

Analysis Report

matrix_mult_kernel_optimised(int*, int*, int*)

Duration	724.031 μ s
Grid Size	[256,256,1]
Block Size	[16,16,1]
Registers/Thread	22
Shared Memory/Block	2 KiB
Shared Memory Requested	48 KiB
Shared Memory Executed	48 KiB
Shared Memory Bank Size	4 B

[0] GRID K520

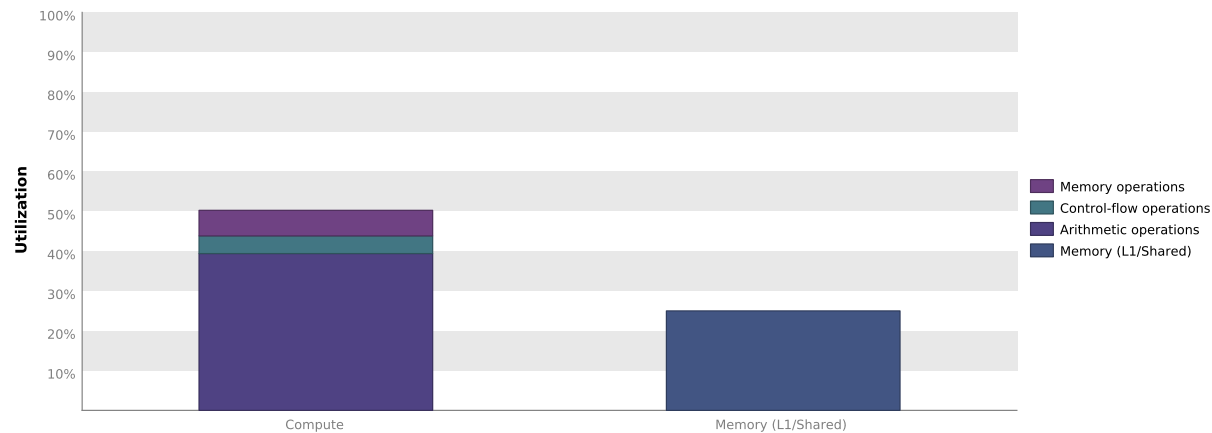
GPU UUID	GPU-b4ee72d2-b156-889f-cccc-dc6aa4a5a894
Compute Capability	3.0
Max. Threads per Block	1024
Max. Shared Memory per Block	48 KiB
Max. Registers per Block	65536
Max. Grid Dimensions	[2147483647, 65535, 65535]
Max. Block Dimensions	[1024, 1024, 64]
Max. Warps per Multiprocessor	64
Max. Blocks per Multiprocessor	16
Single Precision FLOP/s	2.448 TeraFLOP/s
Double Precision FLOP/s	102.016 GigaFLOP/s
Number of Multiprocessors	8
Multiprocessor Clock Rate	797 MHz
Concurrent Kernel	true
Max IPC	7
Threads per Warp	32
Global Memory Bandwidth	160 GB/s
Global Memory Size	4 GiB
Constant Memory Size	64 KiB
L2 Cache Size	512 KiB
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 "matrix_mult_kernel_optimised" is most likely limited by instruction and memory latency. You should first examine the information in the "Instruction And Memory Latency" section to determine how it is limiting performance.

1.1. Kernel Performance Is Bound By Instruction And Memory Latency

This kernel exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of "GRID K520". These utilization levels indicate that the performance of the kernel is most likely limited by the latency of arithmetic or memory operations. Achieved compute throughput and/or memory bandwidth below 60% of peak typically indicates latency issues.




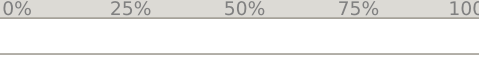


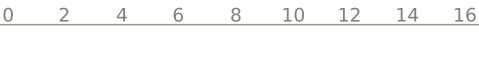


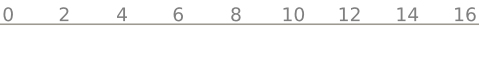




2. Instruction and Memory Latency

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.

