

Universidad Nacional de San Agustín de Arequipa

AÑO DE LA RECUPERACIÓN Y CONSOLIDACIÓN DE LA  
ECONOMÍA PERUANA



ESCUELA PROFESIONAL DE CIENCIA DE LA COMPUTACIÓN  
Temas en Inteligencia Artificial

---

# PRACTICA 1

CUDA PROGRAMMING

A-2025

---

**Profesor: MSc. PERCY MALDONADO QUISPE**

Javier Wilber Quispe Rojas

# Índice

Índice	1
1 Get information about the device where you are running the code.	2
2 Implements vector addition in CUDA.	2
3 Implements matrix multiplication in CUDA.	3
4 Link de colab	5

## 1. Get information about the device where you are running the code.

```
1  #include <stdio.h>
2
3  int main()
4  {
5      int noOfDevices;
6
7      cudaGetDeviceCount (&noOfDevices);
8
9      cudaDeviceProp prop;
10     for (int i = 0; i < noOfDevices; i++)
11     {
12
13         cudaGetDeviceProperties (&prop, i);
14
15         printf("Device-Name:\t%s\n", prop.name);
16         printf("Total-global-memory:\t%d\n", prop.totalGlobalMem);
17         printf("No.-of-SMs:\t%d\n", prop.multiProcessorCount);
18         printf("Shared-memory-/SM:\t%d\n", prop.sharedMemPerBlock);
19         printf("Registers-/SM:\t%d\n", prop.regsPerBlock);
20     }
21
22     return 1;
23 }
```

Listing 1: Codigo ejercicio 1

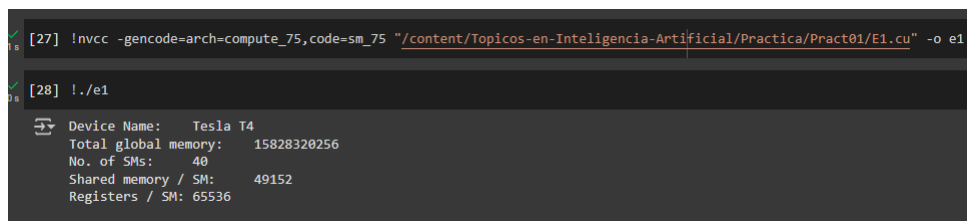


Figura 1: Ejecución ejercicio 1

## 2. Implements vector addition in CUDA.

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <cuda_runtime.h>
4
5
6  __global__ void suma_vectores(float *c, float *a, float *b, int N) {
7      int idx = blockIdx.x * blockDim.x + threadIdx.x;
8      if (idx < N) {
9          c[idx] = a[idx] + b[idx];
10     }
11 }
12
13 int main(void) {
14
15     float *a_h, *b_h, *c_h;
16
17     float *a_d, *b_d, *c_d;
18     int N=10;
19     size_t size = N * sizeof(float);
20
21
22     a_h = (float *)malloc(size);
23     b_h = (float *)malloc(size);
24     c_h = (float *)malloc(size);
25
26 }
```

```

27     for (int i = 0; i < N; i++) {
28         a_h[i] = (float)i;
29         b_h[i] = (float)(i + 1);
30     }
31
32     printf("\nArreglo a:\n");
33     for (int i = 0; i < N; i++) printf("%f-", a_h[i]);
34     printf("\nArreglo b:\n");
35     for (int i = 0; i < N; i++) printf("%f-", b_h[i]);
36
37
38     cudaMalloc((void **) &a_d, size);
39     cudaMalloc((void **) &b_d, size);
40     cudaMalloc((void **) &c_d, size);
41
42
43     cudaMemcpy(a_d, a_h, size, cudaMemcpyHostToDevice);
44     cudaMemcpy(b_d, b_h, size, cudaMemcpyHostToDevice);
45
46
47     int block_size = 8;
48     int n_blocks = (N + block_size - 1) / block_size;
49     suma_vectores<<<n_blocks, block_size>>>(c_d, a_d, b_d, N);
50
51
52     cudaMemcpy(c_h, c_d, size, cudaMemcpyDeviceToHost);
53
54
55     printf("\nResultado c:\n");
56     for (int i = 0; i < N; i++) printf("%f-", c_h[i]);
57     printf("\n");
58
59
60     free(a_h);
61     free(b_h);
62     free(c_h);
63
64
65     cudaFree(a_d);
66     cudaFree(b_d);
67     cudaFree(c_d);
68
69     return 0;
70 }

```

Listing 2: Código ejercicio 2

```

[32] nvcc -gencode=arch=compute_75,code=sm_75 "/content/Topicos-en-Inteligencia-Artificial/Practica/Pract01/E2.cu" -o e2
[33] !./e2
Arreglo a:
0.000000 1.000000 2.000000 3.000000 4.000000 5.000000 6.000000 7.000000 8.000000 9.000000
Arreglo b:
1.000000 2.000000 3.000000 4.000000 5.000000 6.000000 7.000000 8.000000 9.000000 10.000000
Resultado c:
1.000000 3.000000 5.000000 7.000000 9.000000 11.000000 13.000000 15.000000 17.000000 19.000000

```

Figura 2: Ejecución ejercicio 2

### 3. Implements matrix multiplication in CUDA.

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <cuda_runtime.h>
4
5 #define BLOCK_SIZE 16
6
7 int div_up(int x, int y) {
8     return (x + y - 1) / y;
9 }

