# Analysis Report

## matmult_gpu4Kernel(int, int, int, double*, double*, double*)

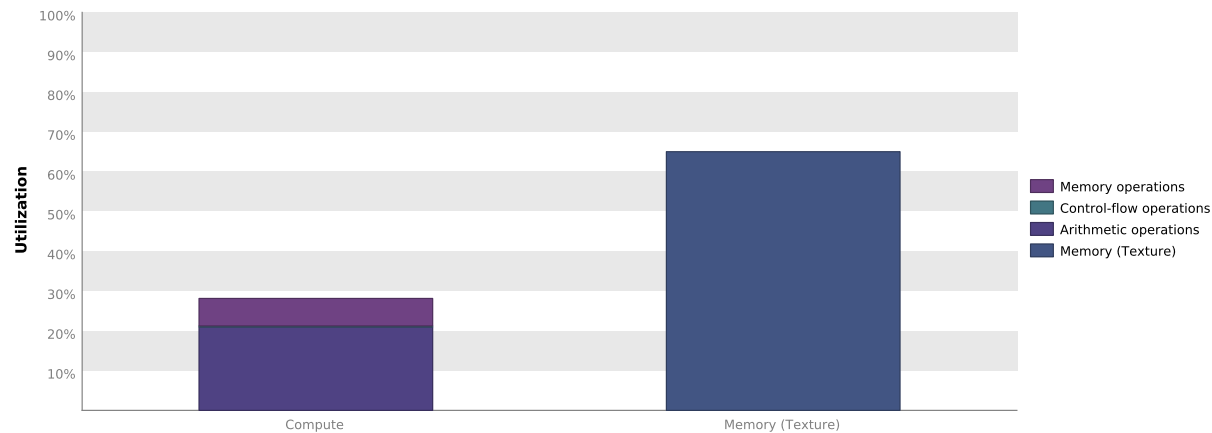| | |
|---|---|
| Duration | 316.27604 ms (316,276,041 ns) |
| Grid Size | [ 313,14,1 ] |
| Block Size | [ 16,16,1 ] |
| Registers/Thread | 72 |
| Shared  Memory/Block | 0 B |
| Shared Memory Requested | 96 KiB |
| Shared Memory Executed | 96 KiB |
| Shared Memory Bank Size | 4 B |

| [0] Tesla V100-SXM2-32GB | |
|---|---|
| GPU UUID | GPU-c2fbdc1c-622a-33f0-430e-76df65ca5880 |
| Compute Capability | 7.0 |
| 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 |
| Half Precision FLOP/s | 31.334 TeraFLOP/s |
| Single Precision FLOP/s | 15.667 TeraFLOP/s |
| Double Precision FLOP/s | 7.834 TeraFLOP/s |
| Number of Multiprocessors | 80 |
| Multiprocessor Clock Rate | 1.53 GHz |
| Concurrent Kernel | true |
| Max IPC | 4 |
| Threads per Warp | 32 |
| Global Memory Bandwidth | 898.048 GB/s |
| Global Memory Size | 31.719 GiB |
| Constant Memory Size | 64 KiB |
| L2 Cache Size | 6 MiB |
| Memcpy Engines | 5 |
| 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 "matmult_gpu4Kernel" 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 V100-SXM2-32GB" 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 Texture 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 unified cache that holds texture, global, and local data.

### 2.1. 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. The analysis is per assembly instruction.

*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.*

### 2.2. High Local Memory Overhead

Local memory loads and stores account for 24% of total memory traffic. High local memory traffic typically indicates excessive register spilling.

*Optimization: Use the -maxrregcount flag or the __launch_bounds__ qualifier to increase the number of registers available to nvcc when compiling the kernel.*

