# Analysis Report

## quadratic_difference

| | |
|---|---|
| Duration | 2.985 ms (2,984,638 ns) |
| Grid Size | [ 15000,47,1 ] |
| Block Size | [ 2,32,1 ] |
| Registers/Thread | 23 |
| Shared Memory/Block | 560 B |
| Shared Memory Requested | 96 KiB |
| Shared Memory Executed | 96 KiB |
| Shared Memory Bank Size | 4 B |

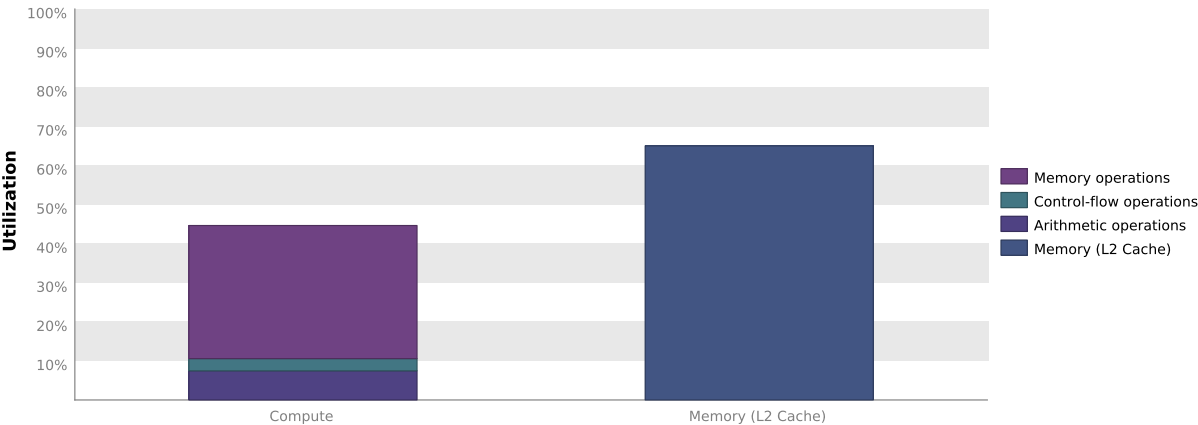| [0] GeForce GTX TITAN X | |
|---|---|
| GPU UUID | GPU-5bcd20eb-5e0d-33cc-a01f-45bb7b69fc82 |
| Compute Capability | 5.2 |
| Max. Threads per Block | 1024 |
| Max. Shared Memory per Block | 48 KiB |
| Max. Registers per Block | 65536 |
| Max. Grid Dimensions | [ 2147483647, 65535, 65535 ] |
| Max. Block Dimensions | [ 1024, 1024, 64 ] |
| Max. Warps per Multiprocessor | 64 |
| Max. Blocks per Multiprocessor | 32 |
| Single Precision FLOP/s | 6.611 TeraFLOP/s |
| Double Precision FLOP/s | 206.592 GigaFLOP/s |
| Number of Multiprocessors | 24 |
| Multiprocessor Clock Rate | 1.076 GHz |
| Concurrent Kernel | true |
| Max IPC | 6 |
| Threads per Warp | 32 |
| Global Memory Bandwidth | 336.48 GB/s |
| Global Memory Size | 12 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. The results below indicate that the performance of kernel "quadratic_difference" 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 "GeForce GTX TITAN X" 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 L2 Cache 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 L2 cache.

