# Analysis Report

## trans_kernel(float*, float*, int, int)

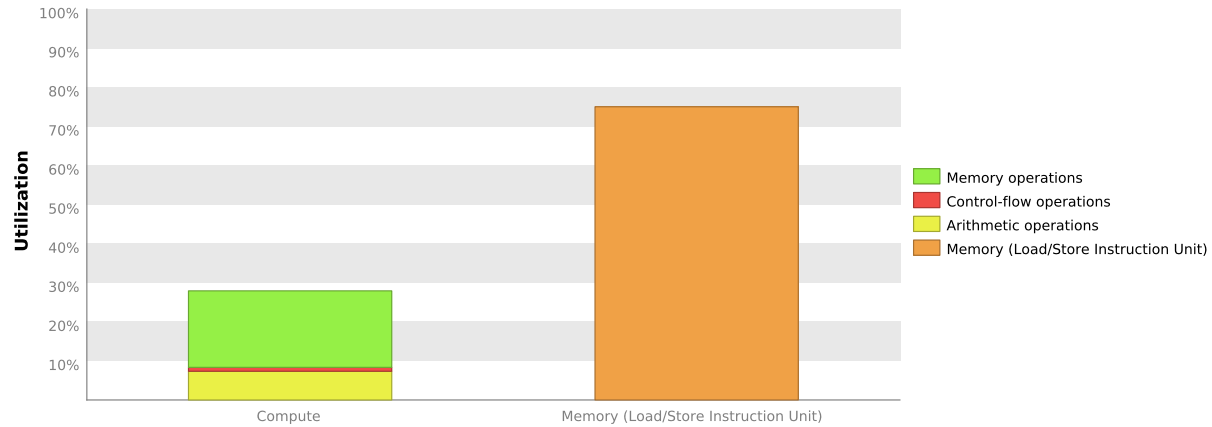| | |
|---|---|
| Duration | 3.326 ms (3,326,472 ns) |
| Grid Size | [ 160,480,1 ] |
| Block Size | [ 16,16,1 ] |
| Registers/Thread | 8 |
| Shared  Memory/Block | 0 B |
| Shared Memory Requested | 16 KiB |
| Shared Memory Executed | 16 KiB |
| Shared Memory Bank Size | 4 B |

| [0] Tesla K20c | |
|---|---|
| GPU UUID | GPU-5af13825-b086-5bfe-f175-a239faf0f1f9 |
| Compute Capability | 3.5 |
| Max. Threads per Block | 1024 |
| Max. Shared Memory per Block | 48 KiB |
| Max. Registers per Block | 65536 |
| 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 | 3.522 TeraFLOP/s |
| Double Precision FLOP/s | 1.174 TeraFLOP/s |
| Number of Multiprocessors | 13 |
| Multiprocessor Clock Rate | 705.5 MHz |
| Concurrent Kernel | true |
| Max IPC | 7 |
| Threads per Warp | 32 |
| Global Memory Bandwidth | 208 GB/s |
| Global Memory Size | 4.687 GiB |
| Constant Memory Size | 64 KiB |
| L2 Cache Size | 1.25 MiB |
| Memcpy Engines | 2 |
| PCIe Generation | 2 |
| PCIe Link Rate | 5 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 "trans_kernel" 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 K20c" 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 load/store instruction units within the multiprocessors.

# 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.

## 2.1. GPU Utilization Is Limited By Memory Instruction Execution

The kernel's performance is potentially limited by the load/store instruction units within the multiprocessors. These units are responsible for executing the instructions that result in accesses to memory. The table below shows the memory bandwidth used by this kernel for the various types of memory on the device.

*Optimization: Examine the compute analysis results for this kernel to determine how to reduce utilization and improve efficiency of the load/store instruction units.*

| Transactions | Bandwidth | Utilization | |
|---|---|---|---|
| **L1/Shared Memory** | | | |
| Local Loads | 0 | 0 B/s | |
| Local Stores | 0 | 0 B/s | |
| Shared Loads | 0 | 0 B/s | |
| Shared Stores | 0 | 0 B/s | |
| Global Loads | 1228800 | 23.314 GB/s | |
| Global Stores | 9830400 | 93.254 GB/s | |
| Atomic | 0 | 0 B/s | |
| L1/Shared Total | 11059200 | 116.568 GB/s | Idle  Low  Medium  High  Max |
| **L2 Cache** | | | |
| L1 Reads | 2457600 | 23.304 GB/s | |
| L1 Writes | 9830436 | 88.557 GB/s | |
| Texture Reads | 0 | 0 B/s | |
| Noncoherent Reads | 0 | 0 B/s | |
| Atomic | 0 | 0 B/s | |
| Total | 12288036 | 111.861 GB/s | Idle  Low  Medium  High  Max |
| **Texture Cache** | | | |
| Reads | 0 | 0 B/s | Idle  Low  Medium  High  Max |
| **Device Memory** | | | |
| Reads | 4599597 | 43.616 GB/s | |
| Writes | 5037127 | 47.692 GB/s | |
| Total | 9636724 | 91.308 GB/s | Idle  Low  Medium  High  Max |
| ECC Overhead | 3871419 | 33.863 GB/s | |
| **System Memory** | | | |
| [ PCIe configuration: Gen2 x16, 5 Gbit/s ] | | | |
| Reads | 0 | 0 B/s | Idle  Low  Medium  High  Max |
| Writes | 1 | 85.376 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 results below indicate that the GPU does not have enough work because instruction execution is stalling excessively.