```

```

10
11 // Funci n Kernel que se ejecuta en el Device
12 __global__ void Multiplica_Matrices_GM(float *C, float *A, float *B, int nfil, int ncol)
13 {
14     int idx = blockIdx.x * blockDim.x + threadIdx.x;
15     int idy = blockIdx.y * blockDim.y + threadIdx.y;
16     int index = idy * ncol + idx;
17
18     if (idy < nfil && idx < ncol) {
19         float sum = 0.0f;
20         for (int k = 0; k < ncol; k++) {
21             sum += A[idy * ncol + k] * B[k * ncol + idx];
22         }
23         C[index] = sum;
24     }
25 }
26
27 int main(void) {
28     float *A_h, *B_h, *C_h;
29
30     float *A_d, *B_d, *C_d;
31
32     int nfil = 5;
33     int ncol = 5;
34     int N = nfil * ncol;
35
36     size_t size = N * sizeof(float);
37
38     cudaEvent_t start, stop;
39     float time;
40     cudaEventCreate(&start);
41     cudaEventCreate(&stop);
42
43
44     A_h = (float *)malloc(size);
45     B_h = (float *)malloc(size);
46     C_h = (float *)malloc(size);
47
48
49     for (int i = 0; i < nfil; i++) {
50         for (int j = 0; j < ncol; j++) {
51             A_h[i * ncol + j] = 1.0f;
52             B_h[i * ncol + j] = 2.0f;
53         }
54     }
55
56
57     cudaMalloc((void **) &A_d, size);
58     cudaMalloc((void **) &B_d, size);
59     cudaMalloc((void **) &C_d, size);
60
61
62     cudaMemcpy(A_d, A_h, size, cudaMemcpyHostToDevice);
63     cudaMemcpy(B_d, B_h, size, cudaMemcpyHostToDevice);
64
65
66     dim3 block_size(BLOCK_SIZE, BLOCK_SIZE);
67     dim3 n_blocks(div_up(ncol, BLOCK_SIZE), div_up(nfil, BLOCK_SIZE));
68     cudaEventRecord(start);
69     Multiplica_Matrices_GM<<<n_blocks, block_size>>>(C_d, A_d, B_d, nfil, ncol);
70     cudaEventRecord(stop);
71
72
73     cudaEventSynchronize(stop);
74
75
76     cudaMemcpy(C_h, C_d, size, cudaMemcpyDeviceToHost);
77
78
79     printf("\nMatriz -A:\n");
80     for (int i = 0; i < nfil; i++) {
81         for (int j = 0; j < ncol; j++) {

```

```

82         printf(" %2.2f-", A.h[i * ncol + j]);
83     }
84     printf("\n");
85 }
86
87 printf("\nMatriz-B:\n");
88 for (int i = 0; i < nfil; i++) {
89     for (int j = 0; j < ncol; j++) {
90         printf(" %2.2f-", B.h[i * ncol + j]);
91     }
92     printf("\n");
93 }
94
95 printf("\nMatriz-C:\n");
96 for (int i = 0; i < nfil; i++) {
97     for (int j = 0; j < ncol; j++) {
98         printf(" %2.2f-", C.h[i * ncol + j]);
99     }
100    printf("\n");
101 }
102
103
104 cudaEventElapsedTime(&time, start, stop);
105 printf("Tiempo-de-ejecuci n:- %3.1f-ms\n", time);
106
107 free(A.h);
108 free(B.h);
109 free(C.h);
110
111
112
113 cudaFree(A_d);
114 cudaFree(B_d);
115 cudaFree(C_d);
116
117 return 0;
118 }

```

Listing 3: Codigo ejercicio 3

```

[34] !nvcc -gencode=arch=compute_75,code=sm_75 "/content/Topicos-en-Inteligencia-Artificial/Practica/Pract01/E3.cu" -o e3
!./e3
Matriz A:
1.00 1.00 1.00 1.00 1.00
1.00 1.00 1.00 1.00 1.00
1.00 1.00 1.00 1.00 1.00
1.00 1.00 1.00 1.00 1.00
1.00 1.00 1.00 1.00 1.00

Matriz B:
2.00 2.00 2.00 2.00 2.00
2.00 2.00 2.00 2.00 2.00
2.00 2.00 2.00 2.00 2.00
2.00 2.00 2.00 2.00 2.00
2.00 2.00 2.00 2.00 2.00

Matriz C:
10.00 10.00 10.00 10.00 10.00
10.00 10.00 10.00 10.00 10.00
10.00 10.00 10.00 10.00 10.00
10.00 10.00 10.00 10.00 10.00
10.00 10.00 10.00 10.00 10.00
Tiempo de ejecución: 0.1 ms

```

Figura 3: Ejecución ejercicio 3

## 4. Link de colab

Para poder ejecutar los codigos en el colab debe cambiar entorno de ejecucion **GPU T4** , luego ejecutar todo.

[https://colab.research.google.com/drive/10eqldS\\_mWLVCoVe\\_-3lTho44W\\_QBz6rC?usp=sharing](https://colab.research.google.com/drive/10eqldS_mWLVCoVe_-3lTho44W_QBz6rC?usp=sharing)