# Machine Problem 3 Report

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

Nathan Goodman:  20228249
Date: April 5th, 2023

# CUDA Code

For the three different requested outputs, three different CUDA scripts were created:

1. Find the time it takes to transfer two matrices worth of data to and from device:

Note: Block width is equal to 16 and gpu execution time includes transfer times.

```
1    #include "cuda_runtime.h"
2    #include <string.h>
3    #include <stdio.h>
4    #include <stdlib.h>
5    #include <time.h>
6    #include <math.h>
7
8    #define debug 0 //If debug = 1, certain print statements will be enabled
9
10   /*
11   1. Code __global__ (GPU) function for matrixMultiplication
12   */
13   __global__ void matrixMultiplication_kernel(float*d_a, float*d_b, float*d_c, int sizeOfMatricies) {
14       int Row = blockIdx.y*blockDim.y + threadIdx.y;
15       int Col = blockIdx.x*blockDim.x + threadIdx.x;
16       if (Row < sizeOfMatricies && Col < sizeOfMatricies) {
17           float Pvalue = 0;
18           for (int k = 0; k < sizeOfMatricies; k++) {
19               Pvalue += d_a[Row*sizeOfMatricies + k] * d_b[k*sizeOfMatricies + Col];
20           }
21           d_c[Row*sizeOfMatricies + Col] = Pvalue;
22       }
23
24   }
25
26   /*
27   1b. Create regular C function (to run on CPU) for matrixMultiplication
28   */
29   void matrixMultiplication(float *a, float *b, float *c, int sizeOfMatricies) {
30       for (int i = 0; i < sizeOfMatricies; i++) {        //i =row
31           for (int j = 0; j < sizeOfMatricies; j++) {    //j =col
32               //For each value of matrix c (c[i][j]):
33               float sum = 0;
34               for (int k = 0; k < sizeOfMatricies; k++) {
35                   sum += *(a + i*sizeOfMatricies + k) * *(b + k*sizeOfMatricies + j);
36                   //sum = sum  + (a[Row][i] * b[i][Col]
37               }
38               *(c + i*sizeOfMatricies + j) = sum;
39           }
40       }
41   }
42
43   //Helpful Functions from MP2:
44   void fillMatrix(float *a, int size) {
45       for (int i = 0; i < size; i++) {
46           for (int j = 0; j < size; j++) {
47               *(a + i*size + j) = rand() % 100; //Every element will be a random number in range of 0 to 100
48           }
49       }
50   }
51   int correct_output(float *a, float *b, int sizeOfMatricies) {
52       //Checks if two matrices are equal returns 1 if yes 0 if no.
53       for (int i = 0; i < sizeOfMatricies; i++) {
54           for (int j = 0; j < sizeOfMatricies; j++) {
55               if (*(a + i*sizeOfMatricies + j) != *(b + i*sizeOfMatricies + j)) {
56                   //If a[i][j] != b[i][j]:
57                   return 0;
58               }
59           }
60       }
61       return 1;
62   }
63   void printMatrix(float *a, int size) {
64       for (int i = 0; i < size; i++) {
65           for (int j = 0; j < size; j++) {
66               printf("%f ", *(a + i*size + j));
67           }
68           printf("\n");
69       }
70   }
```

```c
72    int main(int argc, char *argv[]) {
73        cudaDeviceProp deviceProps;
74
75        //Get Device Name (They did this in tutorial so I'll do it here)
76        cudaGetDeviceProperties(&deviceProps, 0);
77        printf("CUDA device [%s]\n", deviceProps.name);
78        if (debug == 1) {
79            printf("\tNumber of Mulitprocessors: %d\n", deviceProps.multiProcessorCount);
80            printf("\tMax Threads Per Block: %d\n", deviceProps.maxThreadsPerBlock);
81            printf("\tMax Dimension of a Block: %d\n", deviceProps.maxThreadsDim);
82            printf("\tMax Dimension of a Grid: %d\n", deviceProps.maxGridSize);
83        }
84
85        //Set Size of Matricies:
86        const int dimOfMatricies = 2000; //Value determines size of matricies Ex. dimOfMatricies = 5 will result in 5x5 matricies
87        printf("Dimensions of Matricies: %dx%d\n", dimOfMatricies, dimOfMatricies);
88
89        //Calculate amount of memory they take:
90        int nbytes = dimOfMatricies*dimOfMatricies*sizeof(float);
91
92        //Allocate host memory for matricies:
93        float *a = 0;
94        float *b = 0;
95        float *c = 0;
96        float *cCopy = 0;                    //Pointer to a copy of second output matrix, to be used to observe GPU to host tranfer time
97        cudaMallocHost((void**)&a, nbytes); //Allocates host memory for matrix A, and points pointer a to first value.
98        cudaMallocHost((void**)&b, nbytes);
99        cudaMallocHost((void**)&c, nbytes);
100       cudaMallocHost((void**)&cCopy, nbytes);
101
102       if (debug == 1) {
103           printf("Memory Locations of matricies:\n");
104           printf("\ta = %x\n",&a);
105           printf("\tb = %x\n", &b);
106           printf("\tc = %x\n", &c);
107       }
108
