## Machine Problem 2 Report

Presented to Prof. Ahmad Afsahi ELEC 374: Digital Systems Engineering Faculty of Applied Science Queen's University

Nathan Goodman: 20228249

Date: April 4<sup>th</sup>, 2023

## **CUDA Code**

For the three different kernel configurations requested, three different CUDA scripts were created:

1. Each thread only produces a single output element and 16x16 threadblocks.

```
#include "cuda_runtime.h"

#include <string.h>

#include <stdio.h>

#include <stdib.h>

#include <tdib.h>

#include <time.h>
 #define debug 0 //If debug = 1, certain print statements will be enabled
 //2. Write a kernel that has each thread producing one output matrix. Kernel config should be 16x16 thread blocks
 /* MatrixAddition Hernel.

Parameters: pointer to output matrix C, two Pointers to input matricies A and B, dimensions of matricies A and B (remember they're square matriceies so this can be single int) */
*/
global_ void matrixAddition_hernel(float*d_s, float*d_b, float*d_c, int siseOfMatricies) {
    int Row = blockIdx.y*blockDim.y + threadIdx.y;
    int Col = blockIdx.x*blockDim.x + threadIdx.x;
    int idx = Row*siseOfMatricies + Col;
    //0 1 2 4
    //5 6 7 8 9
    if (idx < siseOfMatricies*siseOfMatricies) { //Avoid accessing beyond end off matricies
    d_c[idx] = d_a[idx] + d_b[idx];
    }
 /* Created a matrixAddition function. Should basically be the same as the kernel function.

This function will be used to check whether the kernel function created correct output (used to check)
 */
woid matrixAddition(float *a, float *b, float *c, int siseOfMatricies) { //Note acc
for (int i = 0; i < siseOfMatricies; i++) {
    for (int j = 0; j < siseOfMatricies; j++) {
        *(c + i*siseOfMatricies + j) = *(a + i*siseOfMatricies + j) + *(b + i*siseOfMatricies + j);
        //C[i][j] = A[i][j] + B[i][j];
}</pre>
return 1:
//Additional function to help with debugging:
void printMatrix(float *a, int sice) {
  for (int i = 0; i < size; i++) {
    for (int j = 0; j < size; j++) {
        printf("%f ", *(a + i*size + j));
        }
}</pre>
                printf("\n");
//Get Device Name (They did this in tutorial so I'll do it here)
cudaGetDeviceProperties(ideviceProps, 0);
printf("CUDA device [%s]\n", deviceProps.name);
printf("NNumber of Multiprocessors: %d\n", deviceProps.multiProcessorCount);
printf("NNumber of Multiprocessors: %d\n", deviceProps.mathreadsPerBlock);
printf("NAMar Dimension of a Block: %d\n", deviceProps.mathreadsDim);
printf("NAMar Dimension of a Block: %d\n", deviceProps.mathreadsDim);
printf("NAMar Dimension of a Grid: %d\n", deviceProps.mathreadsDim);
         . 1. Define two square input matricies \boldsymbol{\lambda} and \boldsymbol{B}, and matching output matrix Note: they're floats
         const int dimOfMatricies = 250: //Value determines size of matricies Ex. dimOfMatricies = 5 will result in 5x5 matricies
         //lb. Calculate amount of memory they take:
         int mbytes = dimOfMatricies*dimOfMatricies*simeof(float);
printf("Dimensions of Matricies: %dx%d\n", dimOfMatricies, dimOfMatricies);
          //lc. Allocate host memory for matricies:
         //lc. Allocate nost memory ava maveler.
float 'a = 0;
float 'b = 0;
float 'b = 0;
float 'e = 0;
cudsMallocHost((void*) &a, nbytes); //Allocates host memory for matrix A, and points pointer a to first value.
cudsMallocHost((void*) &b, nbytes);
cudsMallocHost((void*) &b, nbytes);
```

```
if (debug == 1) {
   printf("Memory Locations of matricies:\n");
   printf("\ta = %x\n", sa);
   printf("\ta = %x\n", sb);
   printf("\ta = %x\n", sb);
                      //ld. Store input matricies into memory
