# CSC8014 Topics: GPU Programming
Student: Lina Mi @01377283

## Homework #1

1. Compile and run the code of adding two vectors **a** and **b** in Chapter 4.

   a. What is the result? You can take a screenshot of the result.



```
Solution Explorer        add_vector_GPU_1.cu  X
(Global Scope)                                                          main(void)
    //#include <cstdio>

    #include "../common/book.h"

    #define N   10

    __global__ void add(int *a, int *b, int *c) {
        int tid = blockIdx.x;    // this thread handles the data at its thread id
        if (tid < N)
            c[tid] = a[tid] + b[tid];
    }

    int main(void) {
        int a[N], b[N], c[N];
        int *dev_a, *dev_b, *dev_c;

        // allocate the memory on the GPU
        HANDLE_ERROR(cudaMalloc((void**)&dev_a, N * sizeof(int)));
        HANDLE_ERROR(cudaMalloc((void**)&dev_b, N * sizeof(int)));
        HANDLE_ERROR(cudaMalloc((void**)&dev_c, N * sizeof(int)));

        // fill the arrays 'a' and 'b' on the CPU
        for (int i = 0; i<N; i++) {
            a[i] = -i;
            b[i] = i * i;
        }

        // copy the arrays 'a' and 'b' to the GPU
        cudaMemcpy(dev_a, a, N * sizeof(int),
            cudaMemcpyHostToDevice);
        cudaMemcpy(dev_b, b, N * sizeof(int),
            cudaMemcpyHostToDevice);

        add <<< N, 1 >> >(dev_a, dev_b, dev_c);

        // copy the array 'c' back from the GPU to the CPU
        cudaMemcpy(c, dev_c, N * sizeof(int),
            cudaMemcpyDeviceToHost);

        // display the results
        for (int i = 0; i<N; i++) {
            printf("%d + %d = %d\n", a[i], b[i], c[i]);
        }
    }
```
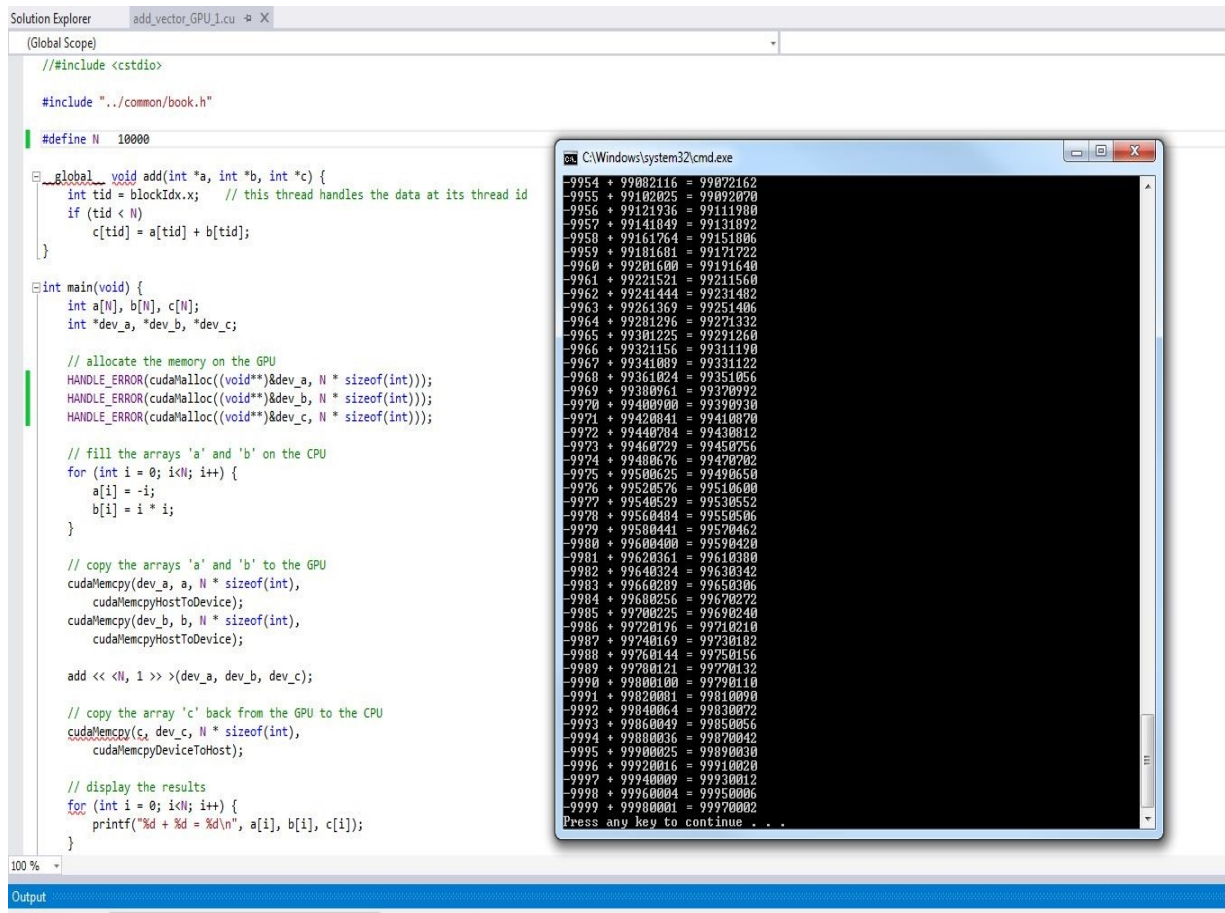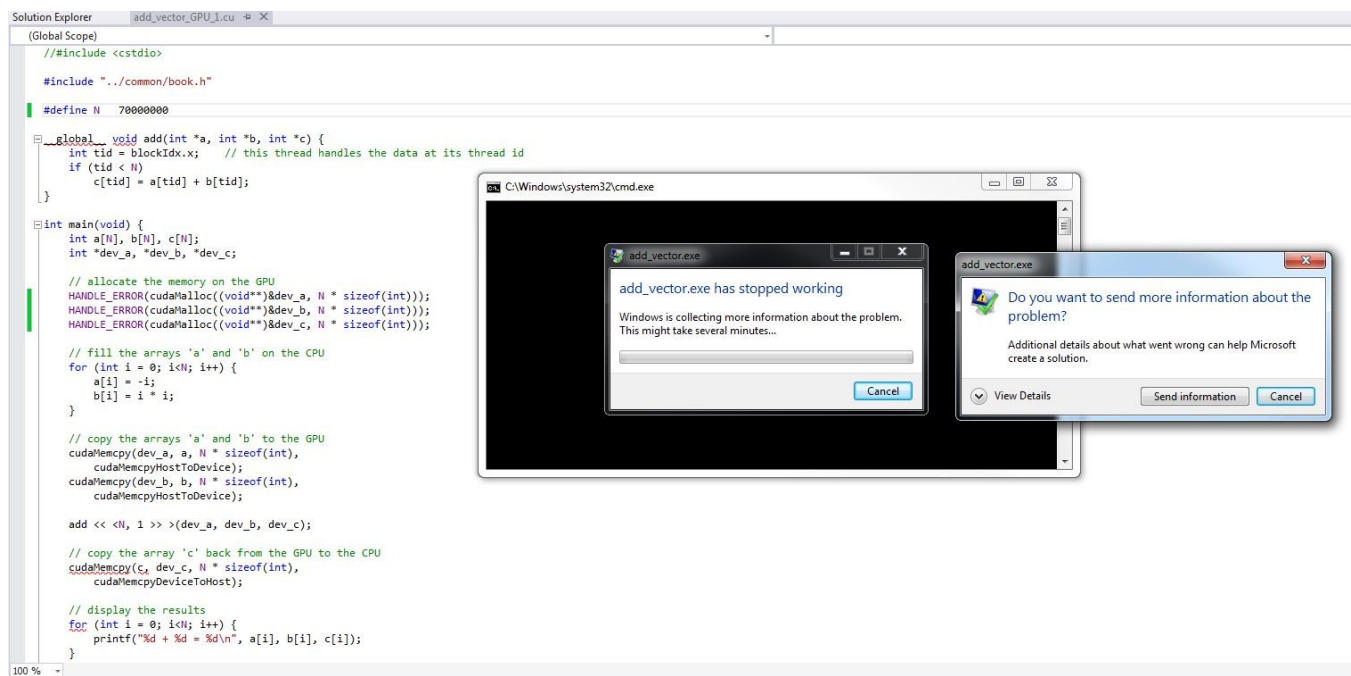
```
C:\Windows\system32\cmd.exe

0 + 0 = 0
-1 + 1 = 0
-2 + 4 = 2
-3 + 9 = 6
-4 + 16 = 12
-5 + 25 = 20
-6 + 36 = 30
-7 + 49 = 42
-8 + 64 = 56
-9 + 81 = 72
Press any key to continue . . .
```