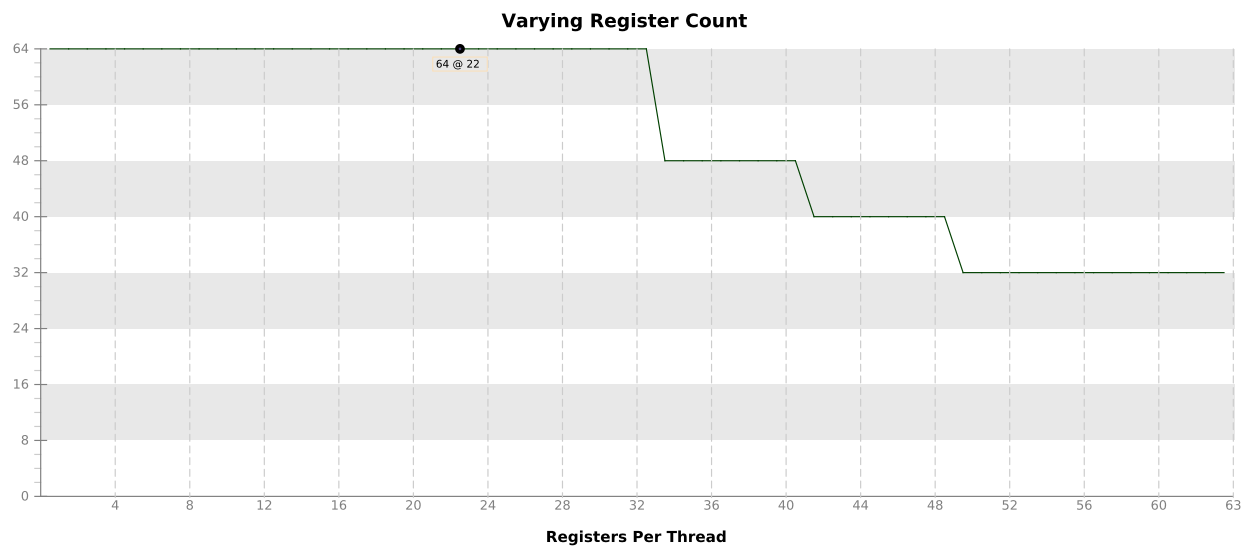
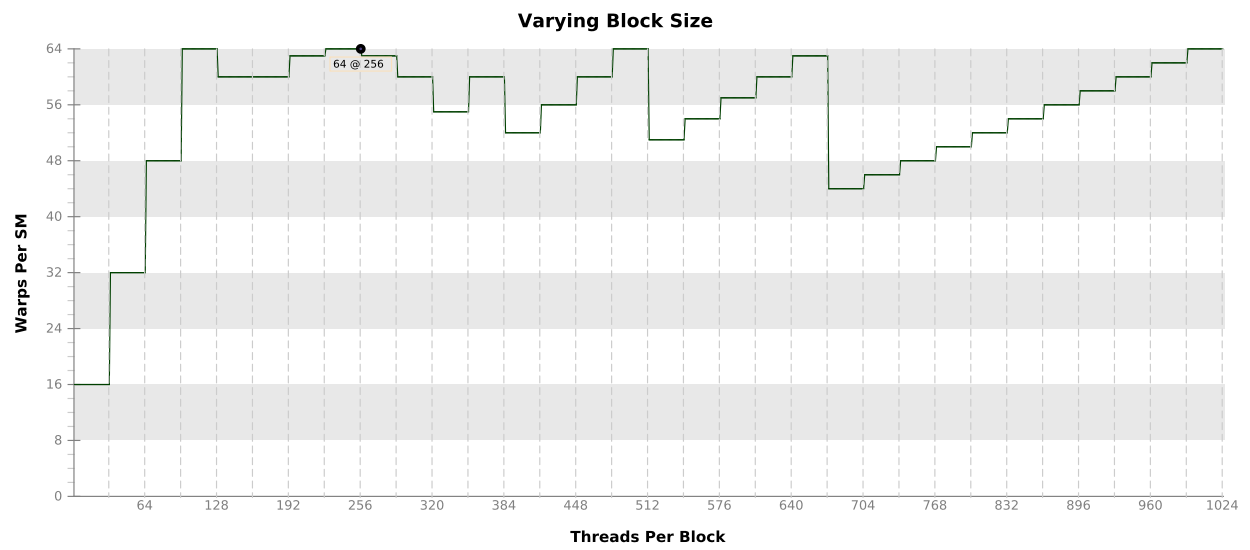
2.1. Occupancy Is Not Limiting Kernel Performance

