# Lab 6: Advanced Topic in CUDA
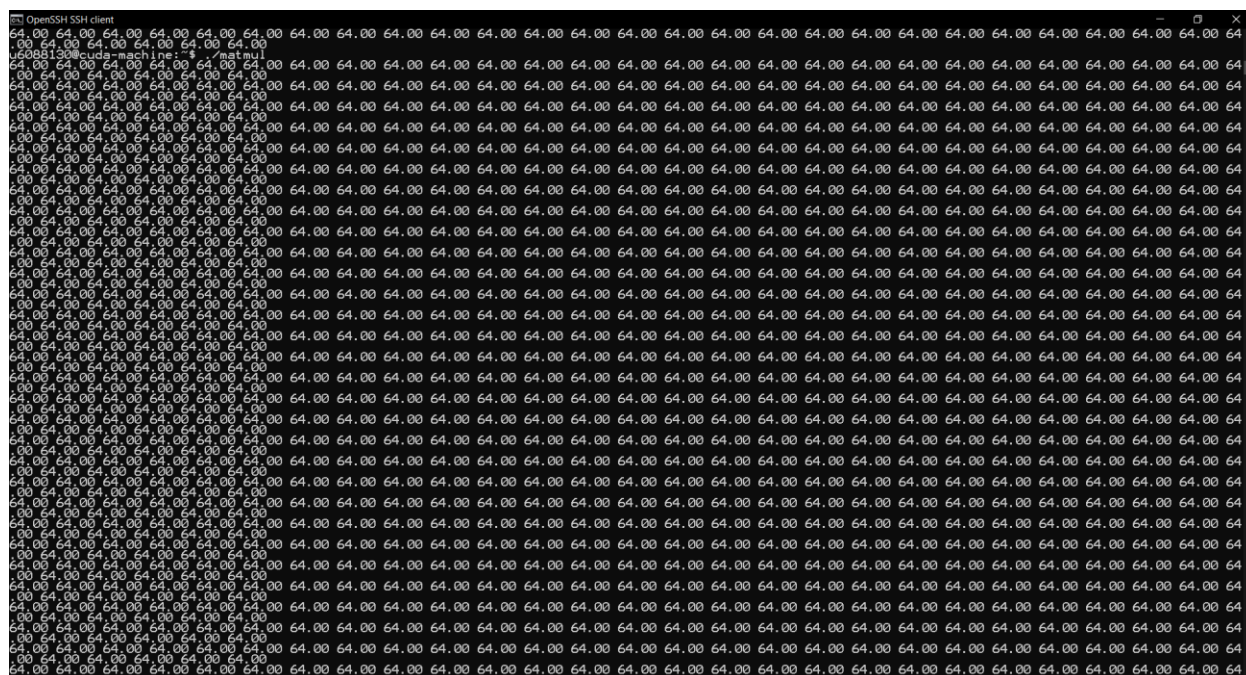
Name:  Sunat Praphanwong                     ID: 6088130            Sec:    1

**Lab: 1.** Result from "matmul.cu"



**Lab: 2.** Result from "matmul_shared.cu"

**Lab: 3.** Result from measuring performance of both "matmul.cu" and "matmul_shared.cu", with Width is set to 512

**Lab: 4.** CUDA Occupancy calculation for "matmul.cu" and "matmul_shared.cu"

Download CUDA Occupancy Calculator from https://docs.nvidia.com/cuda/cuda-occupancy-calculator/CUDA_Occupancy_Calculator.xls

**4.1 "matmul.cu"**

Threads per block = _____256_____

Registers used per thread = _____32_____

Shared memory used per thread block = __0_____

**Screenshot of the result from NVIDIA Occupancy Calculator:**

## 4.2 "matmul_shared.cu"

Threads per block = _____256_____
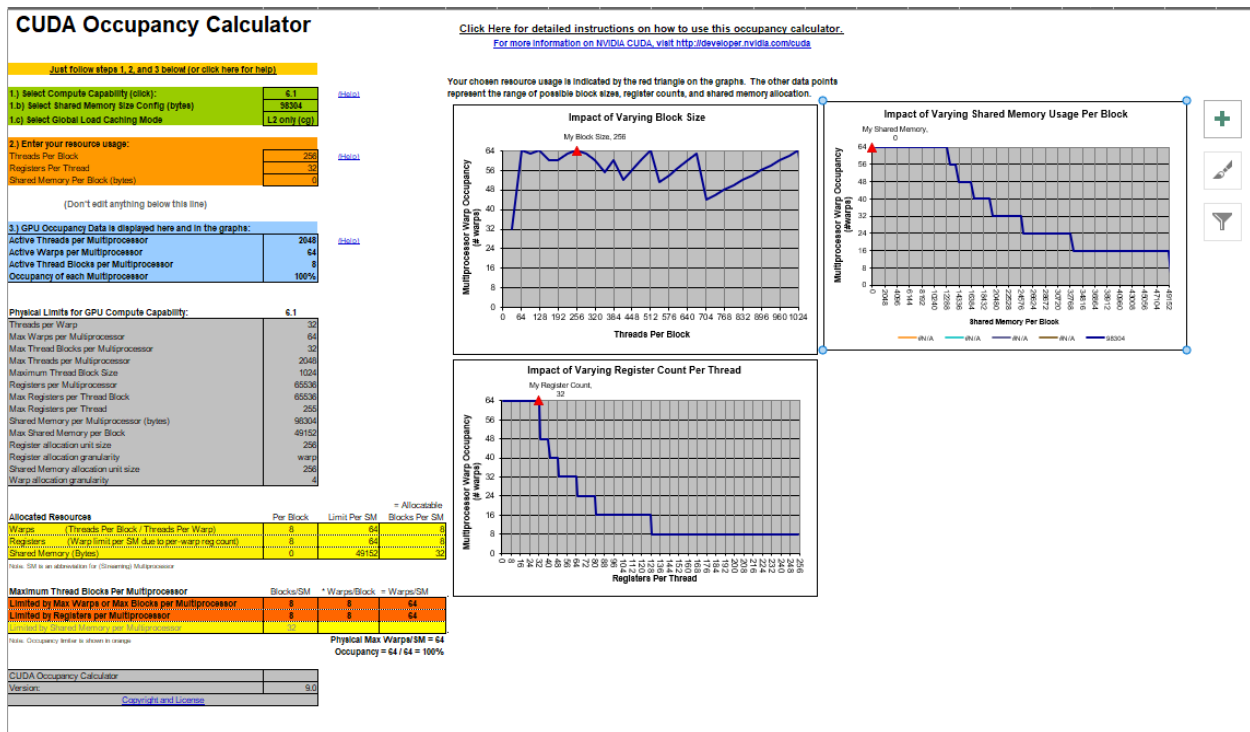
Registers used per thread = _____25_____