109       //Generate input matricies into memory
110
111       srand(time(NULL));
112       fillMatrix(a, dimOfMatricies);
113       fillMatrix(b, dimOfMatricies);
114
115       if (debug == 1) {
116           printf("Matrix A:\n");
117           printMatrix(a, dimOfMatricies);
118           printf("Matrix B:\n");
119           printMatrix(b, dimOfMatricies);
120       }
121
122       //1e. Allocate device memory for matricies:
123       float *d_a = 0;
124       float *d_b = 0;
125       float *d_c = 0;
126       cudaMalloc((void**)&d_a, nbytes); //Allocates memory for matrix A, and points pointer a to first value.
127       cudaMemset(d_a, 255, nbytes);     //Sets all allocated bytes to 255 (they did this in tutorial so i did it here)
128       cudaMalloc((void**)&d_b, nbytes);
129       cudaMemset(d_b, 255, nbytes);
130       cudaMalloc((void**)&d_c, nbytes);
131       cudaMemset(d_c, 255, nbytes);
132
133       //Set kernel launch configuration
134       int blckWidth = 16;
135       int threadsPerBlock = blckWidth*blckWidth;
136       int threadsNeeded = ceil(dimOfMatricies*dimOfMatricies); //Because in this configuration one thread only produces one value
137       int blocksNeeded = ceil(threadsNeeded / threadsPerBlock);
138       if (blocksNeeded < 1) blocksNeeded++;
139       int gridWidth = ceil(sqrt(blocksNeeded)); //Note grid will be square
```

```
141        if (debug == 1) {
142            printf("Block Width: %d\n", blckWidth);
143            printf("Grid Width: %d\n", gridWidth);
144        }
145
146        dim3 dimBlocks = dim3(blckWidth, blckWidth);
147        dim3 dimGrid = dim3(gridWidth, gridWidth);
148
149        //create cuda event handles
150        cudaEvent_t start, stopToDevice, startToHost, stop;
151        cudaEventCreate(&start);
152        cudaEventCreate(&stopToDevice);
153        cudaEventCreate(&startToHost);
154        cudaEventCreate(&stop);
155
156        cudaDeviceSynchronize();
157        float gpu_time = 0.0f;
158        float toGPU_time = 0.0f;
159        float fromGPU_time = 0.0f;
160
161        //asynchronously issue work to the GPU (all stream 0)
162        cudaEventRecord(start, 0);
163
164        //Copy inputs to device
165        cudaMemcpyAsync(d_a, a, nbytes, cudaMemcpyHostToDevice, 0);
166        cudaMemcpyAsync(d_b, b, nbytes, cudaMemcpyHostToDevice, 0);
167
168        //Save time it took to write to device
169        cudaEventRecord(stopToDevice, 0);
170
171        //Call Kernel
172        matrixMultiplication_kernel << <dimGrid, dimBlocks, 0, 0 >> >(d_a, d_b, d_c, dimOfMatricies);
173        //Note: only section after >>> is the actual function parameters
174
175        //Save time before saving to host
176        cudaEventRecord(startToHost, 0);
177
178        //Copy outputs from device
179        cudaMemcpyAsync(c, d_c, nbytes, cudaMemcpyDeviceToHost, 0);
180        cudaMemcpyAsync(cCopy, d_c, nbytes, cudaMemcpyDeviceToHost, 0); //A second matrix is recorded as question asked for tranfer to two matricies
181
182        cudaEventRecord(stop, 0);
183        cudaEventSynchronize(stopToDevice);
184        cudaEventSynchronize(startToHost);
185        cudaEventSynchronize(stop);
186        cudaEventElapsedTime(&gpu_time, start, stop);          //Note: gpu_time includes toGPU_time and fromGPU_time
187        cudaEventElapsedTime(&toGPU_time, start, stopToDevice);
188        cudaEventElapsedTime(&fromGPU_time, startToHost, stop);
189
190        //print the GPU times
191        printf("time spent executing by the GPU: %.2f\n", gpu_time);
192        printf("time spent transferring input data to the GPU: %.2f\n", toGPU_time);
193        printf("time spent transferring output data to host: %.2f\n", fromGPU_time);
194
195        //Calculate matrixMultiplication using CPU:
196        float *d = 0;
197        cudaMallocHost((void**)&d, nbytes);
198        matrixMultiplication(a, b, d, dimOfMatricies); //d is output matrix
199
200                                        //Check Output
201        bool bFinalResults = (bool)correct_output(c, d, dimOfMatricies); //check if c(from GPU) = d(from CPU)
202        if (bFinalResults == true) {
203            printf("Test PASSED\n");
204        }
205
206        if (debug == 1) {
207            printf("Matrix C:\n");
208            printMatrix(c, dimOfMatricies);
209            printf("Correct Matrix:\n");
210            printMatrix(d, dimOfMatricies);
211        }
212
213        //release resources
214        cudaEventDestroy(start);
215        cudaEventDestroy(stop);
216
217        cudaFreeHost(a);
218        cudaFreeHost(b);
219        cudaFreeHost(c);
220        cudaFreeHost(cCopy);
221        cudaFreeHost(d);
222        cudaFree(d_a);
223        cudaFree(d_b);
224        cudaFree(d_c);
225        cudaDeviceReset();
226
227        return 0;
228
229    }
```