## 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.

*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. 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 | 11000256 | 608.159 GB/s | |
| Shared Stores | 13750696 | 760.22 GB/s | |
| Shared Total | 24750952 | 1,368.379 GB/s | Idle　Low　Medium　High　Max |
| **L2 Cache** | | | |
| Reads | 34310927 | 474.228 GB/s | |
| Writes | 5055090 | 69.869 GB/s | |
| Total | 39366017 | 544.097 GB/s | Idle　Low　Medium　High　Max |
| **Unified Cache** | | | |
| Local Loads | 0 | 0 B/s | |
| Local Stores | 0 | 0 B/s | |
| Global Loads | 74118192 | 474.222 GB/s | |
| Global Stores | 5055084 | 69.869 GB/s | |
| Texture Reads | 46742816 | 646.055 GB/s | |
| Unified Total | 125916092 | 1,190.146 GB/s | Idle　Low　Medium　High　Max |
| **Device Memory** | | | |
| Reads | 21209 | 293.14 MB/s | |
| Writes | 2672569 | 36.939 GB/s | |
| Total | 2693778 | 37.232 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 | 69.107 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.
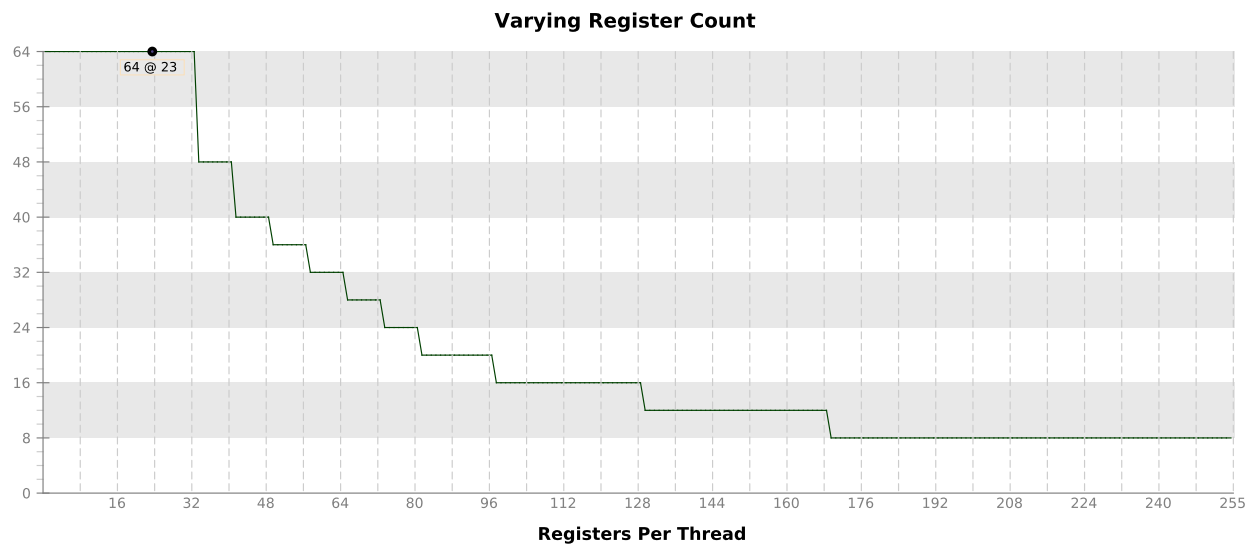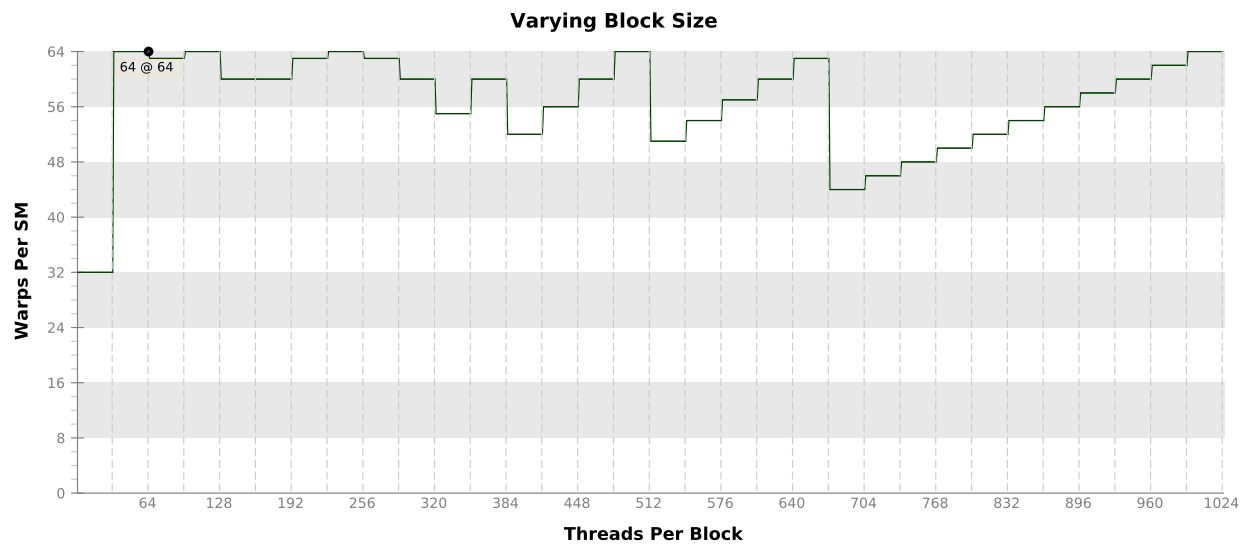
## 3.1. Occupancy Is Not Limiting Kernel Performance