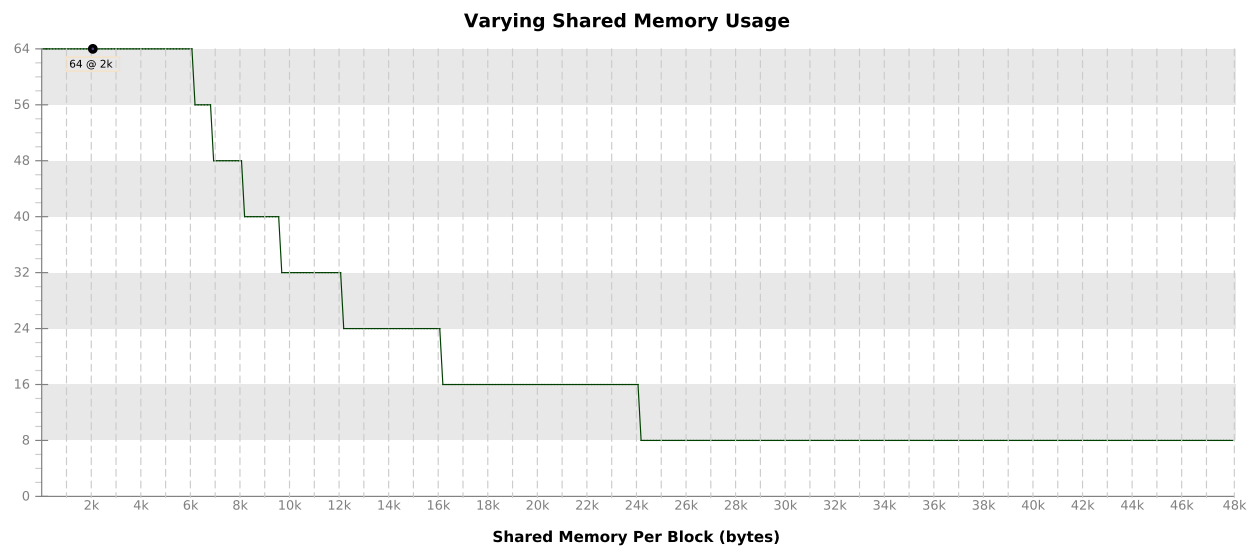
The kernel's block size, register usage, and shared memory usage allow it to fully utilize all warps on the GPU.

Variable	Achieved	Theoretical	Device Limit	Grid Size: [256,256,1] (65536 blocks) Block Size: [16,16,1]
Occupancy Per SM				
Active Blocks		8	16	
Active Warps	44.46	64	64	
Active Threads		2048	2048	
Occupancy	69.5%	100%	100%	
Warps				
Threads/Block		256	1024	
Warps/Block		8	32	
Block Limit		8	16	
Registers				
Registers/Thread		22	63	
Registers/Block		6144	65536	
Block Limit		10	16	
Shared Memory				
Shared Memory/Block		2048	49152	
Block Limit		24	16	

