# Assignment 4
# CUDA C/C++ Acceleration

Arnav Goyal - 251244778
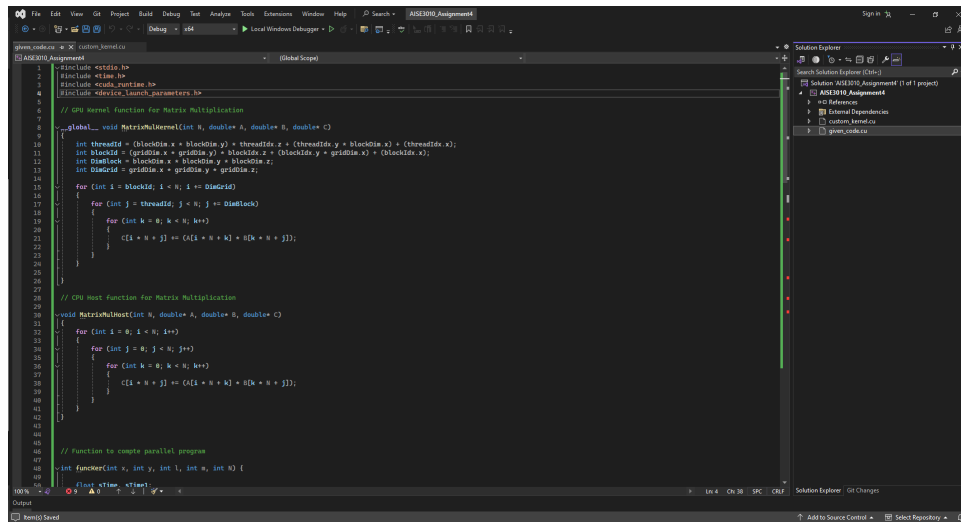
Graeme Watt - 251201502

*April 8, 2024*

# CUDA Environment

My CUDA environment was installed following the docs on NVIDIA's website. Here are my specs:

- Windows 11 Home 22H2
- GTX 1660-Ti 6 GB - Driver Version 552.12
- CUDA Compiler Driver 12.4
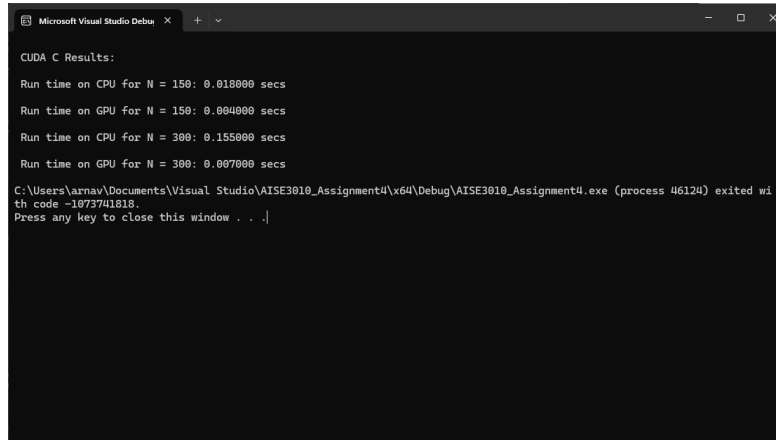- Microsoft Visual Studio 2022 Community Edition (Version 17.9.5)

# The Given Code

The `ipynb` file was uploaded onto colab, and the CUDA script was copied and pasted into my visual studio solution. It was built and ran without any modifications (other than including the neccessary `cuda_runtime.h` and `device_launch_parameters.h` header files), Im assuming that this is what *running smoothly means...* The given code produces the following console output.



**Image 1:** Given code in my CUDA environment

**Image 2:** Given code's output

The given code executes the kernel with the following parameters:

- a 3-dim grid size of 16, 16, 1
- a 3-dim thread block size of 16, 8, 1

The kernel uses these parameters to linearize the blocks and threads, to come up with a unique blockId and threadId for each parallel execution. The kernel then iterates over blocks, then iterates over threads) and performs the dot product of the row in A and the column in B and uses the result as an index of the matrix multiplication result C.