                      srand(time(NULL));
                       fillMatrix(a, dimOfMatricies);
fillMatrix(b, dimOfMatricies);
                      if (debug -- 1) (
                           printf("Matrix A:\n");
printfatrix(a, dimOfMatricles);
printf("Matrix B:\n");
                            printMatrix(b, dimOfMatricies);
                      //le. Allocate device memory for matricies: float *d_a = 0; float *d_b = 0;
                       float *d c = 0;
                      cudaMalloc((void**)&d_a, nbytes); //Allocates memory for matrix A, and points pointer a to first value.

cudaMemset(d_a, 255, nbytes); //Sets all allocated bytes to 255 (they did this in tutorial so i did it here)

cudaMalloc((void**)&d_b, nbytes);
                      cudaMemset(d_b, 255, nbytes);
cudaMalloc((void**)&d_c, nbytes);
                       cudaMemset (d c, 255, nbytes);
                       //Set kernel launch configuration
                      //Set kernel launch configuration
int blckWidth = 16; //block width and kength
int threadsPerBlock = blckWidth*blckWidth;
int threadsNeeded = dimOfMatricles*dimOfMatricles; //Because in this configuration one thread only produces one value
int numBlocks = threadsNeeded / threadsPerBlock;
if (numBlocks < 1) numBlocks+;</pre>
                      if (debug == 1) {
    printf("Block Width: %d\n", blckWidth);
    printf("Grid Width: %d\n", numBlocks);
142
143
144
                      dim3 dimBlocks - dim3(blckWidth, blckWidth);//asks for 16X16-256 thread blocks
                                                                                   //256 threads per block layed out in 16x16
                      dim3 dimGrid = dim3(numBlocks, numBlocks);
148
149
150
151
                      cudaEvent t start, stop;
                       cudaEventCreate(&start);
                       cudaEventCreate(&stop);
                       cudaDeviceSynchronize():
                       float gpu_time = 0.0f;
                      //asynchronously issue work to the GPU (all stream 0)
                        udaEventRecord(start, 0);
                       //Copy inputs to device
                       cudaMemcpyAsync(d_a, a, nbytes, cudaMemcpyHostToDevice, 0); cudaMemcpyAsync(d_b, b, nbytes, cudaMemcpyHostToDevice, 0);
                       //Call Kernel
                       matrixAddition_kernel << dimGrid, dimBlocks, 0, 0 >> <(d_a, d_b, d_c, dimOfMatricies); //Note: only section after >>> is the actual function parameters
                      //Copy outputs from device cudaMemcpyAsync(c, d_c, nbytes, cudaMemcpyDeviceToHost, 0);
                      cudaEventRecord(stop, 0);
cudaEventSynchronize(stop); //stop is updated here
cudaEventElapsedTime(&gpu_time, start, stop);
                       //print the CPU times
                       printf("time spent executing by the GPU: %.2f\n", gpu_time);
                       //Calculate matrixAddition using CPU:
                     183
184
185
186
187
188
                      if (debug -- 1) {
   printf("Matrix C:\n");
                            printMatrix(c, dimOfMatricies);
printf("Correct Matrix:\n");
                            printMatrix(d, dimOfMatricies);
190
191
192
193
194
195
196
197
198
                      //release resources
                      cudaEventDestroy(start);
cudaEventDestroy(stop);
                      cudaFreeHost(a);
                       cudaFreeHost(b);
cudaFreeHost(c);
                       cudaFreeHost (d);
                      cudaFree(d_a);
cudaFree(d_b);
                      cudaFree(d c);
                       cudaDeviceReset():
```