Shared memory used per thread block = _2048__

**Screenshot of the result from NVIDIA Occupancy Calculator:**

# CUDA Occupancy Calculator

| 1.) Select Compute Capability (click): | 6.1 | (Help) |
|---|---|---|
| 1.b) Select Shared Memory Size Config (bytes) | 98304 | |
| 1.c) Select Global Load Caching Mode | L2 only (cg) | |

| 2.) Enter your resource usage: | | |
|---|---|---|
| Threads Per Block | 256 | (Help) |
| Registers Per Thread | 25 | |
| Shared Memory Per Block (bytes) | 2048 | |

(Don't edit anything below this line)

| 3.) GPU Occupancy Data is displayed here and in the graphs: | | |
|---|---|---|
| Active Threads per Multiprocessor | 2048 | (Help) |
| Active Warps per Multiprocessor | 64 | |
| Active Thread Blocks per Multiprocessor | 8 | |
| Occupancy of each Multiprocessor | 100% | |

| Physical Limits for GPU Compute Capability: | 6.1 |
|---|---|
| Threads per Warp | 32 |
| Max Warps per Multiprocessor | 64 |
| Max Thread Blocks per Multiprocessor | 32 |
| Max Threads per Multiprocessor | 2048 |
| Maximum Thread Block Size | 1024 |
| Registers per Multiprocessor | 65536 |
| Max Registers per Thread Block | 65536 |
| Max Registers per Thread | 255 |
| Shared Memory per Multiprocessor (bytes) | 98304 |
| Max Shared Memory per Block | 49152 |
| Register allocation unit size | 256 |
| Register allocation granularity | warp |
| Shared Memory allocation unit size | 256 |
| Warp allocation granularity | 4 |

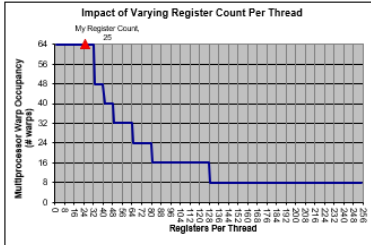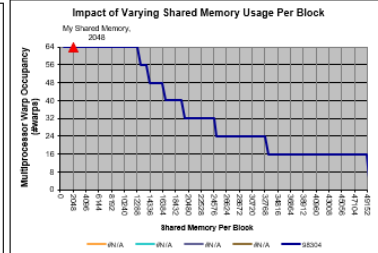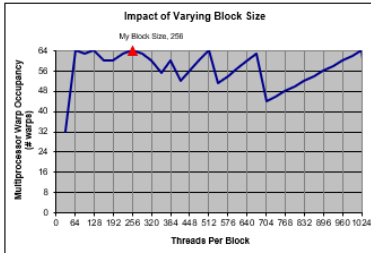| Allocated Resources | | Per Block | Limit Per SM | = Allocatable Blocks Per SM |
|---|---|---|---|---|
| Warps | (Threads Per Block / Threads Per Warp) | 8 | 64 | 8 |
| Registers | (Warp limit per SM due to per-warp reg count) | 8 | 64 | 8 |
| Shared Memory (Bytes) | | 2048 | 49152 | 48 |

Note: SM is an abbreviation for (Streaming) Multiprocessor

| Maximum Thread Blocks Per Multiprocessor | Blocks/SM | * Warps/Block | = Warps/SM |
|---|---|---|---|
| Limited by Max Warps or Max Blocks per Multiprocessor | 8 | 8 | 64 |
| Limited by Registers per Multiprocessor | 8 | 8 | 64 |
| Limited by Shared Memory per Multiprocessor | 48 | | |

Note: Occupancy limiter is shown in orange

Physical Max Warps/SM = 64
Occupancy = 64 / 64 = 100%

| CUDA Occupancy Calculator | |
|---|---|
| Version: | 9.0 |
| Copyright and License | |

Your chosen resource usage is indicated by the red triangle on the graphs. The other data points represent the range of possible block sizes, register counts, and shared memory allocation.

### Impact of Varying Block Size

My Block Size, 256

X-axis: Threads Per Block
Y-axis: Multiprocessor Warp Occupancy (# warps)

### Impact of Varying Shared Memory Usage Per Block

My Shared Memory, 2048

X-axis: Shared Memory Per Block
Y-axis: Multiprocessor Warp Occupancy (# warps)

Legend: #N/A #N/A #N/A #N/A 98304

### Impact of Varying Register Count Per Thread

My Register Count, 25

X-axis: Registers Per Thread
Y-axis: Multiprocessor Warp Occupancy (# warps)

**Do not forget to include the source files into the zip file before submission.**