A big advantage of the code is that it is quite parallelized, and it basically calculates C through a whole 'rod' of the resultant 'cube' in each thread. This is better than performing a less parallelized matrix multiplication, for example writing a kernel that performs one matrix multiplication, and sequentially looping over it.

A big disadvantage of the code is that it is quite inefficient, it has time complexity of $O(n^3)$ (it has 3 nested for-loops), this makes it extremely inefficient and quite slow to larger and larger inputs. My solution is slightly better at $O(n^2)$

# A Generalized Solution

In order to generalize the code I started from scratch and created a new kernel, It works similarly to the given kernel but has some modifications, this kernel can be found in the `custom_kernel.cu` file. It works with any general matrix sizes for A and B that are compatible for the operation of multiplication. I tested its execution time for the following (given) scenarios:

i) At least 100 matrix multiplications of $A \times B$, for $A \in \mathbb{R}^{500 \times 500}$ and $B \in \mathbb{R}^{500 \times 400}$

ii) At least 5000 matrix multiplications of $A \times B$, for $A \in \mathbb{R}^{50 \times 20}$ and $B \in \mathbb{R}^{20 \times 50}$

iii) At least 1000 matrix multiplications of $A \times B$, for $A \in \mathbb{R}^{6 \times 4000}$ and $B \in \mathbb{R}^{4000 \times 9}$

I also can examine the usage of my GPU with `Task Manager - Performance Monitor` and the `nvidia-smi -l 1` command. These two methods give the following results. Notice that the GPU utlization spikes to 100% from close to idle during build and output of my CUDA program.
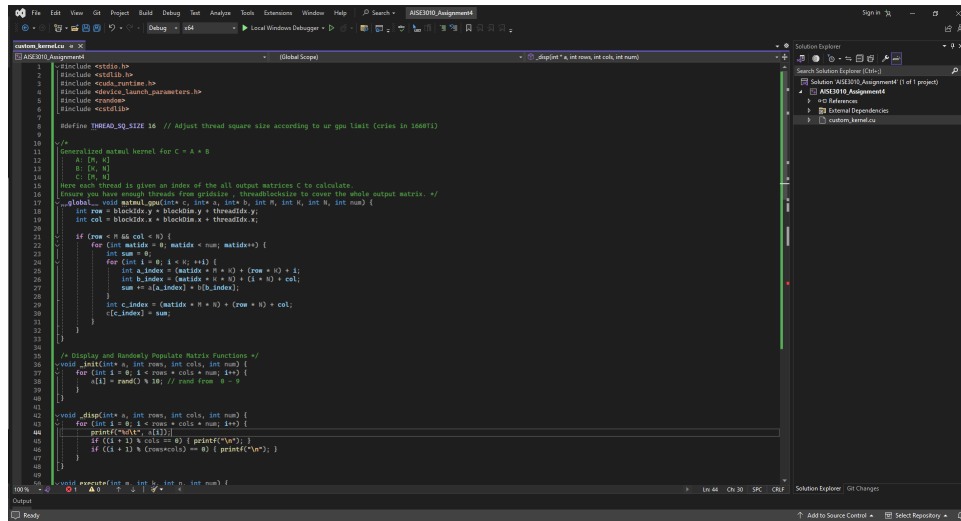


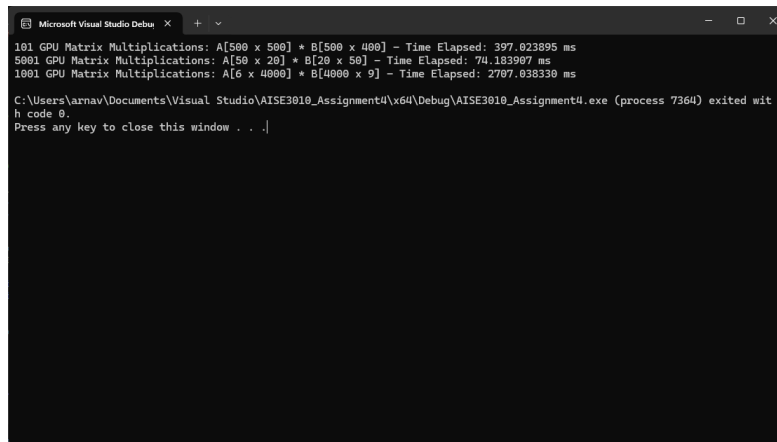**Image 3:** My CUDA program in the VS environment