2.  Find time it takes to execute matrix on GPU with one block and one thread vs CPU

Note: block width is equal to 16. gpu execution time does NOT include transfer times.

```
1    #include "cuda_runtime.h"
2    #include <string.h>
3    #include <stdio.h>
4    #include <stdlib.h>
5    #include <time.h>
6    #include <math.h>
7
8    #define debug 0 //If debug = 1, certain print statements will be enabled
9
10   /*
11       2. Code __global__ (GPU) function for matrixMultiplication but one thread calculates whole output matrix
12   */
13   __global__ void matrixMultiplication_kernel(float*d_a, float*d_b, float*d_c, int sizeOfMatricies) {
14       for (int i = 0; i < sizeOfMatricies; i++) {          //i =row
15           for (int j = 0; j < sizeOfMatricies; j++) {      //j =col
16               float Pvalue = 0;
17               for (int k = 0; k < sizeOfMatricies; k++) {
18                   Pvalue += d_a[i*sizeOfMatricies + k] * d_b[k*sizeOfMatricies + j];
19               }
20               d_c[i*sizeOfMatricies + j] = Pvalue;
21           }
22       }
23
24   }
25
26   //Regular C function (to run on CPU) for matrixMultiplication:
27   void matrixMultiplication(float *a, float *b, float *c, int sizeOfMatricies) {
28       for (int i = 0; i < sizeOfMatricies; i++) {          //i =row
29           for (int j = 0; j < sizeOfMatricies; j++) {      //j =col
30               //For each value of matrix c (c[i][j]):
31               float sum = 0;
32               for (int k = 0; k < sizeOfMatricies; k++) {
33                   sum += *(a + i*sizeOfMatricies + k) * *(b + k*sizeOfMatricies + j);
34                   //sum = sum  + (a[Row][i] * b[i][Col]
35               }
36               *(c + i*sizeOfMatricies + j) = sum;
37           }
38       }
39   }
40
41   //Helpful Functions from MP2:
```

```
131      //Set kernel launch configuration
132      int blckWidth = 1;
133      int gridWidth = 1;
134      /* Below is not needed for this part as it asks for one block and one thread
135      int threadsPerBlock = blckWidth*blckWidth;
136      int threadsNeeded = ceil(dimOfMatricies*dimOfMatricies); //Because in this configuration one thread only produces one value
137      int blocksNeeded = ceil(threadsNeeded / threadsPerBlock);
138      if (blocksNeeded < 1) blocksNeeded++;
139      int gridWidth = ceil(sqrt(blocksNeeded)); //Note grid will be square
140      */
```

```
191      //print the GPU times
192      printf("time spent executing by the GPU: %.2f\n", (gpu_time - toGPU_time - fromGPU_time)); //Time spent executing by the GPU not includeing transfers
193      //printf("time spent transferring input data to the GPU: %.2f\n", toGPU_time);
194      //printf("time spent transferring output data to host: %.2f\n", fromGPU_time);
```

