



Práctica 2

Parte 1: Compilando Programas en C

Máster Universitario en Ingeniería de Computadores y Redes

1 de diciembre de 2025

Pablo González Segarra

Índice

Ejercicio 1	2
Ejercicio 3	3
Ejercicio 4	5
Ejercicio 5	7
Ejercicio 6	7
Ejercicio 7	8
Ejercicio 8	8

Ejercicio 1 y 2

Debido a la similitud entre los ejercicios 1 y 2, se presentan juntos en esta sección. A continuación se muestra el código modificado necesario para ejecutar el programa en la GPU utilizando CUDA y comprobar que los resultados son correctos.

```
1 #define N 1000000000
2
3 #include <stdio.h>
4 #include <stdlib.h>
5 #include <math.h>
6 #include <cuda_runtime.h>
7
8 __global__ void vector_add(float *out, float *a, float *b, int n) {
9     for(int i = 0; i < n; i++){
10         out[i] = a[i] + b[i];
11     }
12 }
13
14 int main(){
15     float *a, *b, *out;
16
17     // Allocate memory
18     a = (float*)malloc(sizeof(float) * N);
19     b = (float*)malloc(sizeof(float) * N);
20     out = (float*)malloc(sizeof(float) * N);
21
22     // Initialize array
23     for(int i = 0; i < N; i++){
24         a[i] = 1.0f; b[i] = 2.0f;
25     }
26
27     float *a_cuda = NULL, *b_cuda = NULL, *out_cuda = NULL;
28
29     // Allocate device memory
30     cudaMalloc((void**)&a_cuda, sizeof(float) * N);
31     cudaMalloc((void**)&b_cuda, sizeof(float) * N);
32     cudaMalloc((void**)&out_cuda, sizeof(float) * N);
33
34     // Copy inputs to device
35     cudaMemcpy(a_cuda, a, sizeof(float) * N, cudaMemcpyHostToDevice);
36     cudaMemcpy(b_cuda, b, sizeof(float) * N, cudaMemcpyHostToDevice);
37
38     // Main function
39     vector_add<<<1,1>>>(out_cuda, a_cuda, b_cuda, N);
40
41     cudaMemcpy(out, out_cuda, sizeof(float) * N,
42     cudaMemcpyDeviceToHost);
43
44     cudaDeviceSynchronize();
45
46     // Verify result
47     const float MAX_ERR = 1e-6f;
48     for(int i = 0; i < N; i++){
49         if(fabsf(out[i] - 3.0f) > MAX_ERR){
50             printf("Error at index %d: %f (diff=%f)\n", i, out[i],
51             fabsf(out[i] - 3.0f));
```

```

50         return -1;
51     }
52 }
53
54 printf("Success! All values are correct.\n");
55
56 // Free memory
57 cudaFree(a_cuda);
58 cudaFree(b_cuda);
59 cudaFree(out_cuda);
60
61 free(a);
62 free(b);
63 free(out);
64 }
```

Listing 1: Código CUDA modificado para ejecutar en GPU

Ejercicio 3

El código modificado para poder obtener los tiempos de ejecución es el siguiente:

```

1 #define N 100000000
2
3 #include <stdio.h>
4 #include <stdlib.h>
5 #include <math.h>
6 #include <cuda_runtime.h>
7
8 __global__ void vector_add(float *out, float *a, float *b, int n) {
9     for(int i = 0; i < n; i++){
10         out[i] = a[i] + b[i];
11     }
12 }
13
14 int main(){
15     float *a, *b, *out;
16
17     // Allocate memory
18     a = (float*)malloc(sizeof(float) * N);
19     b = (float*)malloc(sizeof(float) * N);
20     out = (float*)malloc(sizeof(float) * N);
21
22     // Initialize array
23     for(int i = 0; i < N; i++){
24         a[i] = 1.0f; b[i] = 2.0f;
25     }
26
27     float *a_cuda = NULL, *b_cuda = NULL, *out_cuda = NULL;
28
29     // Create events for timing
30     cudaEvent_t start, end;
31     cudaEventCreate(&start);
32     cudaEventCreate(&end);
```

```

34     // Start timing
35     cudaEventRecord(start);
36
37     // Allocate device memory
38     cudaMalloc((void**)&a_cuda, sizeof(float) * N);
39     cudaMalloc((void**)&b_cuda, sizeof(float) * N);
40     cudaMalloc((void**)&out_cuda, sizeof(float) * N);
41
42     // Copy inputs to device
43     cudaMemcpy(a_cuda, a, sizeof(float) * N, cudaMemcpyHostToDevice);
44     cudaMemcpy(b_cuda, b, sizeof(float) * N, cudaMemcpyHostToDevice);
45
46     // Main function
47     vector_add<<<1,1>>>(out_cuda, a_cuda, b_cuda, N);
48
49     cudaMemcpy(out, out_cuda, sizeof(float) * N,
50     cudaMemcpyDeviceToHost);
51
52     // End timing
53     cudaEventRecord(end);
54     cudaEventSynchronize(end);
55     float milliseconds = 0;
56     cudaEventElapsedTime(&milliseconds, start, end);
57     printf("Time elapsed: %f ms\n", milliseconds);
58
59     // Verify result
60     const float MAX_ERR = 1e-6f;
61     for(int i = 0; i < N; i++){
62         if(fabsf(out[i] - 3.0f) > MAX_ERR){
63             printf("Error at index %d: %f (diff=%f)\n", i, out[i],
64             fabsf(out[i] - 3.0f));
65             return -1;
66         }
67     }
68
69     printf("Success! All values are correct.\n");
70
71     // Free memory
72     cudaFree(a_cuda);
73     cudaFree(b_cuda);
74     cudaFree(out_cuda);
75
76     free(a);
77     free(b);
78     free(out);
79 }
```