The kernel's block size, register usage, and shared memory usage allow it to fully utilize all warps on the GPU.
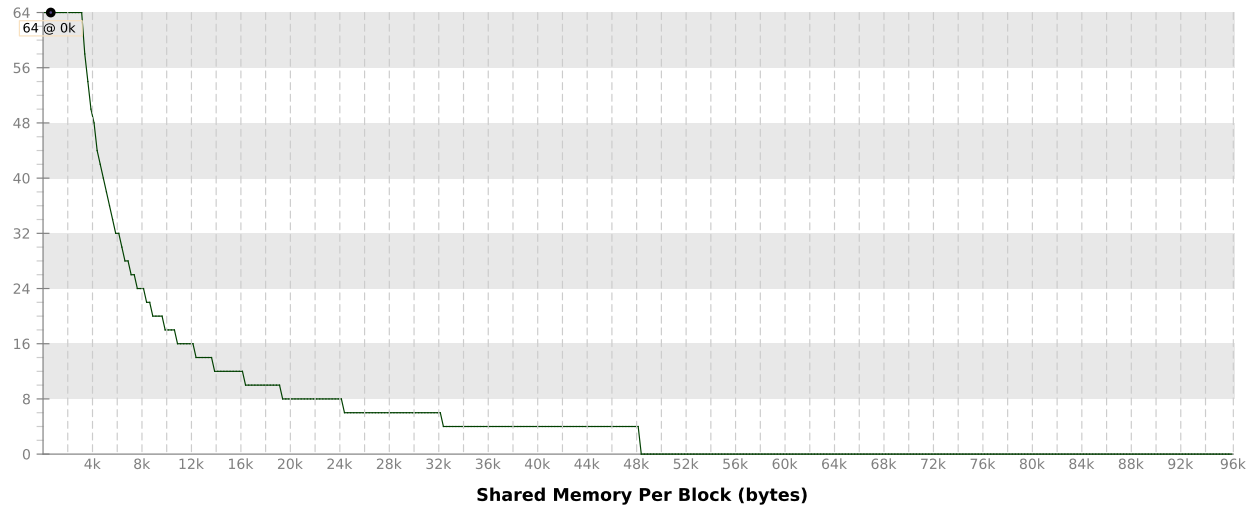
| Variable | Achieved | Theoretical | Device Limit | Grid Size: [ 15000,47,1 ] (705000 blocks) Block Size: [ 2 |
|---|---|---|---|---|
| **Occupancy Per SM** | | | | |
| Active Blocks | | 32 | 32 | |
| Active Warps | 57.71 | 64 | 64 | |
| Active Threads | | 2048 | 2048 | |
| Occupancy | 90.2% | 100% | 100% | |
| **Warps** | | | | |
| Threads/Block | | 64 | 1024 | |
| Warps/Block | | 2 | 32 | |
| Block Limit | | 32 | 32 | |
| **Registers** | | | | |
| Registers/Thread | | 23 | 255 | |
| Registers/Block | | 1536 | 65536 | |
| Block Limit | | 42 | 32 | |
| **Shared Memory** | | | | |
| Shared Memory/Block | | 560 | 98304 | |
| Block Limit | | 128 | 32 | |

## 3.2. Occupancy Charts

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

## Varying Block Size



Y-axis: **Warps Per SM** — 0, 8, 16, 24, 32, 40, 48, 56, 64

X-axis: **Threads Per Block** — 64, 128, 192, 256, 320, 384, 448, 512, 576, 640, 704, 768, 832, 896, 960, 1024

64 @ 64

## Varying Register Count



Y-axis: 0, 8, 16, 24, 32, 40, 48, 56, 64

X-axis: **Registers Per Thread** — 16, 32, 48, 64, 80, 96, 112, 128, 144, 160, 176, 192, 208, 224, 240, 255

64 @ 23

## Varying Shared Memory Usage



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. Compute resources are used most efficiently when all threads in a warp have the same branching and predication behavior. The results below indicate that a significant fraction of the available compute performance is being wasted because branch and predication behavior is differing for threads within a warp.

## 4.1. Kernel Profile - Instruction Execution

The Kernel Profile - Instruction Execution shows the execution count, inactive threads, and predicated threads for each source and assembly line of the kernel. Using this information you can pinpoint portions of your kernel that are making inefficient use of compute resource due to divergence and predication.

*Examine portions of the kernel that have high execution counts and inactive or predicated threads to identify optimization opportunities.*

Cuda Fuctions :

quadratic_difference

Maximum instruction execution count in assembly: 1410000
Average instruction execution count in assembly: 1156138
Instructions executed for the kernel: 124862941
Thread instructions executed for the kernel: 3566422111
Non-predicated thread instructions executed for the kernel: 3081119656
Warp non-predicated execution efficiency of the kernel: 77.1%
Warp execution efficiency of the kernel: 89.3%

## 4.2. Low Warp Execution Efficiency

Warp execution efficiency is the average percentage of active threads in each executed warp. Increasing warp execution efficiency will increase utilization of the GPU's compute resources. The kernel's warp execution efficiency of 77.4% is less than 100% due to divergent branches and predicated instructions. If predicated instructions are not taken into account the warp execution efficiency for these kernels is 89.4%.