Note: All parts not shown were identical to the parts from the first configuration. Full script can be viewed from attached .zip

3.  Increase block_width as matrix size increases.

Note: Gpu execution time does NOT include transfer times.

```
138      //Set kernel launch configuration
139      int blckWidth = 25;
140      int threadsPerBlock = blckWidth*blckWidth;
141      int threadsNeeded = ceil(dimOfMatricies*dimOfMatricies); //Because in this configuration one thread only produces one value
142      int blocksNeeded = ceil(threadsNeeded / threadsPerBlock);
143      if (blocksNeeded < 1) blocksNeeded++;
144      int gridWidth = ceil(sqrt(blocksNeeded)); //Note: grid will be square
145
146      printf("Block Width: %d\n", blckWidth);
147      printf("Grid Width: %d\n", gridWidth);
148
```

```
191        //print the GPU times
192        printf("time spent executing by the GPU: %.2f\n", (gpu_time - toGPU_time - fromGPU_time)); //Time spent executing by the GPU not includeing transfers
193        //printf("time spent transferring input data to the GPU: %.2f\n", toGPU_time);
194        //printf("time spent transferring output data to host: %.2f\n", fromGPU_time);
```

Note: All parts not shown were identical to the parts from the first configuration. Full scripts can be viewed from attached .zip

## Output

For each question, outputs were recorded for the following matrix sizes: 125x125, 250x250, 500x500, 1000x1000, 2000x2000; with each configuration from above.

1.  Find the time it takes to transfer two matrices worth of data to and from device:

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 125x125
time spent executing by the GPU: 0.28
time spent transferring input data to the GPU: 0.08
time spent transferring output data to host: 0.03
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 250x250
time spent executing by the GPU: 1.20
time spent transferring input data to the GPU: 0.13
time spent transferring output data to host: 0.08
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 500x500
time spent executing by the GPU: 7.75
time spent transferring input data to the GPU: 0.38
time spent transferring output data to host: 0.32
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 1000x1000
time spent executing by the GPU: 58.40
time spent transferring input data to the GPU: 1.35
time spent transferring output data to host: 1.24
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 2000x2000
time spent executing by the GPU: 444.63
time spent transferring input data to the GPU: 5.27
time spent transferring output data to host: 4.92
Test PASSED
Press any key to continue . . . _
```

2.  Find time it takes to execute matrix on GPU with one block and one thread vs CPU

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 125x125
time spent executing by the GPU: 811.50
time spent executing by CPU: 10.00
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 250x250
time spent executing by the GPU: 9092.94
time spent executing by CPU: 106.00
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 500x500
time spent executing by the GPU: 73857.30
time spent executing by CPU: 1001.00
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 1000x1000
time spent executing by the GPU: 593185.50
time spent executing by CPU: 12923.00
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 2000x2000
time spent executing by the GPU: 4745166.00
time spent executing by CPU: 119458.00
Test PASSED
Press any key to continue . . . _
```