Listing 2: Código CUDA modificado para medir tiempos de ejecución

Se ha ejecutado 5 veces para evitar variaciones ajenas al rendimiento del código. El tiempo de ejecucion tiene en cuenta tanto la transferencia de datos entre host y device como la ejecución del kernel.

Resultado de las ejecuciones:

1. 5887.4 ms
2. 5864.1 ms

3. 5918.6 ms
4. 5856.9 ms
5. 5868.1 ms

Siendo el mejor tiempo 5856.9 ms.

Ejercicio 4

Para comparar la aceleración obtenida al ejecutar el código en la GPU frente a la CPU, primero se ha compilado y ejecutado la versión original en CPU, obteniendo un tiempo mejor de 992.6 ms.

Para obtener los tiempos de ejecución en GPU con distintas configuraciones numero de bloques y hilos por bloque, se ha modificado el código del ejercicio 3 para variar estos parámetros. Todas las configuraciones se han ejecutado multiples veces y se ha tomado el mejor tiempo de ejecución. En la tabla 1 se muestran los resultados obtenidos.

Versión	Tiempo de ejecución (ms)	Speedup GPU vs CPU
CPU (original)	992.6	1.00
1 thread / block	532.3	1.86
16 threads / block	438.2	2.27
32 threads / block	431.8	2.30
256 threads / block	428.4	2.32

Tabla 1: Tiempos de ejecución y speedup para distintas configuraciones de hilos/ bloque

Como se puede observar la parallelización en GPU ofrece una mejora significativa en el tiempo de ejecución en comparación con la versión en CPU. El mejor rendimiento se obtiene con 256 hilos por bloque, logrando un speedup de aproximadamente 2.32 veces respecto a la ejecución en CPU. Sin embargo, la mejora a partir de 16 hilos por bloque es marginal.

A continuación se muestra el código modificado para variar el número de bloques e hilos por bloque:

```

1 #define N 100000000
2
3 #include <stdio.h>
4 #include <stdlib.h>
5 #include <math.h>
6 #include <cuda_runtime.h>
7
8 __global__ void vector_add(float *out, float *a, float *b, int n) {
9     int i = blockIdx.x * blockDim.x + threadIdx.x;
10    if (i < n) {
11        out[i] = a[i] + b[i];
12    }
13 }
14
15 int main(int argc, char **argv){
16