b. Change the number N from 10 to 10000. Compile and run the code again. What is the result?



c. Change the number N to 70000000. Then compile and run the code one more time. Does the code still execute? Explain what happens.
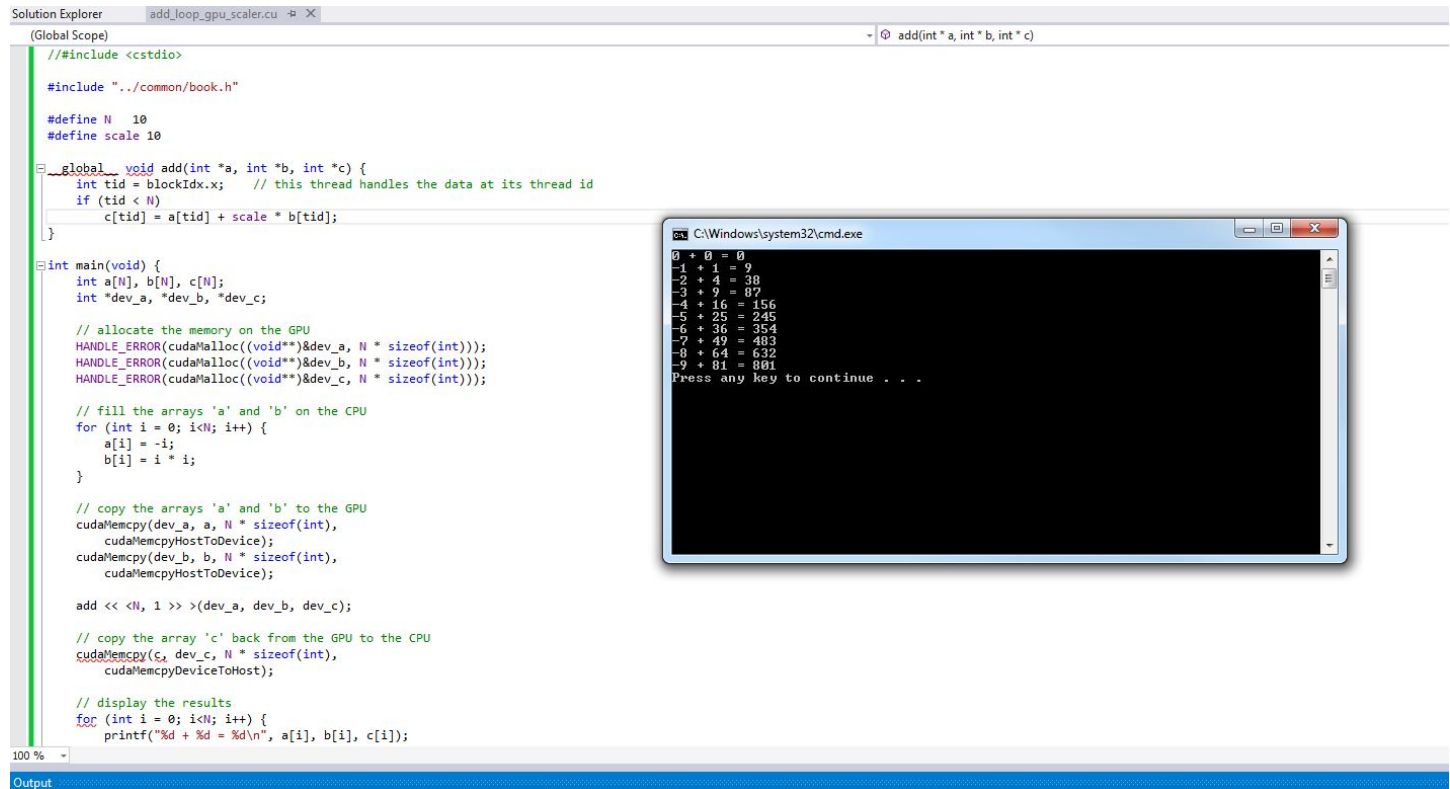
The code stopped executing, because the number of 70000000 is very large and square of it is too large to be hold with a *int* type variable (32bit int type variable can hold maximum value is 2^32=4,294,967,296<70000000^2)