2.2. Occupancy Charts

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





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

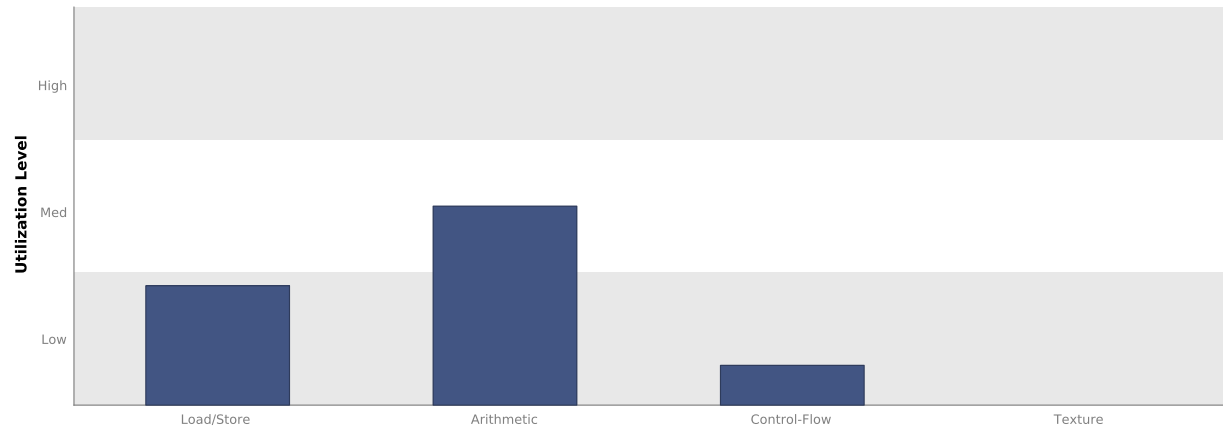
Different types of instructions are executed on different function units within each SM. Performance can be limited if a function unit is over-used by the instructions executed by the kernel. The following results show that the kernel's performance is not limited by overuse of any function unit.

Load/Store - Load and store instructions for local, shared, global, constant, etc. memory.

Arithmetic - All arithmetic instructions including integer and floating-point add and multiply, logical and binary operations, etc.

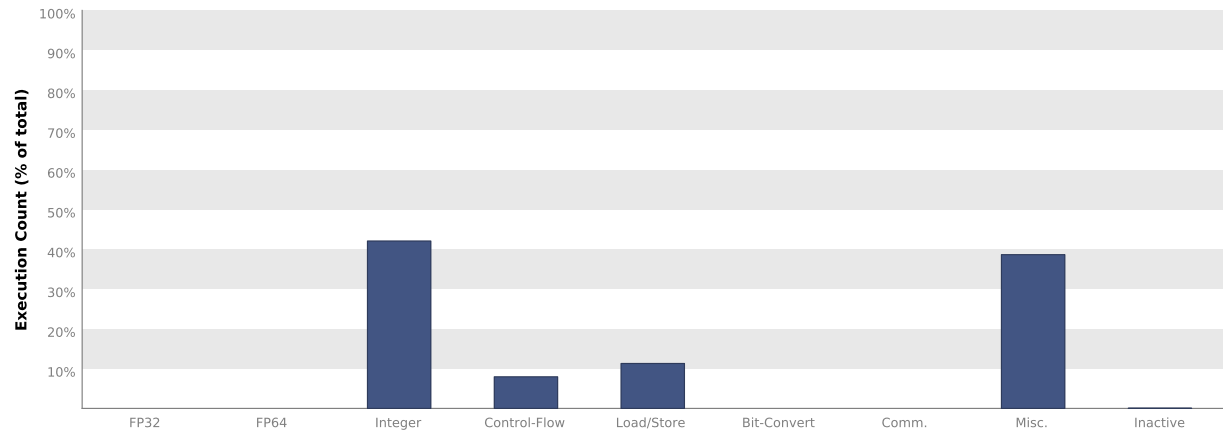
Control-Flow - Direct and indirect branches, jumps, and calls.

Texture - Texture operations.



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.

4.1. Memory Bandwidth And Utilization

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.

Transactions	Bandwidth	Utilization	
L1/Shared Memory			
Local Loads	0	0 B/s	
Local Stores	0	0 B/s	
Shared Loads	786432	362.312 GB/s	
Shared Stores	66738	30.746 GB/s	
Global Loads	131072	15.096 GB/s	
Global Stores	4096	471.76 MB/s	
Atomic	0	0 B/s	
L1/Shared Total	988338	408.626 GB/s	
L2 Cache			
L1 Reads	262144	15.096 GB/s	
L1 Writes	8416	484.66 MB/s	
Texture Reads	0	0 B/s	
Atomic	0	0 B/s	
Total	270560	15.581 GB/s	
Texture Cache			
Reads	0	0 B/s	
Device Memory			
Reads	8199	472.163 MB/s	
Writes	8192	471.76 MB/s	
Total	16391	943.924 MB/s	
System Memory			
[PCIe configuration: Gen3 x16, 8 Gbit/s]			
Reads	0	0 B/s	
Writes	3	172.763 kB/s	