```

```
17     float *a, *b, *out;
18
19     // Allocate memory
20     a = (float*)malloc(sizeof(float) * N);
21     b = (float*)malloc(sizeof(float) * N);
22     out = (float*)malloc(sizeof(float) * N);
23
24     // Initialize array
25     for(int i = 0; i < N; i++){
26         a[i] = 1.0f; b[i] = 2.0f;
27     }
28
29     float *a_cuda = NULL, *b_cuda = NULL, *out_cuda = NULL;
30
31     // Create events for timing
32     cudaEvent_t start, end;
33     cudaEventCreate(&start);
34     cudaEventCreate(&end);
35
36     // Start timing
37     cudaEventRecord(start);
38
39     // Allocate device memory
40     cudaMalloc((void**)&a_cuda, sizeof(float) * N);
41     cudaMalloc((void**)&b_cuda, sizeof(float) * N);
42     cudaMalloc((void**)&out_cuda, sizeof(float) * N);
43
44     // Copy inputs to device
45     cudaMemcpy(a_cuda, a, sizeof(float) * N, cudaMemcpyHostToDevice);
46     cudaMemcpy(b_cuda, b, sizeof(float) * N, cudaMemcpyHostToDevice);
47
48     // Main function
49     // threadsPerBlock is taken from argv[1] (assume valid integer
50     // provided)
51     int threadsPerBlock = atoi(argv[1]);
52     int blocks = (N + threadsPerBlock - 1) / threadsPerBlock;
53     // Print configuration
54     printf("Using %d blocks of %d threads\n", blocks, threadsPerBlock);
55
56     vector_add<<<blocks, threadsPerBlock>>>(out_cuda, a_cuda, b_cuda,
57     N);
58
59     cudaMemcpy(out, out_cuda, sizeof(float) * N,
60     cudaMemcpyDeviceToHost);
61
62     // End timing
63     cudaEventRecord(end);
64     cudaEventSynchronize(end);
65     float milliseconds = 0;
66     cudaEventElapsedTime(&milliseconds, start, end);
67     printf("Time elapsed: %f ms\n", milliseconds);
68
69     // Verify result
70     const float MAX_ERR = 1e-6f;
71     for(int i = 0; i < N; i++){
72         if(fabsf(out[i] - 3.0f) > MAX_ERR){
73             printf("Error at index %d: %f (diff=%f)\n", i, out[i],
74             fabsf(out[i] - 3.0f));
```

```

70         return -1;
71     }
72 }
73 printf("Success! All values are correct.\n");
75
76 // Free memory
77 cudaFree(a_cuda);
78 cudaFree(b_cuda);
79 cudaFree(out_cuda);
80
81 free(a);
82 free(b);
83 free(out);
84 }
```

Listing 3: Código CUDA modificado para variar bloques e hilos por bloque

Ejercicio 5

En este ejercicio se ha creado un kernel para la función `void heavy_cpu(float* data, int n)`, el código se puede ver a continuación:

```

1 __global__ void heavy_cpu(float* data, int n) {
2     int i = blockIdx.x * blockDim.x + threadIdx.x;
3
4     if (i < n) {
5         float x = data[i];
6         for (int j = 0; j < 10000; j++) {
7             x = sinf(x) * 1.00001f + cosf(x) * 0.99999f;
8         }
9         data[i] = x;
10    }
11 }
```

Listing 4: Código CUDA del kernel para la función `heavy_cpu`

Ejercicio 6

Para este ejercicio se ha usado la misma metodología que en los ejercicios anteriores, pero esta vez se ha medido solo el tiempo de ejecución del kernel, sin contar las transferencias de datos entre host y device. Se ha ejecutado la función en CPU y GPU, obteniendo los siguientes tiempos:

- CPU: 224,69 s
- GPU: 10,33 s

Suponiendo que la función `heavy_cpu` contiene 5 operaciones de coma flotante por iteración del bucle, y considerando que el bucle se ejecuta 10.000 veces por cada elemento del array, podemos calcular el número total de operaciones de coma flotante realizadas

durante la ejecución de la función. Dado que el array tiene un tamaño de $N = 1.000.000$, el número total de operaciones de coma flotante es:

$$\text{Total de operaciones} = N \times 10,000 \times 5 = 1,000,000 \times 10,000 \times 5 = 5 \times 10^{10} \text{ operaciones}$$

Con estos datos, podemos calcular el rendimiento en GFLOPS para ambas ejecuciones:

- CPU: $\frac{5 \times 10^{10} \text{ operaciones}}{224,69 \text{ s}} \approx 0,2226 \text{ GFLOPS}$
- GPU: $\frac{5 \times 10^{10} \text{ operaciones}}{10,33 \text{ s}} \approx 4,8385 \text{ GFLOPS}$

Estos resultados muestran una mejora significativa en el rendimiento al ejecutar la función en la GPU en comparación con la CPU. Obteniendo una aceleración de aproximadamente 21.6 veces.

Ejercicio 7

Cada hilo ejecuta el procesamiento de un elemento del array, es decir, obtiene su índice, lee `data[i]`, aplica el bucle de 10.000 iteraciones con el grueso de la carga computacional, y escribe el resultado de vuelta.

En cuanto al paralelismo, se podría mejorar ligeramente retocando algunos detalles como los parámetros de threads per block o blocks, pero en general no es posible aumentar significativamente el paralelismo porque cada iteración del bucle interno depende del valor anterior, por lo que esta dependencia impide parallelizar el bucle.

Ejercicio 8