### 2.3. 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.*
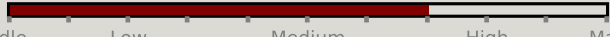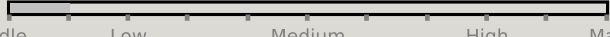 *Shared Memory - If possible use 64-bit accesses to shared memory and 8-byte bank mode to achieved 2x throughput.*
 *L2 Cache - Align and block kernel data to maximize L2 cache efficiency.*
 *Unified Cache - Reallocate texture data to shared or global memory. Resolve alignment and access pattern issues 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 | 0 | 0 B/s | |
| Shared Stores | 0 | 0 B/s | |
| Shared Total | 0 | 0 B/s | Idle　Low　Medium　High　Max |
| **L2 Cache** | | | |
| Reads | 1143758861 | 290.87 GB/s | |
| Writes | 28562754 | 7.264 GB/s | |
| Total | 1172321615 | 298.133 GB/s | Idle　Low　Medium　High　Max |
| **Unified Cache** | | | |
| Local Loads | 341760 | 86.913 MB/s | |
| Local Stores | 537600 | 136.717 MB/s | |
| Global Loads | 10333967424 | 2,628.033 GB/s | |
| Global Stores | 25621360 | 6.516 GB/s | |
| Texture Reads | 9912695192 | 10,083.597 GB/s | |
| Unified Total | 20273163336 | 12,718.37 GB/s | Idle　Low　Medium　High　Max |
| **Device Memory** | | | |
| Reads | 155293419 | 39.493 GB/s | |
| Writes | 9448409 | 2.403 GB/s | |
| Total | 164741828 | 41.896 GB/s | Idle　Low　Medium　High　Max |
| **System Memory** | | | |
| [ PCIe configuration: Gen3 x16, 8 Gbit/s ] | | | |
| Reads | 4000 | 1.017 MB/s | Idle　Low　Medium　High　Max |
| Writes | 3845 | 977.822 kB/s | Idle　Low　Medium　High　Max |

# 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. The results below indicate that occupancy can be improved by reducing the number of registers used by the kernel.

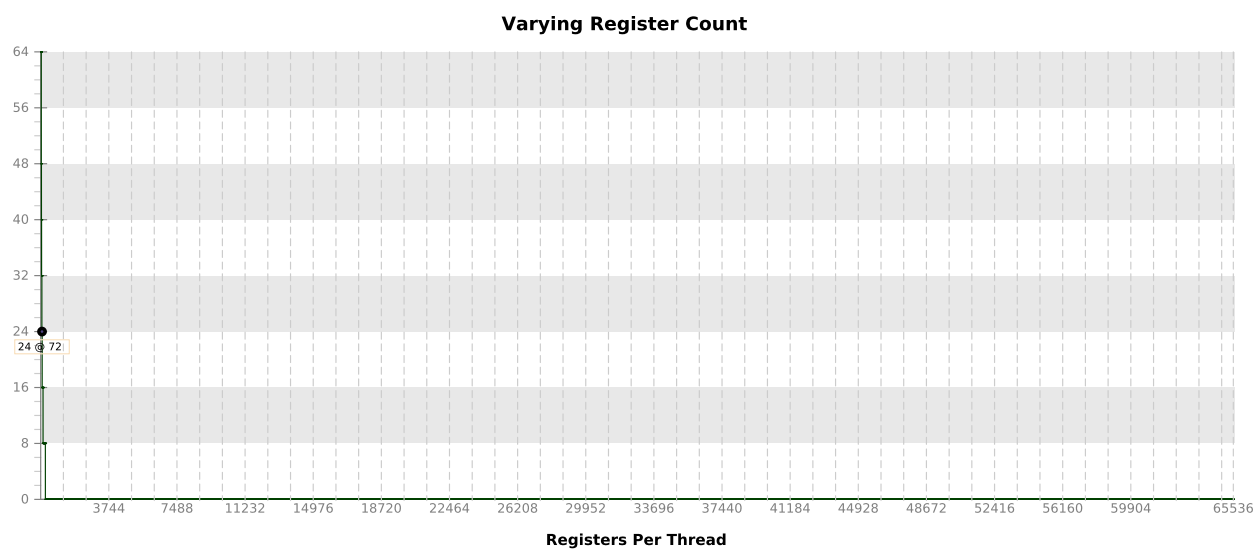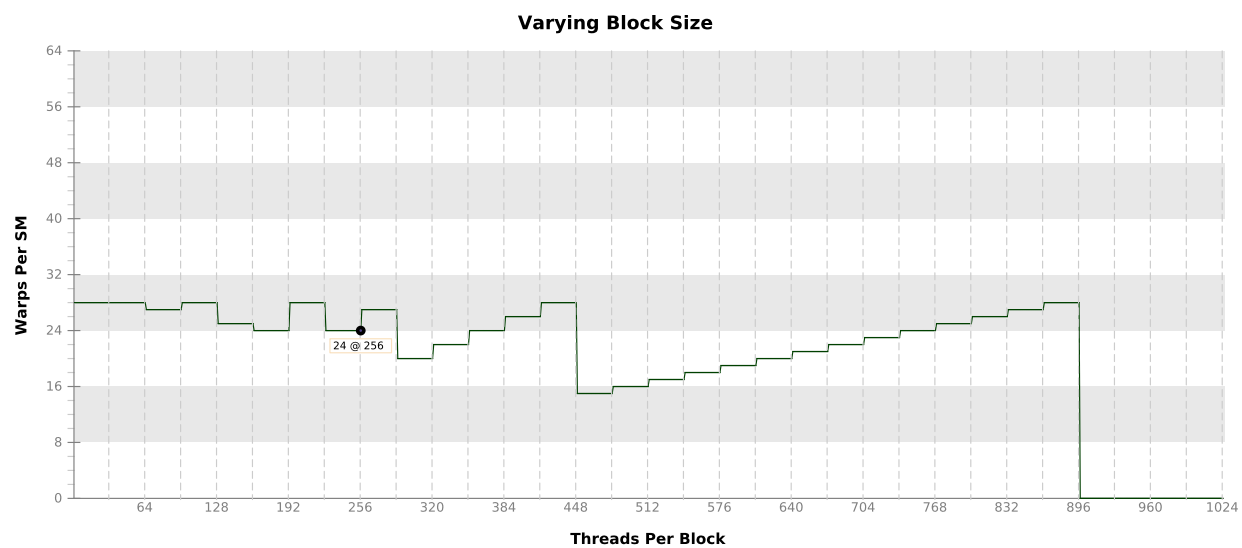## 3.1. GPU Utilization Is Limited By Register Usage

