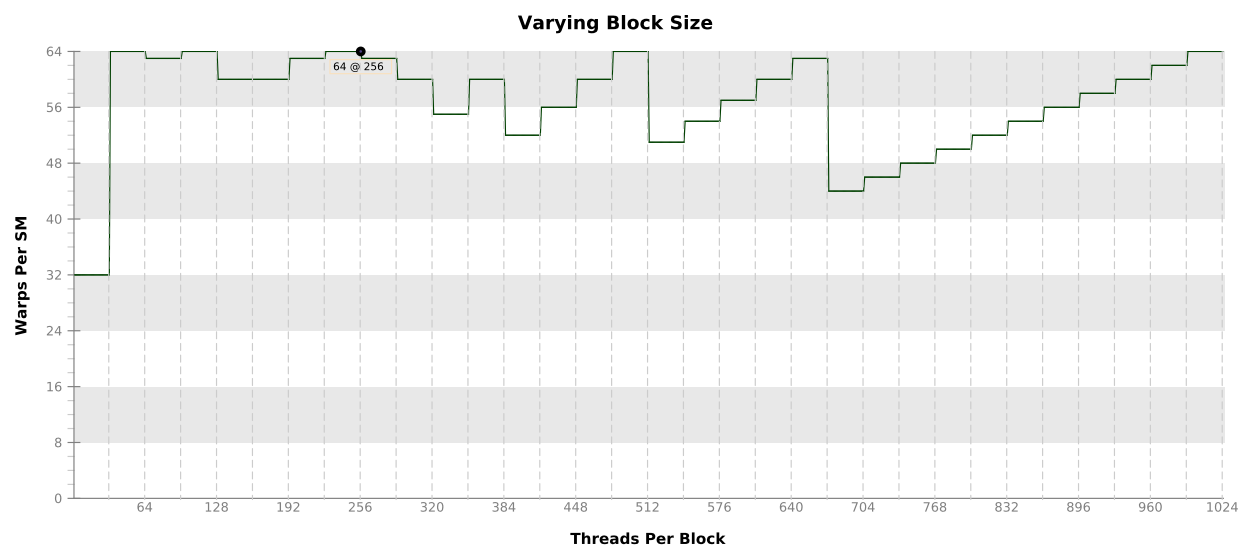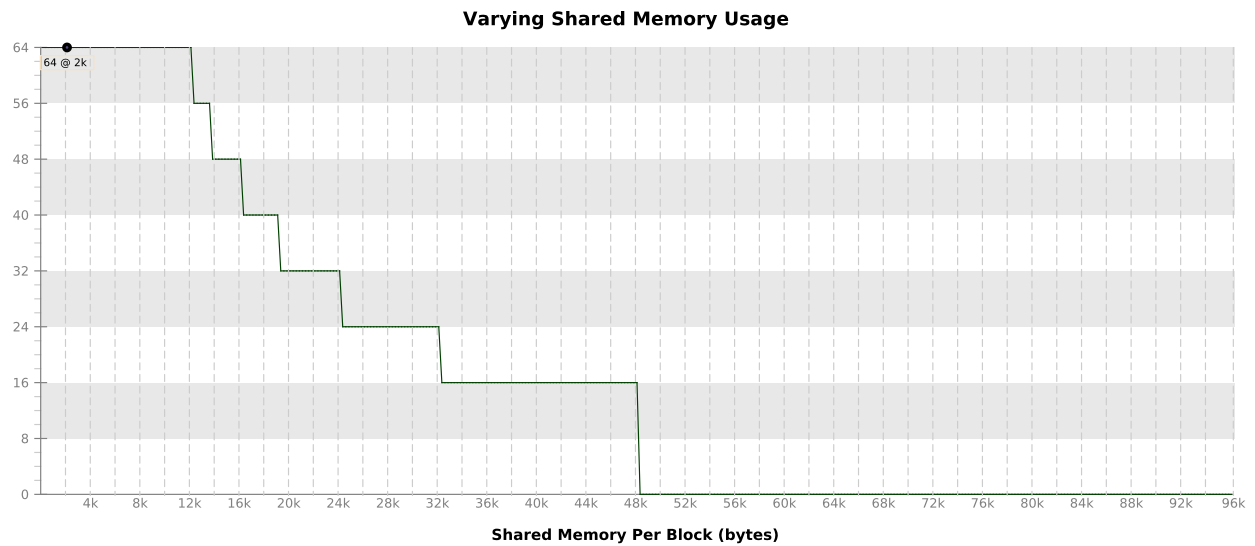
## 3.1. Occupancy Is Not Limiting Kernel Performance