2. Change the code (Also attached as add_loop_gpu.zip) in Chapter 4 to perform the following operation for vectors **a** and **b**:

c[i] = a[i] + scale * b[i]

where scale is an integer.

    a. Compile and run the code with N = 10 and scale = 10.



```
Solution Explorer    add_loop_gpu_scaler.cu  ↔ X
(Global Scope)                                                    add(int * a, int * b, int * c)
//#include <cstdio>

#include "../common/book.h"

#define N    10
#define scale 10

__global__ void add(int *a, int *b, int *c) {
    int tid = blockIdx.x;    // this thread handles the data at its thread id
    if (tid < N)
        c[tid] = a[tid] + scale * b[tid];
}

int main(void) {
    int a[N], b[N], c[N];
    int *dev_a, *dev_b, *dev_c;

    // allocate the memory on the GPU
    HANDLE_ERROR(cudaMalloc((void**)&dev_a, N * sizeof(int)));
    HANDLE_ERROR(cudaMalloc((void**)&dev_b, N * sizeof(int)));
    HANDLE_ERROR(cudaMalloc((void**)&dev_c, N * sizeof(int)));

    // fill the arrays 'a' and 'b' on the CPU
    for (int i = 0; i<N; i++) {
        a[i] = -i;
        b[i] = i * i;
    }

    // copy the arrays 'a' and 'b' to the GPU
    cudaMemcpy(dev_a, a, N * sizeof(int),
        cudaMemcpyHostToDevice);
    cudaMemcpy(dev_b, b, N * sizeof(int),
        cudaMemcpyHostToDevice);

    add <<< N, 1 >> >(dev_a, dev_b, dev_c);

    // copy the array 'c' back from the GPU to the CPU
    cudaMemcpy(c, dev_c, N * sizeof(int),
        cudaMemcpyDeviceToHost);

    // display the results
    for (int i = 0; i<N; i++) {
        printf("%d + %d = %d\n", a[i], b[i], c[i]);
    }
```

```
C:\Windows\system32\cmd.exe
0 + 0 = 0
-1 + 1 = 9
-2 + 4 = 38
-3 + 9 = 87
-4 + 16 = 156
-5 + 25 = 245
-6 + 36 = 354
-7 + 49 = 483
-8 + 64 = 632
-9 + 81 = 801
Press any key to continue . . .
```

Please submit the source code in blackboard.