## 3.1. 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.
 Synchronization - The warp is blocked at a __syncthreads() call.
 Memory Throttle - Large number of pending memory operations prevent further forward progress. These can be reduced by combining several memory transactions into one.
 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.
 Execution Dependency - An input required by the instruction is not yet available. Execution dependency stalls can potentially be reduced by increasing instruction-level parallelism.
 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.
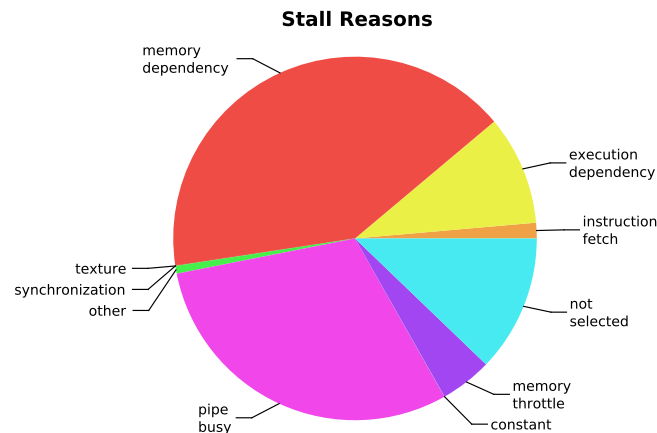 Pipeline Busy - The compute resource(s) required by the instruction is not yet available.
 Instruction Fetch - The next assembly instruction has not yet been fetched.
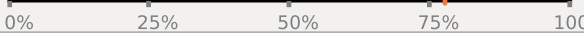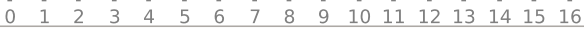 Constant - A constant load is blocked due to a miss in the constants cache.
 Texture - The texture sub-system is fully utilized or has too many outstanding requests.

*Optimization: Resolve the primary stall issue; memory dependency.*
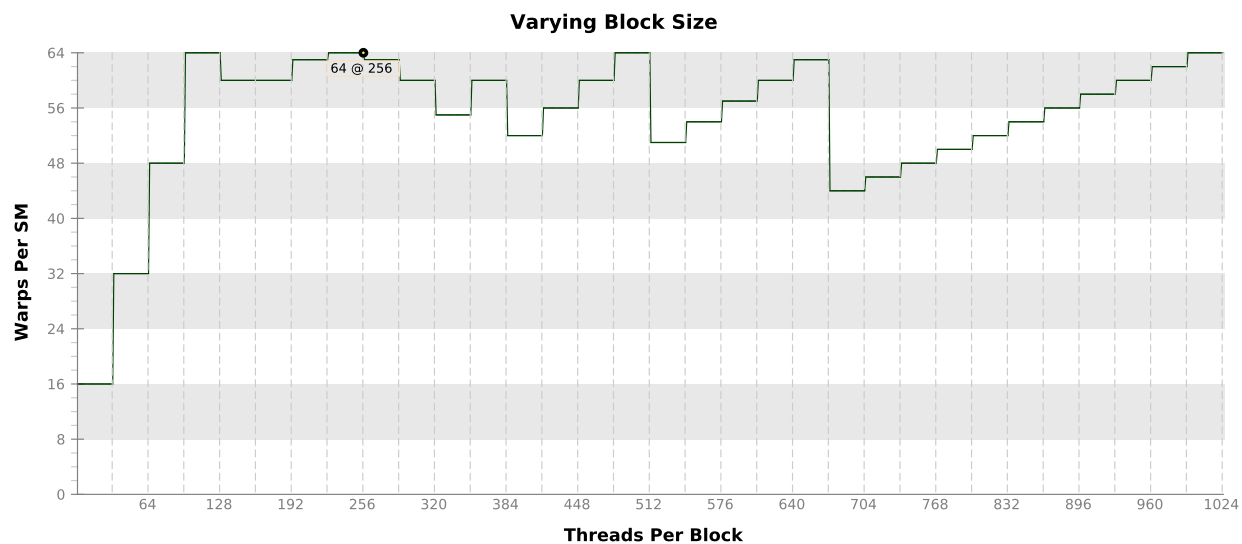


