# Practice

## Lab 05

Make a program in c/c++ using CUDA in which you implement a kernel that adds two arrays (A and B) of integers generated randomly, considering the requirements:

- 32 threads

- A 1D grid with 4 blocks

- Use 1D blocks with 8 threads in each block

- Array size must be of 32

### Solution

```
1  #include "cuda_runtime.h"
2  #include "device_launch_parameters.h"
3  #include <stdio.h>
4  #include <stdlib.h> /* srand, rand */
5  #include <time.h> /* time */
6  __global__ void sumaVectores(int* vecA, int* vecB, int* vecRes
     ) {
7      int gId = threadIdx.x + blockDim.x * blockIdx.x;
8      vecRes[gId] = vecA[gId] + vecB[gId];
9  }
10 int main()
11 {
12     const int vectorSize = 32;
13     int* host_a = (int*)malloc(sizeof(int) * vectorSize);
14     int* host_b = (int*)malloc(sizeof(int) * vectorSize);
15     int* host_c = (int*)malloc(sizeof(int) * vectorSize);
16
17     int* dev_a, * dev_b, * dev_c;
18
19     cudaMalloc((void**)&dev_a, sizeof(int) * vectorSize);
20     cudaMalloc((void**)&dev_b, sizeof(int) * vectorSize);
21     cudaMalloc((void**)&dev_c, sizeof(int) * vectorSize);
22
23     srand(time(NULL));
24
25     for (int i = 0; i < vectorSize; i++) {
26         int num = rand() % vectorSize + 1;
27         host_a[i] = num;
28         num = rand() % vectorSize + 1;
29         host_b[i] = num;
30     }
31
```

```
32        cudaMemcpy(dev_a, host_a, sizeof(int) * vectorSize,
              cudaMemcpyHostToDevice);
33        cudaMemcpy(dev_b, host_b, sizeof(int) * vectorSize,
              cudaMemcpyHostToDevice);
34        dim3 grid(4, 1, 1);
35        dim3 block(8, 1, 1);
36
37        sumaVectores << < grid, block >> > (dev_a, dev_b, dev_c);
38
39        cudaMemcpy(host_c, dev_c, sizeof(int) * vectorSize,
              cudaMemcpyDeviceToHost);
40
41        printf("Vector A: \n");
42        for (int i = 0; i < vectorSize; i++) {
43            printf("%d ", host_a[i]);
44        }
45        printf("\nVector B: \n");
46        for (int i = 0; i < vectorSize; i++) {
47            printf("%d ", host_b[i]);
48        }
49        printf("\nVector C: \n");
50        for (int i = 0; i < vectorSize; i++) {
51            printf("%d ", host_c[i]);
52        }
53
54        free(host_a);
55        free(host_b);
56        free(host_c);
57        cudaFree(dev_a);
58        cudaFree(dev_b);
59        cudaFree(dev_c);
60        return 0;
61  }
```

## Lab 06

Make a program in c/c++ using CUDA in which you implement a kernel that assigns the different identification indexes for the threads into three arrays (A, B and C) considering the requirements:

- Use 3 arrays of integers with 64 elements each

- The kernel must manage 64 threads

- Each thread must write in an array a thread index:

    - The thread id (threadIdx.x) in array A

    - The block index (blockIdx.x) in array B

- The global index (globalId = thread.x + blockDim.x * blockIdx.x) in the C array

- The kernel must execute 3 times, each with the structure:

    - 1 block of 64 threads

    - 64 blocks of 1 thread

    - 4 blocks of 16 threads

- Print the results from main function, with a text indicating the config of each execution.