3. Increase block_width as matrix size increases.

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 125x125
Block Width: 2
Grid Width: 63
time spent executing by the GPU: 2.86
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 250x250
Block Width: 4
Grid Width: 63
time spent executing by the GPU: 6.13
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 500x500
Block Width: 10
Grid Width: 50
time spent executing by the GPU: 9.00
Test PASSED
Press any key to continue . . . _
```

```
CUDA device [Tesla C2075]
Dimensions of Matricies: 1000x1000
Block Width: 20
Grid Width: 50
time spent executing by the GPU: 62.11
Test PASSED
Press any key to continue . . . _
```
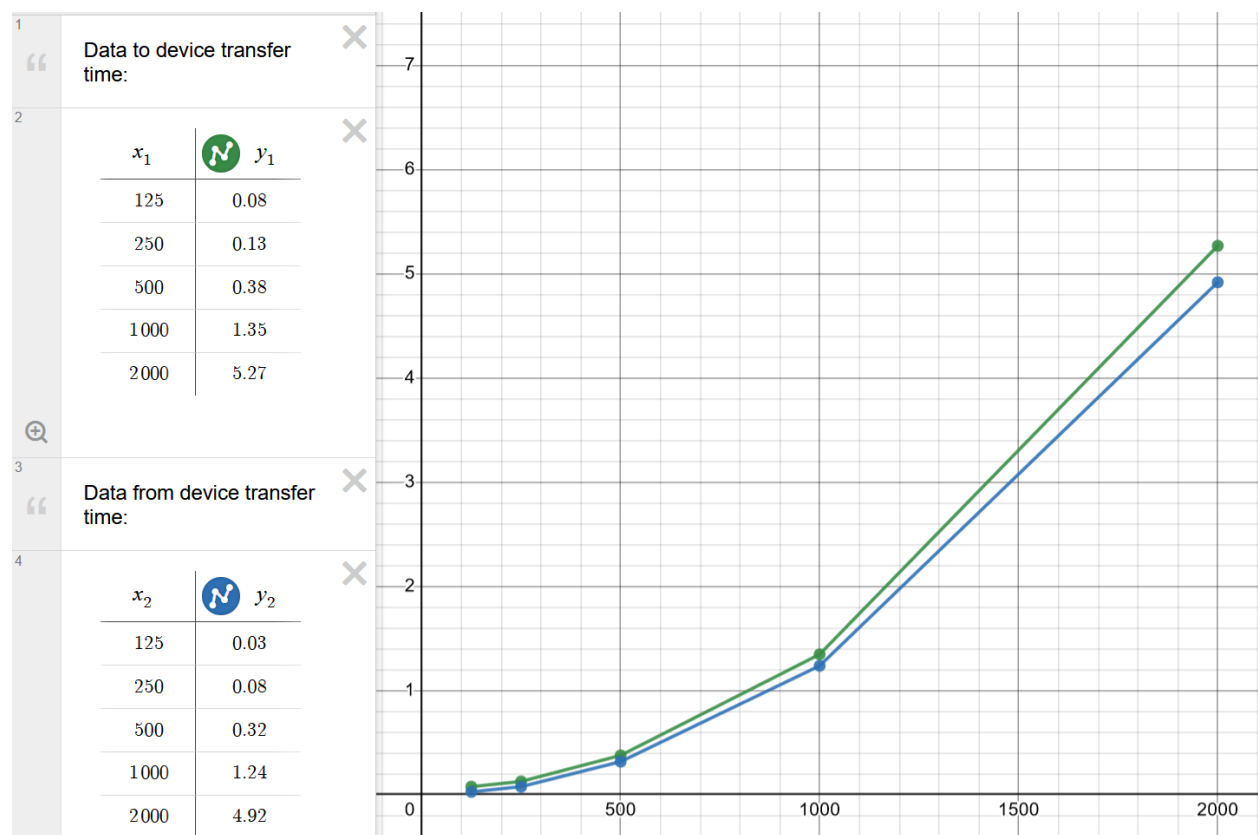
```
CUDA device [Tesla C2075]
Dimensions of Matricies: 2000x2000
Block Width: 25
Grid Width: 80
time spent executing by the GPU: 516.63
Test PASSED
Press any key to continue . . . _
```

## Analysis

1. For the first configuration, the following graph showing transfer times to and from device versus matrix size, were generated from the outputs above. The x-axis is block width and the y-axis is transfer time (in milliseconds):

Data to device transfer time:

| $x_1$ | $y_1$ |
|-------|-------|
| 125   | 0.08  |
| 250   | 0.13  |
| 500   | 0.38  |
| 1000  | 1.35  |
| 2000  | 5.27  |

Data from device transfer time:

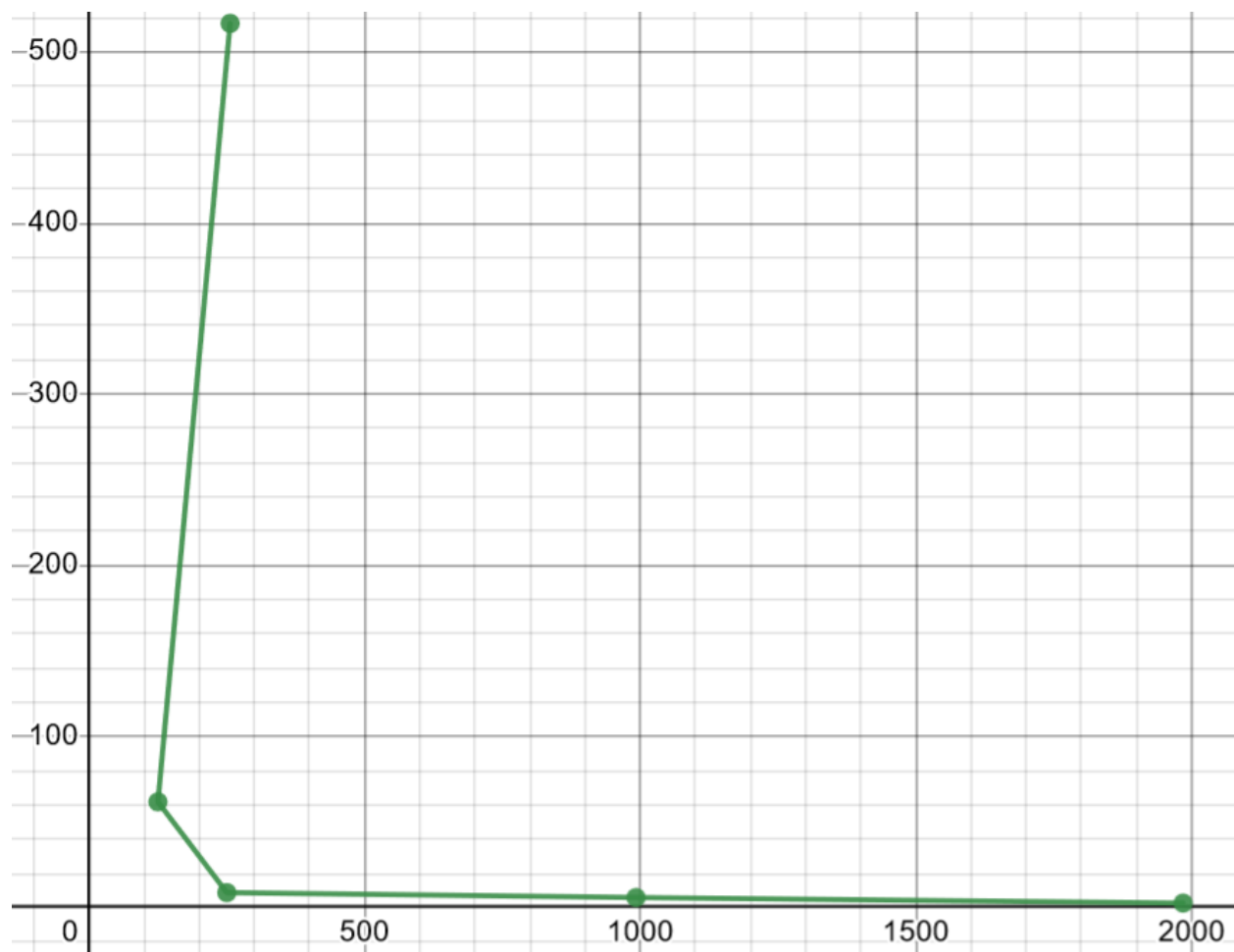| $x_2$ | $y_2$ |
|-------|-------|
| 125   | 0.03  |
| 250   | 0.08  |
| 500   | 0.32  |
| 1000  | 1.24  |
| 2000  | 4.92  |



There is for the most part no difference in transfer speeds to and from the gpu device, with only a slight—seemingly proportional difference being seen (difference < 0.4 ms for 2000x2000). This small difference is likely due to a difference in memory speeds between the device and host.

