

## Machine Problem 1 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.

```
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  //2. Write a kernel that has each thread producing one output matrix. Kernel config should be 16x16 thread blocks
9
10 /* MatrixAddition Kernel.
11 Parameters: pointer to output matrix C, two Pointers to input matrices A and B,
12 dimensions of matrices A and B (remember they're square matrices so this can be single int)
13 */
14 __global__ void matrixAddition_kernel(float*d_a, float*d_b, float*d_c, int sizeofMatrices) {
15     int Row = blockIdx.y*blockDim.y + threadIdx.y;
16     int Col = blockIdx.x*blockDim.x + threadIdx.x;
17     int idx = Row* sizeofMatrices + Col;
18     //0 1 2 3 4
19     //5 6 7 8 9
20     if (idx < sizeofMatrices* sizeofMatrices) { //Avoid accessing beyond end of matrices
21         d_c[idx] = d_a[idx] + d_b[idx];
22     }
23 }
24
25
26 /* Created a matrixAddition function. Should basically be the same as the kernel function.
27 This function will be used to check whether the kernel function created correct output (used to check)
28 */
29 void matrixAddition(float *a, float *b, float *c, int sizeofMatrices) { //Note acc
30     for (int i = 0; i < sizeofMatrices; i++) {
31         for (int j = 0; j < sizeofMatrices; j++) {
32             *(c + i* sizeofMatrices + j) = *(a + i* sizeofMatrices + j) + *(b + i* sizeofMatrices + j);
33             //C[i][j] = A[i][j] + B[i][j];
34         }
35     }
36 }
37
38 //3b. Create function to check if two matrices are equal (to be used to compare outputs)
39 int correct_output(float *a, float *b, int sizeofMatrices) {
40     for (int i = 0; i < sizeofMatrices; i++) {
41         for (int j = 0; j < sizeofMatrices; j++) {
42             if (*(a + i* sizeofMatrices + j) != *(b + i* sizeofMatrices + j)) {
43                 //If a[i][j] != b[i][j]:
44                 return 0;
45             }
46         }
47     }
48     return 1;
49 }
50
51 //Additional function to help with debugging:
52 void printMatrix(float *a, int size) {
53     for (int i = 0; i < size; i++) {
54         for (int j = 0; j < size; j++) {
55             printf("%f ", *(a + i* size + j));
56         }
57         printf("\n");
58     }
59 }
60
61 //Function for generating randomly initialized square matrices of a given length
62 void fillMatrix(float *a, int size) {
63     for (int i = 0; i < size; i++) {
64         for (int j = 0; j < size; j++) {
65             *(a + i* size + j) = rand() % 100; //Every element will be a random number in range of 0 to 100
66         }
67     }
68 }
69 }
```

```

70 int main(int argc, char *argv[]) {
71     cudaDeviceProp deviceProps;
72
73     //Get Device Name (They did this in tutorial so I'll do it here)
74     cudaGetDeviceProperties(&deviceProps, 0);
75     printf("CUDA device [%s]\n", deviceProps.name);
76     printf("\tNumber of Multiprocessors: %d\n", deviceProps.multiProcessorCount);
77     printf("\tMax Threads Per Block: %d\n", deviceProps.maxThreadsPerBlock);
78     printf("\tMax Dimension of a Block: %d\n", deviceProps.maxThreadsDim);
79     printf("\tMax Dimension of a Grid: %d\n", deviceProps.maxGridSize);
80
81     /*
82     1. Define two square input matrices A and B, and matching output matrix
83     Note: they're floats
84     */
85     const int dimOfMatrices = 250; //Value determines size of matrices Ex. dimOfMatrices = 5 will result in 5x5 matrices
86
87     float C[dimOfMatrices][dimOfMatrices];
88
89     //1b. Calculate amount of memory they take:
90     int nbytes = dimOfMatrices*dimOfMatrices*sizeof(float);
91
92     //1c. Allocate host memory for matrices:
93     float *a = 0;
94     float *b = 0;
95     float *c = 0;
96     cudaMallocHost((void**)&a, nbytes); //Allocates host memory for matrix A, and points pointer a to first value.
97     cudaMallocHost((void**)&b, nbytes);
98     cudaMallocHost((void**)&c, nbytes);
99
100     //printf("a = %x\n", &a);
101     //printf("b = %x\n", &b);
102     //printf("c = %x\n", &c);
103
104     //1d. Store input matrices into memory
105
106     srand(time(NULL));
107     fillMatrix(a, dimOfMatrices);
108     fillMatrix(b, dimOfMatrices);
109
110     //printf("%f\n", *b);
111     //printf("Matrix A\n");
112     //printMatrix(a, dimOfMatrices);
113     //printf("Matrix B\n");
114     //printMatrix(b, dimOfMatrices);
115
116     //1e. Allocate device memory for matrices:
117     float *d_a = 0;
118     float *d_b = 0;
119     float *d_c = 0;
120     cudaMalloc((void**)&d_a, nbytes); //Allocates memory for matrix A, and points pointer a to first value.
121     cudaMemcpy(d_a, 255, nbytes); //Sets all allocated bytes to 255 (they did this in tutorial so i did it here)
122     cudaMalloc((void**)&d_b, nbytes);
123     cudaMemcpy(d_b, 255, nbytes);
124     cudaMalloc((void**)&d_c, nbytes);
125     cudaMemcpy(d_c, 255, nbytes);
126
127     //Set kernel launch configuration
128     int blkWidth = 16; //block width and kength
129     int threadsPerBlock = blkWidth*blkWidth;
130     int threadsNeeded = dimOfMatrices*dimOfMatrices; //Because in this configuration one thread only produces one value
131     int numBlocks = threadsNeeded / threadsPerBlock;
132     if (numBlocks < 1) numBlocks++;
133
134     dim3 dimBlocks = dim3(blkWidth, blkWidth); //asks for 16x16=256 thread blocks
135     //256 threads per block layed out in 16x16
136     dim3 dimGrid = dim3(numBlocks, numBlocks);