2. Each thread produces an output row and 16 threads per block.

Note: Due to scripts for the most part being identical, only changed parts are shown. Full file is found within attached .zip

```
/* 3. Create a matrixAddition Kernel where each thread calculates a row of output values.
Parameters: pointer to output matrix C, two Pointers to input matricies A and B,
dimensions of matricies A and B (remember they're square matriceies so this can be single int)
        _ void matrixAddition_kernel(float*d_a, float*d_b, float*d_c, int sizeOfMatricies) {
   int idx = blockIdx.x*blockDim.x + threadIdx.x;
   idx = idx*sizeOfMatricies; //To account for already computed indicies
                           //Ex. thread 1 (idx=0) will compute output elements 0 (5*0),1(5*0+1),2,3,4
                           //
                              thread 2 (idx=1) will compute output elements 5 (5*1),6 (5*1+1),7,8,9
   int n = sizeOfMatricies*sizeOfMatricies;
   if (idx < n) { //Avoid accessing beyond end off matricies
       for (int i = 0; i < sizeOfMatricies; i++) {
          //Each thread will calculate a row of output matrix:
          d_c[idx + i] = d_a[idx + i] + d_b[idx + i];
 //Set kernel launch configuration
 int blckWidth = 4; //block width and length (16 threads per block)
 int threadsPerBlock = blckWidth*blckWidth;
int threadsNeeded = dimOfMatricies*dimOfMatricies;
 int numBlocks = threadsNeeded / threadsPerBlock;
 if (numBlocks < 1) numBlocks++;
 dim3 \ dimBlocks = dim3(16, 1);
                                                           // 16 threads per block
 dim3 dimGrid = dim3(numBlocks, 1);
   3. Each thread produces a column of output and 16 threads per block.
```

```
/* 4. Create a matrixAddition Kernel where each thread calculates a column of output values..
Parameters: pointer to output matrix C, two Pointers to input matricies A and B,
dimensions of matricies A and B (remember they're square matriceies so this can be single int)
__global
        void matrixAddition_kernel(float*d_a, float*d_b, float*d_c, int sizeOfMatricies) {
   int idx = blockIdx.x*blockDim.x + threadIdx.x;
                                                   //Assume block
   int n = sizeOfMatricies*sizeOfMatricies;
   if (idx < n) { //Avoid accessing beyond end off matricies
      for (int i = 0; i < sizeOfMatricies; i++) {
          //Each thread will calculate a column of output matrix:
          d_c[idx + sizeOfMatricies*i] = d_a[idx + sizeOfMatricies*i] + d_b[idx + sizeOfMatricies*i];
 //Set kernel launch configuration
 int blckWidth = 4; //block width and length (16 threads per block)
 int threadsPerBlock = blckWidth*blckWidth;
int threadsNeeded = dimOfMatricies*dimOfMatricies;
 int numBlocks = threadsNeeded / threadsPerBlock;
if (numBlocks < 1) numBlocks++;
dim3 \ dimBlocks = dim3(16, 1);
                                                        // 16 threads per block
dim3 dimGrid = dim3(numBlocks, 1);
```

## Output

Following guidelines from the Machine Problem 2 Document, the dimOfMatricies variable (equivalent to BLOCK WIDTH from the lecture slides) was varied and outputs were recorded.

Unfortunately for whatever reason, when attempting to input a value of 500 or greater (500, 1000, or 2000), into the configuration for the third question, the script failed to calculate GPU execution time despite no compilation errors:

```
CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 14351708

Max Dimension of a Grid: 14351720

Dimensions of Matricies: 500x500

time spent executing by the GPU: 0.00

Press any key to continue . . . _
```

The cause of this issue was unable to be determined, as everything in third configuration is identical to the second (except the matrixAddition function), and the script works fine on smaller matrices.

As such, instead of the values 500x500, 1000x1000, and 2000x2000; the values: 60x60, 175x175, 300x300 were added in their place. The outputs for these matrix sizes for the third kernel configuration can be viewed below (under number 3). The outputs from top to bottom are: 60x60, 125x125, 175x,175, 250x250, 300x300:

1. Each thread only produces a single output element and 16x16 threadblocks.

```
CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 10288532

Max Dimension of a Grid: 10288544

Dimensions of Matricies: 125x125

time spent executing by the GPU: 0.27

Test PASSED

Press any key to continue . . . _

CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 5635188

Max Dimension of a Grid: 5635200

Dimensions of Matricies: 250x250

time spent executing by the GPU: 1.44

Test PASSED

Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 12319008

Max Dimension of a Grid: 12319020

Dimensions of Matricies: 500x500

time spent executing by the GPU: 20.66

Test PASSED

Press any key to continue . . . _

CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 2489040

Max Dimension of a Grid: 2489052

Dimensions of Matricies: 1000x1000

time spent executing by the GPU: 320.10

Test PASSED

Press any key to continue . . . _

CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 11074152

Max Dimension of a Grid: 11074164

Dimensions of Matricies: 2000x2000

time spent executing by the GPU: 5226.32

Test PASSED
```