The kernel uses 72 registers for each thread (18432 registers for each block). This register usage is likely preventing the kernel from fully utilizing the GPU. Device "Tesla V100-SXM2-32GB" provides up to 65536 registers for each block. Because the kernel uses 18432 registers for each block each SM is limited to simultaneously executing 3 blocks (24 warps). Chart "Varying Register Count" below shows how changing register usage will change the number of blocks that can execute on each SM.
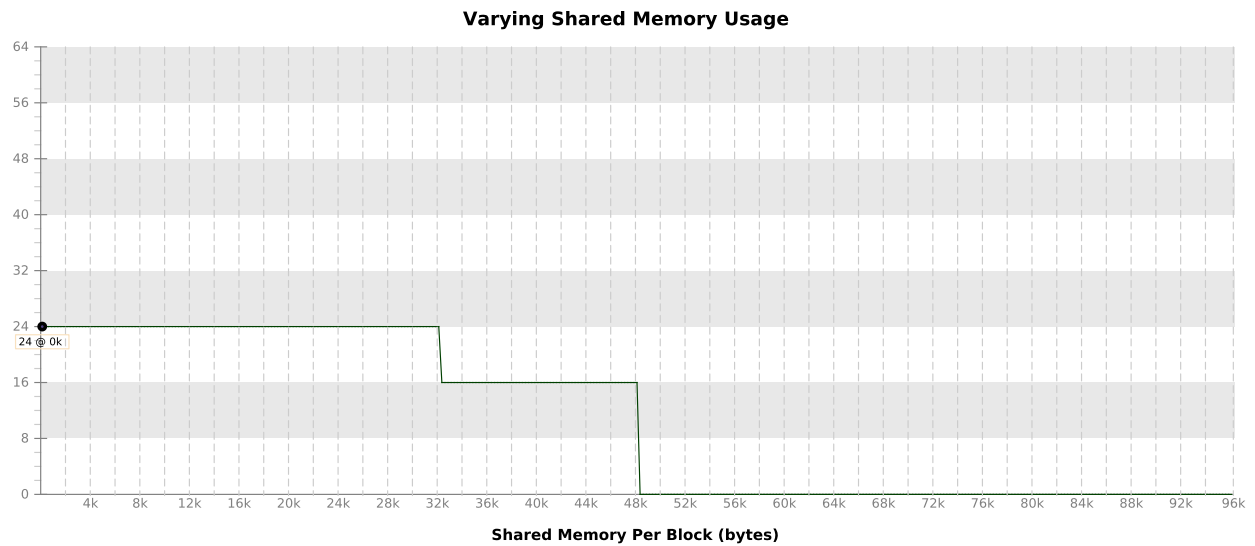
*Optimization: Use the -maxrregcount flag or the __launch_bounds__ qualifier to decrease the number of registers used by each thread. This will increase the number of blocks that can execute on each SM. On devices with Compute Capability 5.2 turning global cache off can increase the occupancy limited by register usage.*