```

```

138     cudaEvent_t start, stop;
139     cudaEventCreate(&start);
140     cudaEventCreate(&stop);
141
142     cudaDeviceSynchronize();
143     float gpu_time = 0.0f;
144
145     //asynchronously issue work to the GPU (all stream 0)
146     cudaEventRecord(start, 0);
147
148     //Copy inputs to device
149     cudaMemcpyAsync(d_a, a, nbytes, cudaMemcpyHostToDevice, 0);
150     cudaMemcpyAsync(d_b, b, nbytes, cudaMemcpyHostToDevice, 0);
151
152     //Call Kernel
153     matrixAddition_kernel << <dimGrid, dimBlocks, 0, 0 >> >(d_a, d_b, d_c, dimOfMatrices);
154     //Note: only section after >>> is the actual function parameters
155
156     //Copy outputs from device
157     cudaMemcpyAsync(c, d_c, nbytes, cudaMemcpyDeviceToHost, 0);
158
159     //2b. Load output matrices from memory
160     for (int i = 0; i < dimOfMatrices; i++) {
161         for (int j = 0; j < dimOfMatrices; j++) {
162             C[i][j] = *(c + i*dimOfMatrices + j);
163         }
164     }
165
166     cudaEventRecord(stop, 0);
167     cudaEventSynchronize(stop); //stop is updated here
168     cudaEventElapsedTime(&gpu_time, start, stop);
169
170     //print the GPU times
171     printf("time spent executing by the GPU: %.2f\n", gpu_time);
172
173     //Calculate matrixAddition using CPU:
174     float D[dimOfMatrices][dimOfMatrices];
175     float *d = 0;
176     cudaMallocHost((void**)&d, nbytes);
177     matrixAddition(a, b, d, dimOfMatrices); //d is output matrix
178
179     //Check Output
180     bool bFinalResults = (bool)correct_output(c, d, dimOfMatrices); //check if c(from GPU) = d(from CPU)
181     if (bFinalResults == true) {
182         printf("Test PASSED\n");
183     }
184
185     //printf("Matrix C\n");
186     //printMatrix(c, dimOfMatrices);
187     //printf("\n");
188     //printMatrix(d, dimOfMatrices);
189
190     //release resources
191     cudaEventDestroy(start);
192     cudaEventDestroy(stop);
193
194     cudaFreeHost(a);
195     cudaFreeHost(b);
196     cudaFreeHost(c);
197     cudaFreeHost(d);
198     cudaFree(d_a);
199     cudaFree(d_b);
200     cudaFree(d_c);
201     cudaDeviceReset();
202
203     return 0;
204
205 }

```

- 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 matrices A and B,
dimensions of matrices A and B (remember they're square matrices so this can be single int)
*/
__global__ void matrixAddition_kernel(float*d_a, float*d_b, float*d_c, int sizeOfMatrices) {
    int idx = blockIdx.x*blockDim.x + threadIdx.x; //Assume block

    idx = idx*sizeOfMatrices; //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 = sizeOfMatrices*sizeOfMatrices;
    if (idx < n) { //Avoid accessing beyond end of matrices
        for (int i = 0; i < sizeOfMatrices; 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 blkWidth = 4; //block width and length (16 threads per block)
    int threadsPerBlock = blkWidth*blkWidth;
    int threadsNeeded = dimOfMatrices*dimOfMatrices;
    int numBlocks = threadsNeeded / threadsPerBlock;
    if (numBlocks < 1) numBlocks++;