## 3.2. Occupancy Is Not Limiting Kernel Performance

The kernel's block size, register usage, and shared memory usage allow it to fully utilize all warps on the GPU.
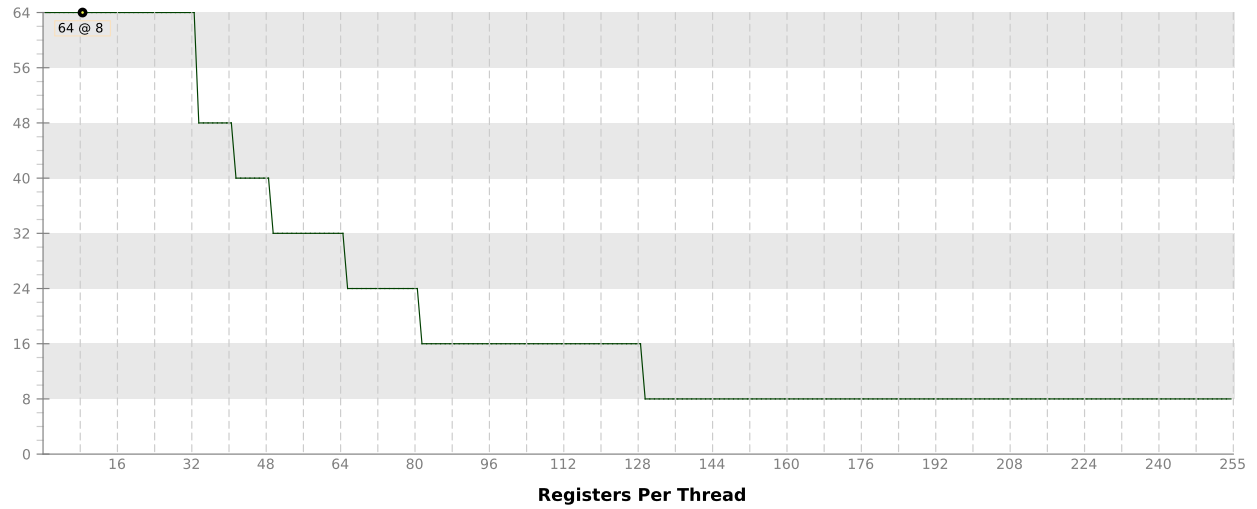
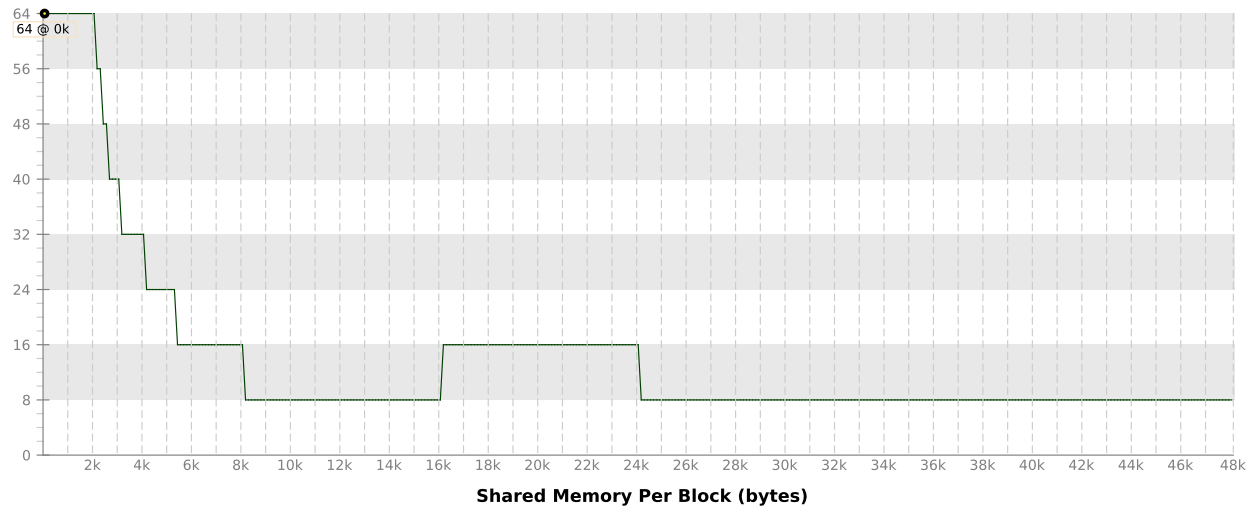| Variable | Achieved | Theoretical | Device Limit | Grid Size: [ 160,480,1 ] (76800 blocks) Block Size: [ 16,16,1 ] (256 |
|---|---|---|---|---|
| Occupancy Per SM | | | | |
| Active Blocks | | 8 | 16 | |
| Active Warps | 49.51 | 64 | 64 | |
| Active Threads | | 2048 | 2048 | |
| Occupancy | 77.4% | 100% | 100% | |
| Warps | | | | |
| Threads/Block | | 256 | 1024 | |
| Warps/Block | | 8 | 32 | |
| Block Limit | | 8 | 16 | |
| Registers | | | | |
| Registers/Thread | | 8 | 255 | |
| Registers/Block | | 2048 | 65536 | |
| Block Limit | | 32 | 16 | |
| Shared Memory | | | | |
| Shared Memory/Block | | 0 | 16384 | |
| Block Limit | | | 16 | |