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: [ 64,64,1 ] (4096 blocks) Block Size: [ 16,16,1 |
|---|---|---|---|---|
| **Occupancy Per SM** | | | | |
| Active Blocks | | 8 | 32 | |
| Active Warps | 63.4 | 64 | 64 | |
| Active Threads | | 2048 | 2048 | |
| Occupancy | 99.1% | 100% | 100% | |
| **Warps** | | | | |
| Threads/Block | | 256 | 1024 | |
| Warps/Block | | 8 | 32 | |
| Block Limit | | 8 | 32 | |
| **Registers** | | | | |
| Registers/Thread | | 30 | 65536 | |
| Registers/Block | | 8192 | 65536 | |
| Block Limit | | 8 | 32 | |
| **Shared Memory** | | | | |
| Shared Memory/Block | | 2048 | 98304 | |
| Block Limit | | 48 | 32 | |

## 3.2. Occupancy Charts

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

**Varying Block Size**

Warps Per SM

64 @ 256

Threads Per Block

**Varying Register Count**

64 @ 30

Registers Per Thread

**Varying Shared Memory Usage**

64 @ 2k

X-axis: Shared Memory Per Block (bytes)

### 3.3. 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.

Y-axis: Utilization

X-axis: Multiprocessor (SM 0 through SM 15)

# 4. Compute Resources

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

## 4.1. 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.
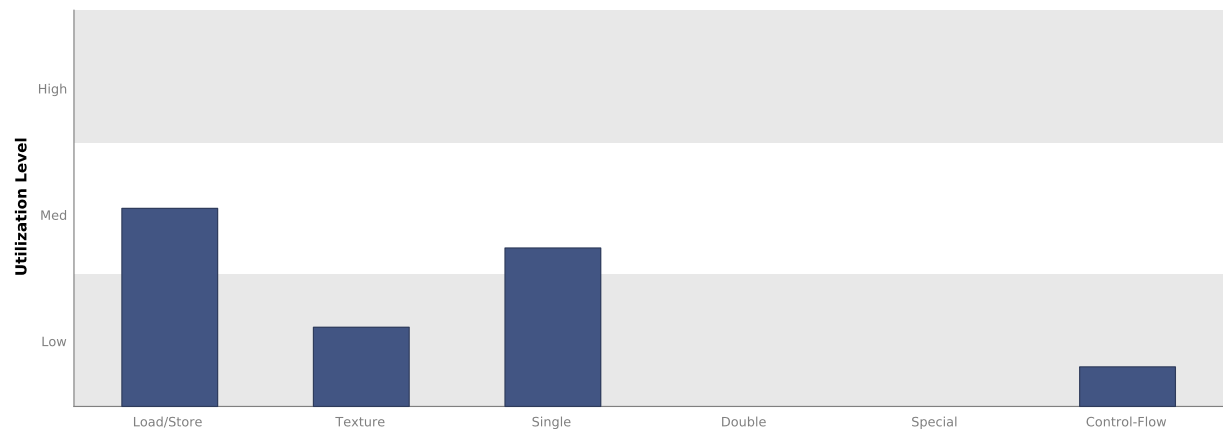 Load/Store - Load and store instructions for shared and constant memory.
 Texture - Load and store instructions for local, global, and texture memory.
 Single - Single-precision integer and floating-point arithmetic instructions.
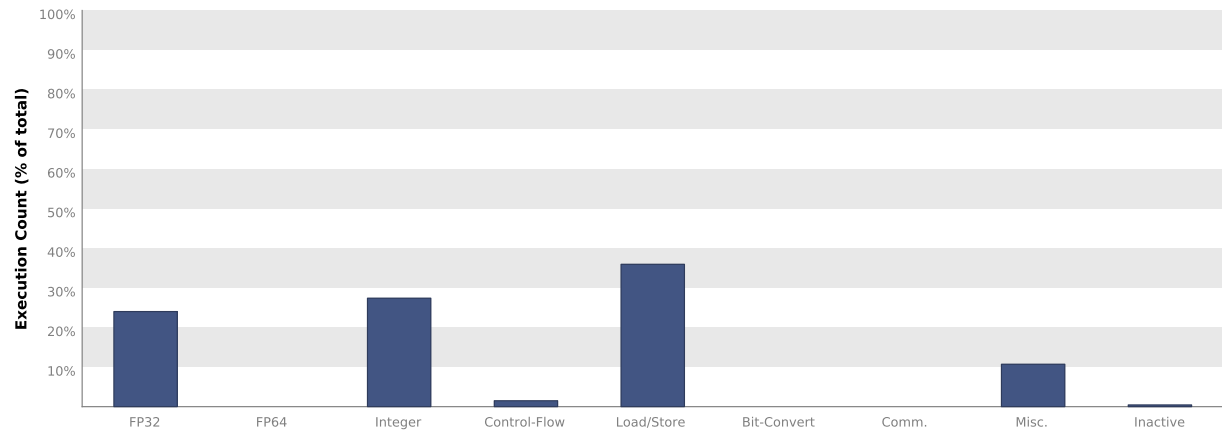 Double - Double-precision floating-point arithmetic instructions.
 Special - Special arithmetic instructions such as sin, cos, popc, etc.
 Control-Flow - Direct and indirect branches, jumps, and calls.



## 4.2. 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.

## 4.3. 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.