*Optimization: Reduce the amount of intra-warp divergence and predication in the kernel.*

## 4.3. Divergent Branches

Compute resource are used most efficiently when all threads in a warp have the same branching behavior. When this does not occur the branch is said to be divergent. Divergent branches lower warp execution efficiency which leads to inefficient use of the GPU's compute resources.

*Optimization: Each entry below points to a divergent branch within the kernel. For each branch reduce the amount of intra-warp divergence.*

## 4.4. 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.

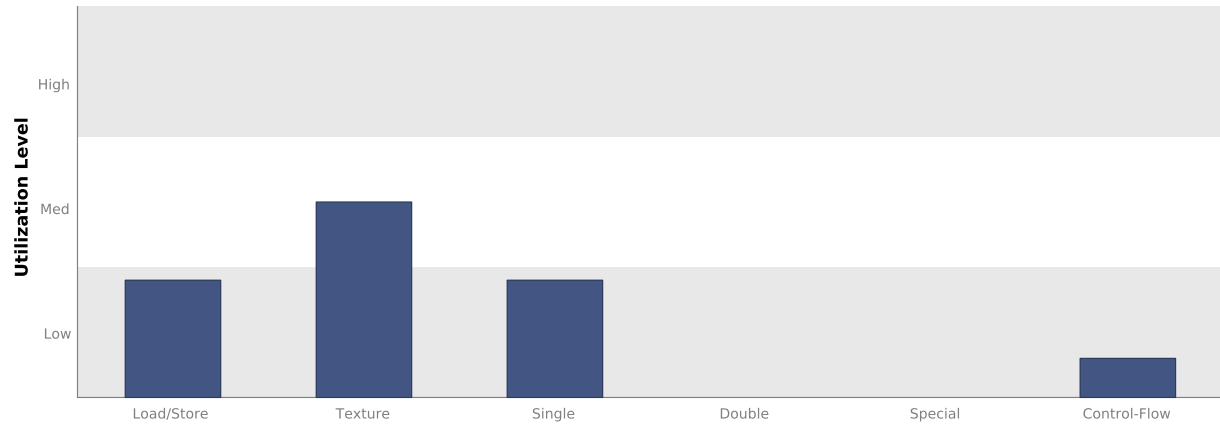Load/Store - Load and store instructions for shared and constant memory.
Texture - Load and store instructions for local, global, and texture memory.
Single - Single-precision integer and floating-point arithmetic instructions.
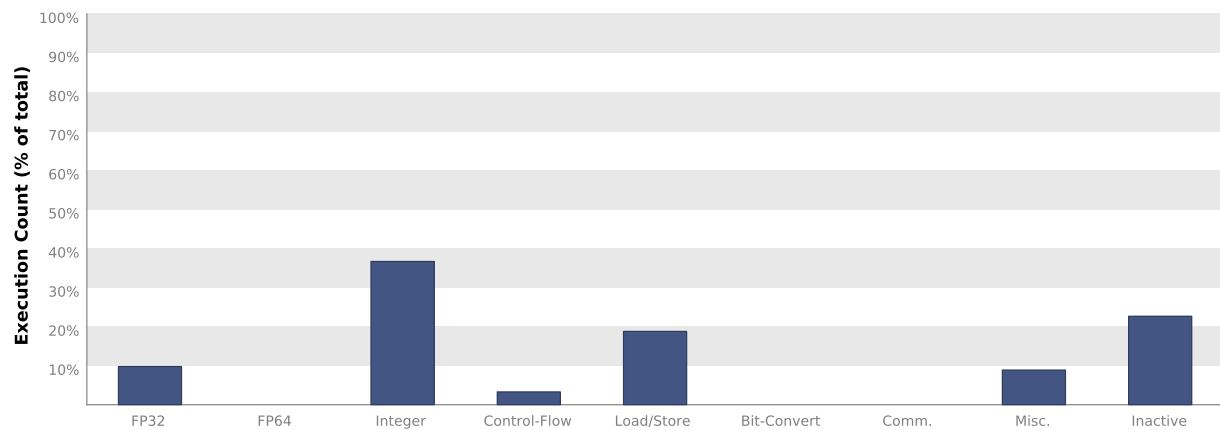Double - Double-precision floating-point arithmetic instructions.
Special - Special arithmetic instructions such as sin, cos, popc, etc.
Control-Flow - Direct and indirect branches, jumps, and calls.

## 4.5. 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.



## 4.6. 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.