## 3.3. Occupancy Charts

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



**Varying Block Size**

64 @ 256

*Warps Per SM* (y-axis)

*Threads Per Block* (x-axis)

## Varying Register Count



64 @ 8

Registers Per Thread

## Varying Shared Memory Usage



64 @ 0k

Shared Memory Per Block (bytes)

# 4. Compute Resources

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

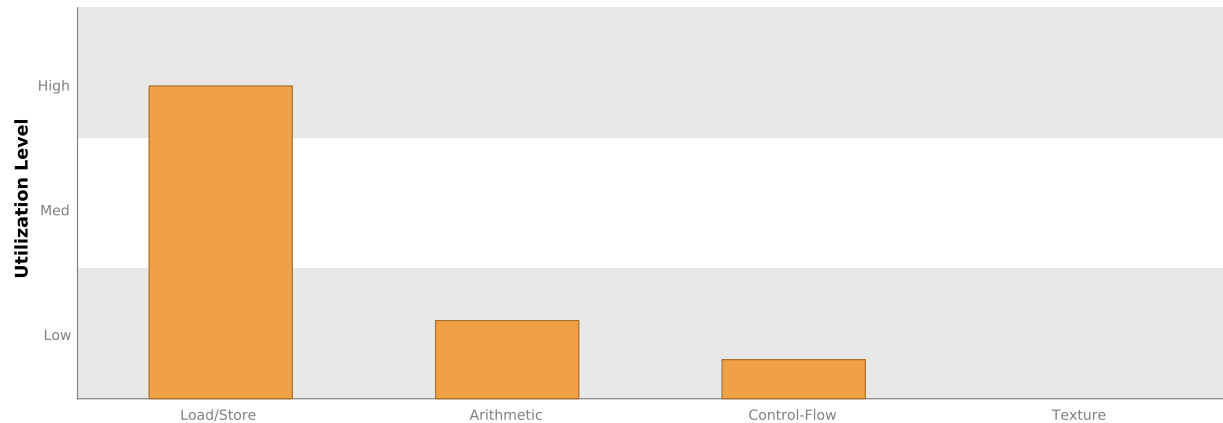## 4.1. Function Unit Utilization

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 local, shared, global, constant, etc. memory.
Arithmetic - All arithmetic instructions including integer and floating-point add and multiply, logical and binary operations, etc.
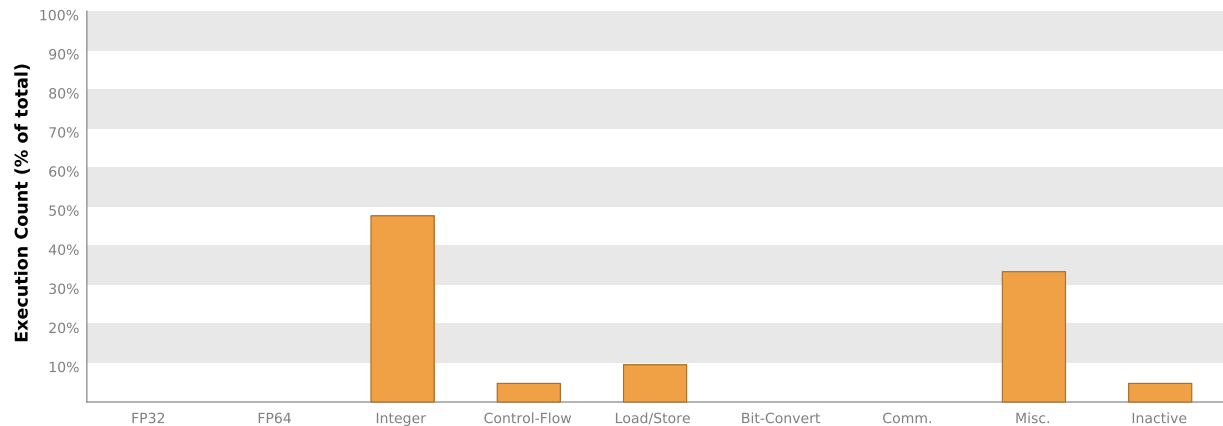Control-Flow - Direct and indirect branches, jumps, and calls.
Texture - Texture operations.



## 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.