    dim3 dimBlocks = dim3(16, 1); // 16 threads per block
    dim3 dimGrid = dim3(numBlocks, 1);
}
```

- 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 matrices A and B,
dimensions of matrices A and B (remember they're square matrices so this can be single int)
*/
__global__ void matrixAddition_kernel(float*d_a, float*d_b, float*d_c, int sizeOfMatrices) {
    int idx = blockIdx.x*blockDim.x + threadIdx.x; //Assume block
    int n = sizeOfMatrices*sizeOfMatrices;
    if (idx < n) { //Avoid accessing beyond end of matrices
        for (int i = 0; i < sizeOfMatrices; i++) {
            //Each thread will calculate a column of output matrix:
            d_c[idx + sizeOfMatrices*i] = d_a[idx + sizeOfMatrices*i] + d_b[idx + sizeOfMatrices*i];
        }
    }

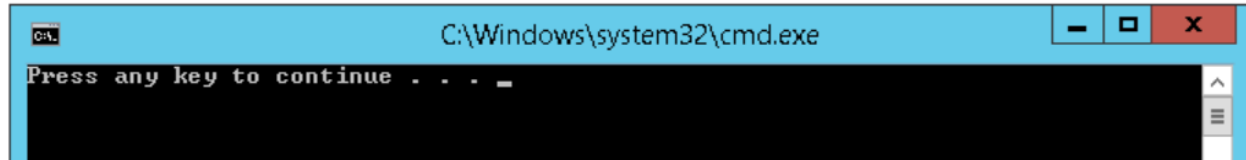
    //Set kernel launch configuration
    int blkWidth = 4; //block width and length (16 threads per block)
    int threadsPerBlock = blkWidth*blkWidth;
    int threadsNeeded = dimOfMatrices*dimOfMatrices;
    int numBlocks = threadsNeeded / threadsPerBlock;
    if (numBlocks < 1) numBlocks++;

    dim3 dimBlocks = dim3(16, 1); // 16 threads per block
    dim3 dimGrid = dim3(numBlocks, 1);
}
```

## Output

Following guidelines from the Machine Problem 2 Document, the `dimOfMatrices` 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, the script failed to execute despite no compilation errors:

```
const int dimOfMatrices = 500; //Value determines size of matrices Ex. dimOfMatrices = 5 will
```



While I could see this being a possibility for the first set of launch configurations (each thread only produces a single output) as to calculate 500\*500 output elements you'd need 500\*500 threads. This issue continued to persist despite improved resource management in second and third launch configurations.

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 each kernel configuration can be viewed below. 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: 1899400
  Max Dimension of a Grid: 1899412
time spent executing by the GPU: 0.18
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: 14481928
  Max Dimension of a Grid: 14481940
time spent executing by the GPU: 41.36
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: 12777900
  Max Dimension of a Grid: 12777912
time spent executing by the GPU: 170.85
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: 4847812
    Max Dimension of a Grid: 4847824
time spent executing by the GPU: 1.48
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: 17823960
    Max Dimension of a Grid: 17823972
time spent executing by the GPU: 2.91
Test PASSED
Press any key to continue . . . _

```

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

```

CUDA device [Tesla C2075]
    Number of Mulitprocessors: 14
    Max Threads Per Block: 1024
    Max Dimension of a Block: 15531332
    Max Dimension of a Grid: 15531344
time spent executing by the GPU: 0.24
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: 9107908
    Max Dimension of a Grid: 9107920
time spent executing by the GPU: 0.34
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: 6945344
    Max Dimension of a Grid: 6945356
time spent executing by the GPU: 0.51
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: 3601924
    Max Dimension of a Grid: 3601936
time spent executing by the GPU: 0.93
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: 16448768
    Max Dimension of a Grid: 16448780
time spent executing by the GPU: 0.77
Test PASSED
Press any key to continue . . . _

```

- 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 . . . _
```