| Variable | Achieved | Theoretical | Device Limit | Grid Size: [ 313,14,1 ] (4382 blocks) Block Size: [ 16,16,1 ] (256 th |
|---|---|---|---|---|
| **Occupancy Per SM** | | | | |
| Active Blocks | | 3 | 32 | |
| Active Warps | 23.61 | 24 | 64 | |
| Active Threads | | 768 | 2048 | |
| Occupancy | 36.9% | 37.5% | 100% | |
| **Warps** | | | | |
| Threads/Block | | 256 | 1024 | |
| Warps/Block | | 8 | 32 | |
| Block Limit | | 8 | 32 | |
| **Registers** | | | | |
| Registers/Thread | | 72 | 65536 | |
| Registers/Block | | 18432 | 65536 | |
| Block Limit | | 3 | 32 | |
| **Shared Memory** | | | | |
| Shared Memory/Block | | 0 | 98304 | |
| Block Limit | | 0 | 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 vs Threads Per Block. Annotation: 24 @ 256

## Varying Register Count



Warps Per SM vs Registers Per Thread. Annotation: 24 @ 72

**Varying Shared Memory Usage**



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.
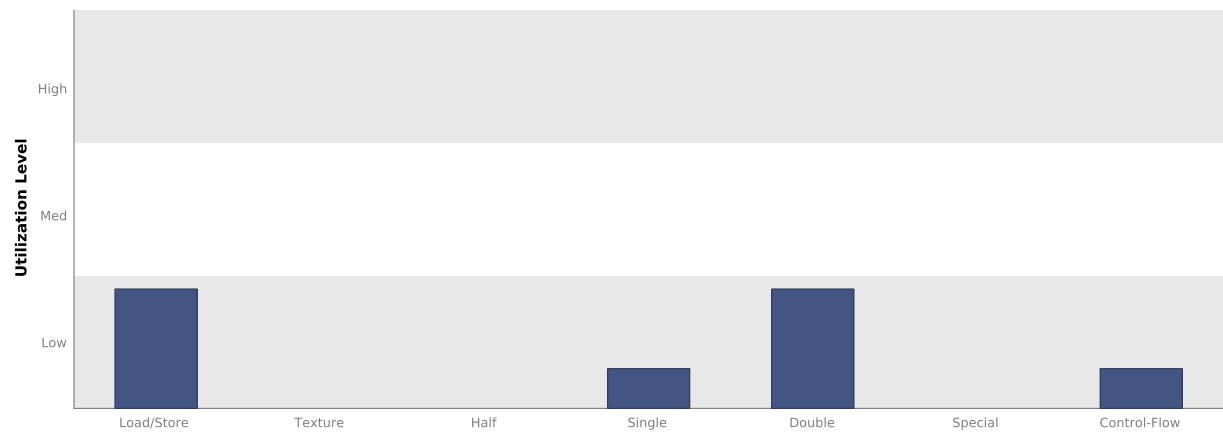
# 4. Compute Resources

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

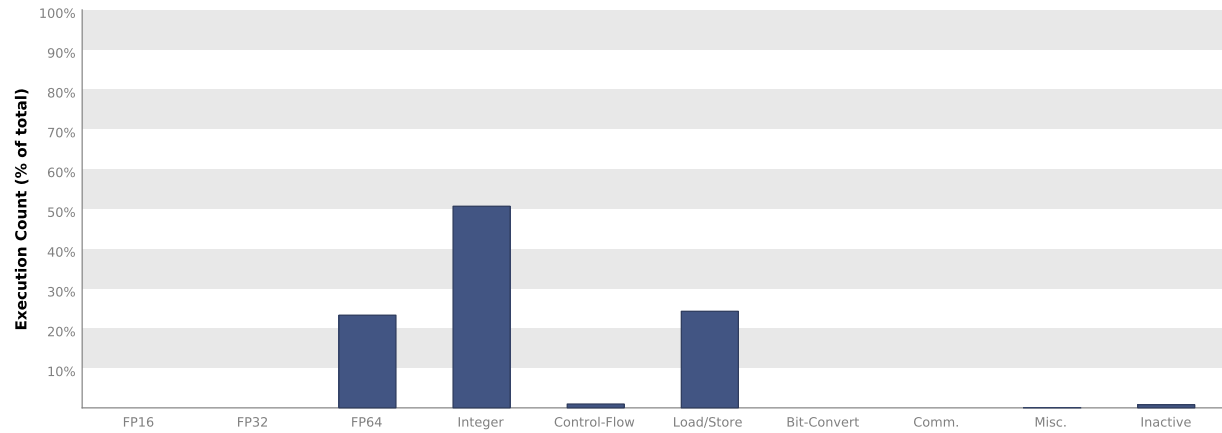## 4.1. Function Unit Utilization

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.
Half - Half-precision floating-point arithmetic instructions.
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.