2. Each thread produces an output row and 16 threads per block.

Press any key to continue  $\dots$ 

```
CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 8058740

Max Dimension of a Grid: 8058752

Dimensions of Matricies: 125x125

time spent executing by the GPU: 0.30

Test PASSED

Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 9829500

Max Dimension of a Grid: 9829512

Dimensions of Matricies: 250x250

time spent executing by the GPU: 0.65

Test PASSED

Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 9631800

Max Dimension of a Grid: 9631812

Dimensions of Matricies: 500x500

time spent executing by the GPU: 1.51

Test PASSED

Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 12844380

Max Dimension of a Grid: 12844392

Dimensions of Matricies: 1000x1000

time spent executing by the GPU: 6.30

Test PASSED

Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 4520320

Max Dimension of a Grid: 4520332

Dimensions of Matricies: 2000x2000

time spent executing by the GPU: 36.17

Press any key to continue . . . _
```

3. Each thread produces a column of output and 16 threads per block.

```
CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 6157800

Max Dimension of a Grid: 6157812

time spent executing by the GPU: 0.19

Test PASSED

Press any key to continue . . . _

CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 13106548

Max Dimension of a Grid: 13106560

time spent executing by the GPU: 0.87

Test PASSED

Press any key to continue . . . _

CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 8124572

Max Dimension of a Grid: 8124584

time spent executing by the GPU: 2.30

Test PASSED

Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 17169232

Max Dimension of a Grid: 17169244

time spent executing by the GPU: 6.11

Test PASSED

Press any key to continue . . . _

CUDA device [Tesla C2075]

Number of Mulitprocessors: 14

Max Threads Per Block: 1024

Max Dimension of a Block: 3930852

Max Dimension of a Grid: 3930864

time spent executing by the GPU: 13.41

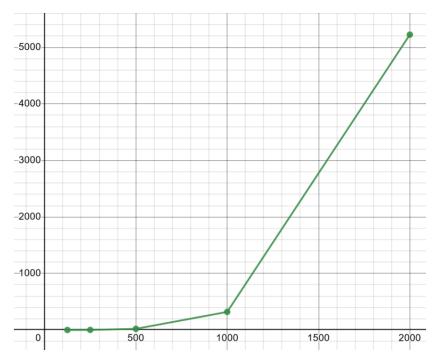
Test PASSED

Press any key to continue . . . _
```

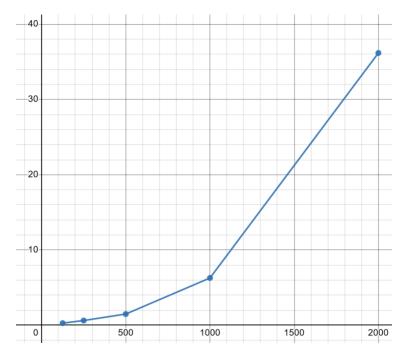
## **Analysis**

The following graphs were generated for each of the launch configurations with the x-axis being the width of the matrix (in floats 4 bytes) and the y-axis being the GPU execution time (in ms):

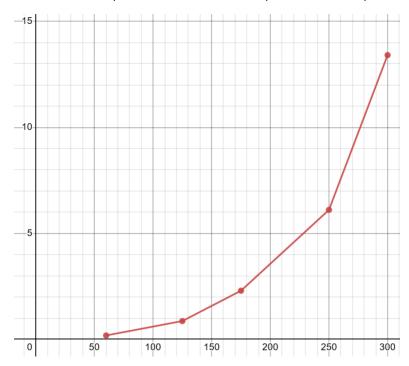
1. Each thread produces a single output element and 16x16 threadblocks:



2. Each thread produces a row of output and threads per block is 16:



3. Each thread produces a column of output and threads per block is 16:



As you can see above, as the width of the matrix grows, the execution time appears to increase exponentially. This makes sense after all the number of output elements or calculations performed is quadratically related to the width of the output matrix—elements to be calculated = matrix width \* matrix width.

From the graphs above it is clear to see that the execution times for the 'row of output per thread and threads per block is 16' were the best, with them still being extremely low even as matrix width

increased. It is expected that this outperformed the 'column of output' matrix addition as there were less operations being performed during each loop iteration (within device load)—refer to CUDA code section.

On top of the quick times, both the column per thread and row per thread strategies save system resources as less threads have to be allocated for a task when compared to the element per thread algorithm.

As for the first launch configuration, which sought to increase use of system resources to maximize the total number of tasks in parallel across all threads—through having a higher thread per block count, it is unexpected that the times from this strategy would (1) be greater than the other two methods and (2) be so substantially greater than the other two. This likely shows that the work was dispersed to too great an extent, to the point where the cost of creating a new block and thread was greater than assigning an extra task to each existing thread.