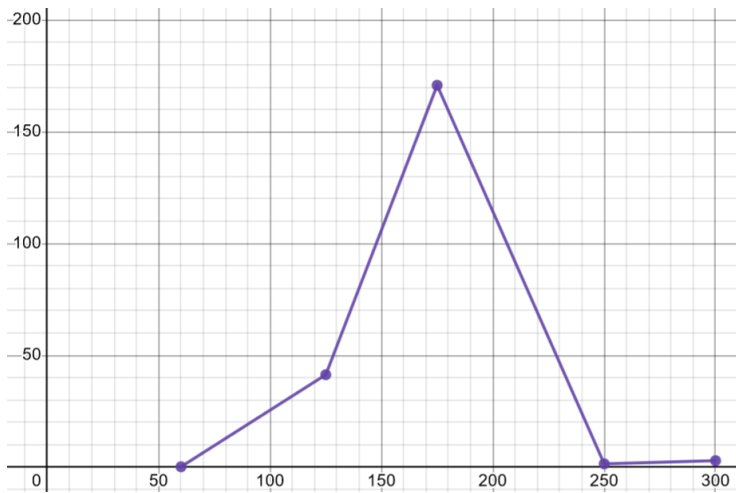
```
CUDA device [Tesla C2075]
  Number of Mulitprocessors: 14
  Max Threads Per Block: 1024
  Max Dimension of a Block: 3930852
  Max Dimension of a Grid: 3930864
time spent executing by the GPU: 13.41
Test PASSED
Press any key to continue . . . _
```

## Analysis

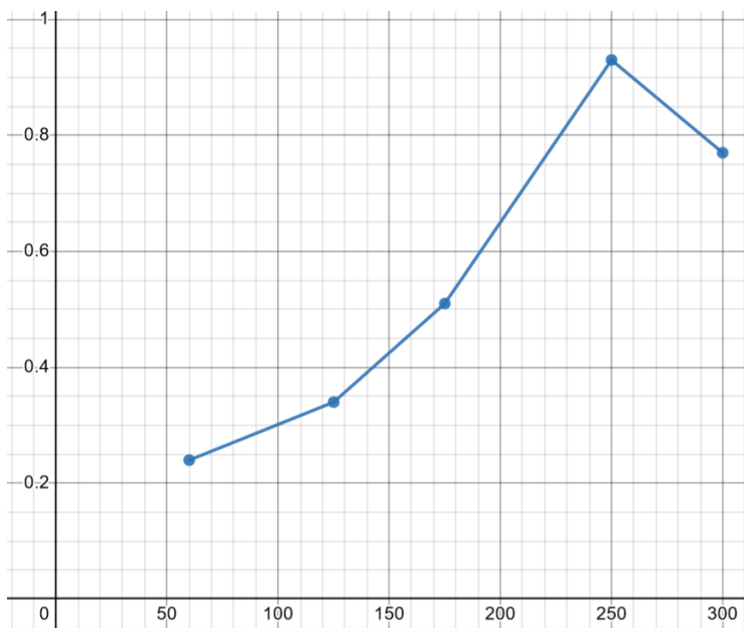
The following graphs were generated for each of the launch configurations with the x-axis being the width of the matrix (in floats 4 bytes) and the y-axis being the GPU execution time (in ms):

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

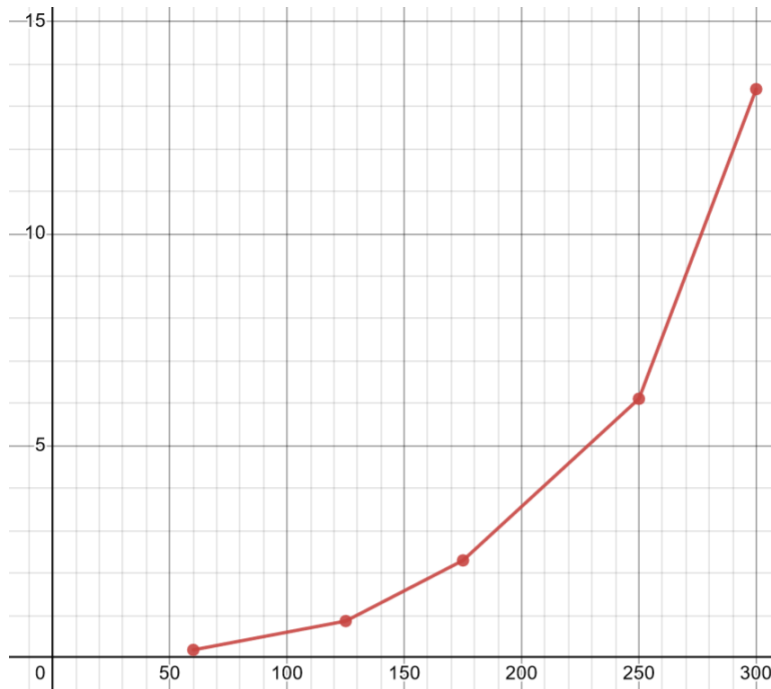




2. Each thread produces a row of output and threads per block is 16:



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



As you can see above, (with the exception of the first launch configuration) 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 kernel function).

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, I believe the unusual behavior and super long execution times were due to the measurements being recorded during a period of high traffic to the virtual machines (around 6pm the day before deadline). As mentioned previously, this methodology seeks to increase use of system resources to maximize the total number of tasks in parallel across all threads. It is unexpected then, 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.