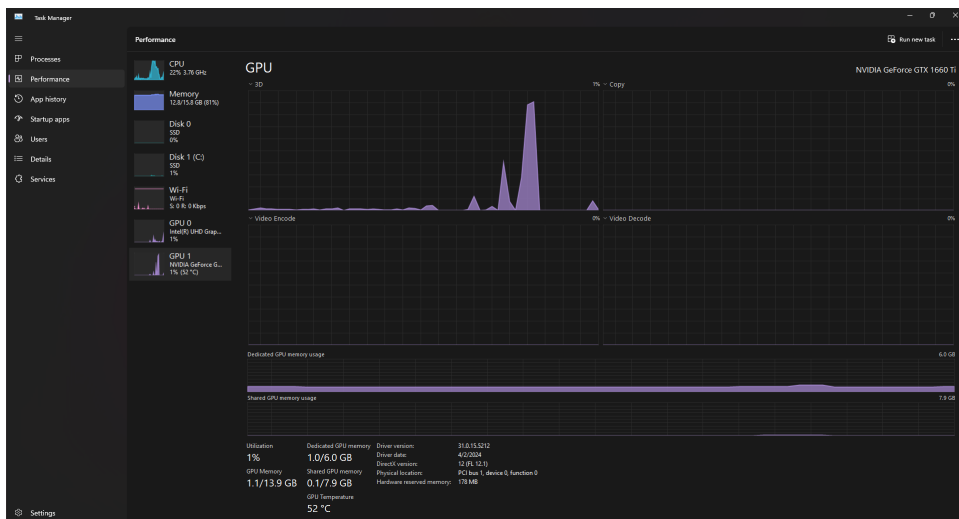**Image 4:** output of my CUDA program on all three scenarios

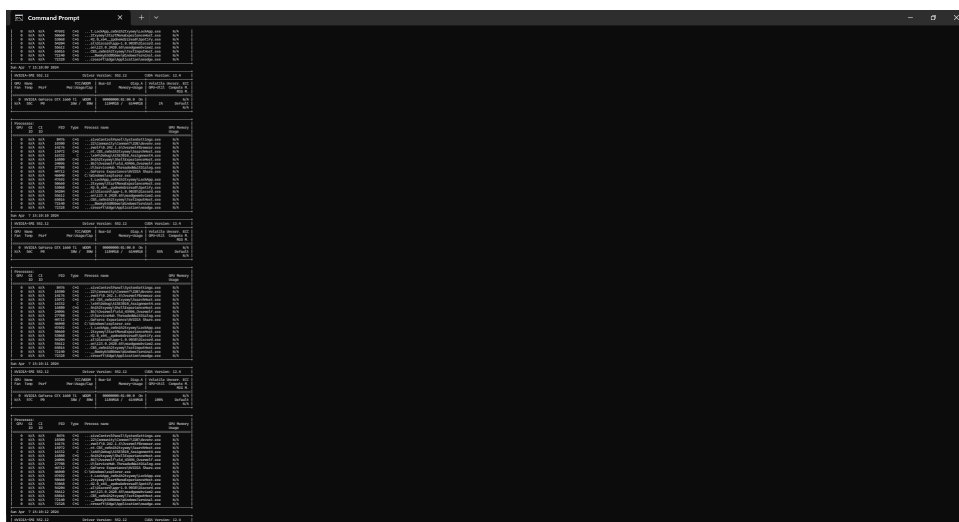**Image 5:** GPU Utiliization spike in Task Manager's Performance Monitor



**Image 6:** nvidia-smi command output looped every 1 second (you might need to zoom in)