## Solution

```cpp
#include "cuda_runtime.h"
#include "device_launch_parameters.h"
#include <stdio.h>
#include <stdlib.h> /* srand, rand */
#include <time.h> /* time */


#include<iostream>
using namespace std;
__global__ void idKernel(int* vecA, int* vecB, int* vecC) {
    int gId = threadIdx.x + blockDim.x * blockIdx.x;

    vecA[gId] = threadIdx.x;
    vecB[gId] = blockIdx.x;
    vecC[gId] = gId;
}

void printArray(int* arr, int size, char * msg) {
    cout << msg << ": ";
    for (int i = 0; i < size; i++) {
        printf("%d ", arr[i]);
    }
    printf("\n");
}

int main()
{
    const int vectorSize = 64;
    int* host_a = (int*)malloc(sizeof(int) * vectorSize);
    int* host_b = (int*)malloc(sizeof(int) * vectorSize);
    int* host_c = (int*)malloc(sizeof(int) * vectorSize);

    int* dev_a, * dev_b, * dev_c;

    cudaMalloc((void**)&dev_a, sizeof(int) * vectorSize);
```

```
35        cudaMalloc((void**)&dev_b, sizeof(int) * vectorSize);
36        cudaMalloc((void**)&dev_c, sizeof(int) * vectorSize);
37
38        srand(time(NULL));
39
40        for (int i = 0; i < vectorSize; i++) {
41            host_a[i] = 0;
42            host_b[i] = 0;
43            host_c[i] = 0;
44        }
45
46        cudaMemcpy(dev_a, host_a, sizeof(int) * vectorSize,
              cudaMemcpyHostToDevice);
47        cudaMemcpy(dev_b, host_b, sizeof(int) * vectorSize,
              cudaMemcpyHostToDevice);
48        cudaMemcpy(dev_c, host_c, sizeof(int) * vectorSize,
              cudaMemcpyHostToDevice);
49
50        dim3 grid(1, 1, 1);
51        dim3 block(64, 1, 1);
52        idKernel << < grid, block >> > (dev_a, dev_b, dev_c);
53        cudaDeviceSynchronize(); // wait until kernel finishes and
                then come back to following code
54        cudaMemcpy(host_a, dev_a, sizeof(int) * vectorSize,
              cudaMemcpyDeviceToHost);
55        cudaMemcpy(host_b, dev_b, sizeof(int) * vectorSize,
              cudaMemcpyDeviceToHost);
56        cudaMemcpy(host_c, dev_c, sizeof(int) * vectorSize,
              cudaMemcpyDeviceToHost);
57
58        printf("Execution 1: 1 block 64 threads \n");
59        printArray(host_a, vectorSize, "threadIdx.x");
60        printArray(host_b, vectorSize, "blockIdx.x");
61        printArray(host_c, vectorSize, "globalId");
62
63        grid.x = 64; // (64, 1, 1)
64        block.x = 1; // (1, 1, 1)
65        idKernel << < grid, block >> > (dev_a, dev_b, dev_c);
66        cudaDeviceSynchronize();
67        cudaMemcpy(host_a, dev_a, sizeof(int) * vectorSize,
              cudaMemcpyDeviceToHost);
68        cudaMemcpy(host_b, dev_b, sizeof(int) * vectorSize,
              cudaMemcpyDeviceToHost);
69        cudaMemcpy(host_c, dev_c, sizeof(int) * vectorSize,
              cudaMemcpyDeviceToHost);
70
71        printf("\nExecution 2: 64 blocks 1 thread \n");
72        printArray(host_a, vectorSize, "threadIdx.x");
73        printArray(host_b, vectorSize, "blockIdx.x");
74        printArray(host_c, vectorSize, "globalId");
75
```

```
76        grid.x = 4;
77        block.x = 16;
78        idKernel << < grid , block >> > (dev_a , dev_b , dev_c);
79        cudaDeviceSynchronize ();
80        cudaMemcpy(host_a , dev_a , sizeof(int) * vectorSize ,
                cudaMemcpyDeviceToHost);
81        cudaMemcpy(host_b , dev_b , sizeof(int) * vectorSize ,
                cudaMemcpyDeviceToHost);
82        cudaMemcpy(host_c , dev_c , sizeof(int) * vectorSize ,
                cudaMemcpyDeviceToHost);
83
84        printf("\nExecution 3: 4 block 16 threads \n");
85        printArray(host_a , vectorSize , "threadIdx.x");
86        printArray(host_b , vectorSize , "blockIdx.x");
87        printArray(host_c , vectorSize , "globalId");
88
89        free(host_a);
90        free(host_b);
91        free(host_c);
92        cudaFree(dev_a);
93        cudaFree(dev_b);
94        cudaFree(dev_c);
95        return 0;
96  }
```
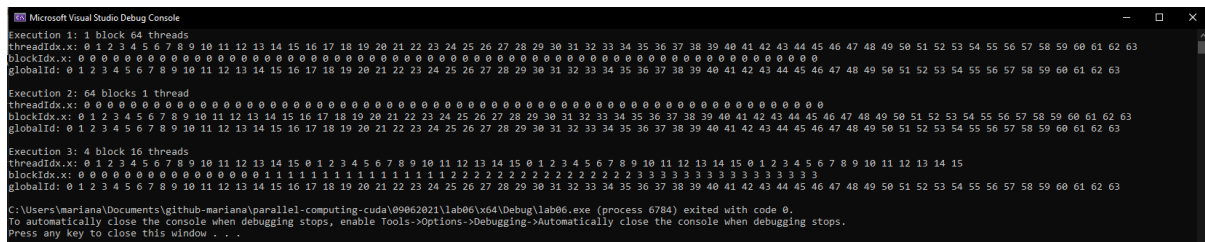
**Output**



**Figure 1:** Image

**Observations**

- To use one kernel, you need to calculate the globalId as `threadIdx.x` + `blockDim.x` * `blockIdx.x` for **all** configurations, because this small formula works for a config of a 1D block along 1 dimension (axis), or more 1D blocks along that same dimension (axis).

- Whenever we launch a kernel, the instruction of the kernel start to execute in the device, but the program execution stream in main continues with the next lines, so the launch of following kernels can get mixed up, especially in printings. In order to

tell the program to **wait until a kernel execution finishes** in order to continue with the next host lines, you use the function `cudaDeviceSynchronize()` right after the launch of each kernel. With this, we are telling the host to make a pause in that line until the kernel finishes so that it can continue with the following lines. This pause means to **synchronize the executions**. Thus, this function is used whenever we want a program to wait for a kernel to finish before it continues its stream.