2. Observing the outputs from the second configuration (above in Outputs), it becomes obvious that setting the grid width and block width to one had a detrimental impact on the GPU's performance. When comparing this performance to the performance of the matrix multiplication on the CPU, its clear that it is not always beneficial to offload your matrix multiplication to the GPU; it is only optimal to do so if the kernel configuration is optimized correctly.

3. The outputs for the third configuration were generated from the following matrixWidth and blockWidth pairs: (125, 2), (250, 4), (500, 10), (1000,20), (2000,25).

The following graph of kernel execution time vs (numBlocks/blockWidth) were generated from the outputs above (Note: numBlocks = gridWidth$^2$). The x-axis is numBlocks/blockWidth and the y-axis is execution time (in milliseconds):



This clearly shows that when there are more blocks with less threads per block, the execution time will be faster than when there is less blocks with more threads per block.

a) Using our knowledge of the matrix multiplication algorithm, we know that to calculate an element in output matrix C, we must access every element from input matrix A sharing a row with the output element, and we must access every element from input matrix B sharing a column with the output element. This means for every element in matrix C we are accessing 2*matrixWidth number of input elements:

number of total accesses = amount of output matrix elements * 2*matrixWidth

Since C is a square matrix 'matrixWidth' wide, the number of output elements will be: matrixWidth$^2$

If we are trying to find the amount of times each input element is accessed:

$$\text{average number of times an input element is accessed} = \frac{number\ of\ total\ accesses}{number\ of\ input\ elements}$$

The number of input elements is equal to all the elements in A plus all the elements from B. Since both of these square matrices are a 'matrixWidth' wide:

Number of input elements = 2(matrixWidth$^2$)

Combining all of these into one equation:

$$avgInputElementAccesses = \frac{totalAccesses}{numberOfInputElements}$$

$$avgInputElementAccesses = \frac{matrixWidth^2 * (2 * matrixWidth)}{2 * matrixWidth^2}$$

$$avgInputElementAccesses = matrixWidth$$

Therefore, for this configuration where the output matrixWidth is equal to the input matrices matrixWidth, the number of times a given input matrix element will be accessed is equal to the matrix width. This makes sense as a given element in A or B needs to be accessed once for each element in C, in the same row/column (=matrixWidth).

b) CGMA ratio is the number of operations a thread performs over the number of memory fetches. Looking at our matrix multiplication kernel for the third configuration (identical to the one for the first configuration), for each loop iteration within the kernel:

```
for (int k = 0; k < sizeOfMatricies; k++) {
    Pvalue += d_a[Row*sizeOfMatricies + k] * d_b[k*sizeOfMatricies + Col];
}
```

There's two global memory accesses d_a[…] and d_b[…] for one floating-point multiplication and one floating-point addition operation. Giving a CGMA ratio of floating-point operations to bytes accessed from global memory of 0.25 FLOP/B.