

Proyecto final Histograma realizado con CUDA

Arquitecturas Avanzadas y de Propósito Específico

Cheuk Kelly Ng Pante (alu0101364544@ull.edu.es)

11 de junio de 2024

Índice general

1. Introducción	1
2. Desarrollo del proyecto	1
2.1. Implementación base	1
2.2. Segunda implementación	1
3. Pruebas realizadas	2
3.1. Tiempos promedio de ejecución	3
3.2. Tiempos máximo de ejecución	4
3.3. Tiempos mínimo de ejecución	5

1. Introducción

Consiste realizar un histograma de un vector V de un número elevado N de elementos enteros aleatorios. El histograma consiste en un vector H que tiene M elementos que representan “cajas”. En cada caja se cuenta el número de veces que ha aparecido un elemento del vector V con el valor adecuado para asignarlo a esa caja (normalmente cada caja representa un rango o intervalo de valores). En nuestro caso, para simplificar la asignación del elemento de V a su caja correspondiente del histograma, vamos a realizar la operación ValorElementoV Módulo M , que nos da directamente el índice de la caja del histograma a la que pertenecerá ese elemento y cuyo contenido deberemos incrementar. Se sugiere como N un valor del orden de millones de elementos y como M , 8 cajas.

2. Desarrollo del proyecto

2.1. Implementación base

Como implementación base se pide crear tantos hilos como elementos V para cada uno se encargue de ir al elemento que le corresponda V e incremente la caja correcta en el vector histograma H (posiblemente de forma atómica).

El código base es el siguiente:

```
1 __global__ void histogram() {
2     int i = blockIdx.x * blockDim.x + threadIdx.x;
3     if (i < N) {
4         atomicAdd(&vector_H[vector_V[i] % M], 1);
5     }
6 }
```

La implementación base lo que hace es crear un único histograma compartido por todos los hilos. Cada hilo se encarga de incrementar la caja correspondiente al elemento V que le corresponda. Para ello, se utiliza la función `atomicAdd` para incrementar de forma atómica el valor de la caja correspondiente al elemento V . La implementación se encuentra en el archivo *histogram_1.cu*.

2.2. Segunda implementación

Como segunda implementación, se va a dividir la operación en dos fases. En la primera, en lugar de trabajar sobre un único histograma global, repartiremos el cálculo realizando un cierto número de histogramas que llamaremos “locales”, cada uno calculado sobre una parte del vector de datos. La idea es reducir el número de hilos que escriben sobre la misma posición del histograma, ya que dicha operación debe ser atómica y se serializan dichos accesos. La segunda fase realizará la suma de los histogramas locales en un único histograma global final. Se debe intentar llevar a cabo esta suma de la forma más paralela o eficiente, posiblemente utilizando el método de reducción. La implementación se encuentra en el archivo *histogram_2.cu*.

El código es el siguiente:

```
1 __global__ void histogram() {
2
3     extern __shared__ int local_histogram[];
4
5     int tid = threadIdx.x;
```

```

6  for (int i = tid; i < M; i += blockDim.x) {
7      local_histogram[i] = 0;
8  }
9  __syncthreads();
10
11  int i = blockIdx.x * blockDim.x + tid;
12  if (i < N) {
13      atomicAdd(&local_histogram[vector_V[i] % M], 1);
14  }
15  __syncthreads();
16
17  for (int j = tid; j < M; j += blockDim.x) {
18      atomicAdd(&vector_H[j], local_histogram[j]);
19  }
20 }

```

La segunda implementación se divide en dos fases.

3. Pruebas realizadas

Durante las pruebas realizadas, se he ejecutado cada implementación 10000 veces para obtener un tiempo promedio significativo de tiempo. También, se ha obtenido el tiempo máximo y mínimo de la ejecución.

3.1. Tiempos promedio de ejecución

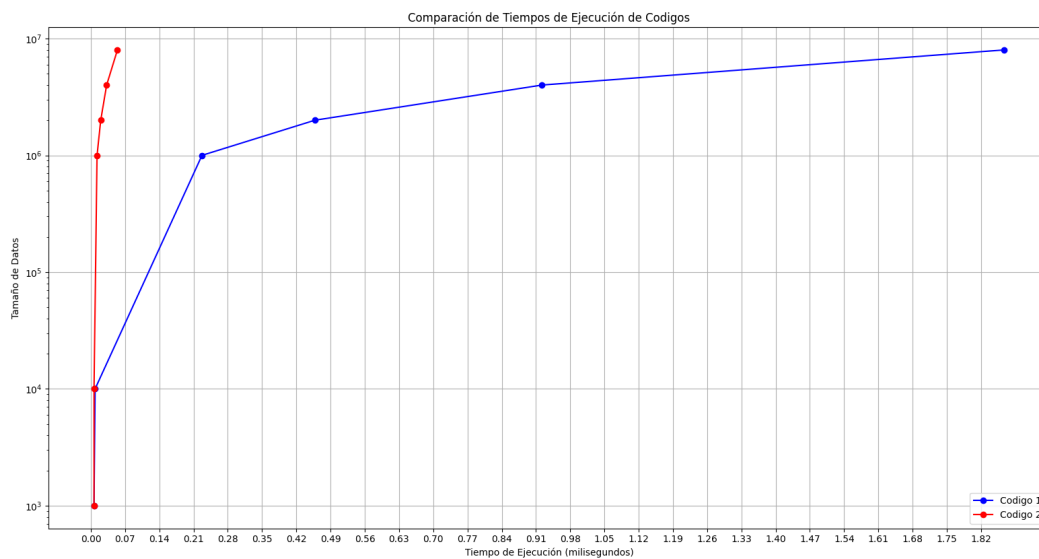


Figura 3.1: Tiempo promedio de ejecución

Tamaño de datos	Tiempo Código 1	Tiempo Código 2
1000	0.00642909	0.00644173
10000	0.00865969	0.00644241
1000000	0.227177	0.0126207
2000000	0.458186	0.0201534
4000000	0.922377	0.0322431
8000000	1.86746	0.0543328

Cuadro 3.1: Tiempo promedio de ejecución

3.2. Tiempos máximo de ejecución

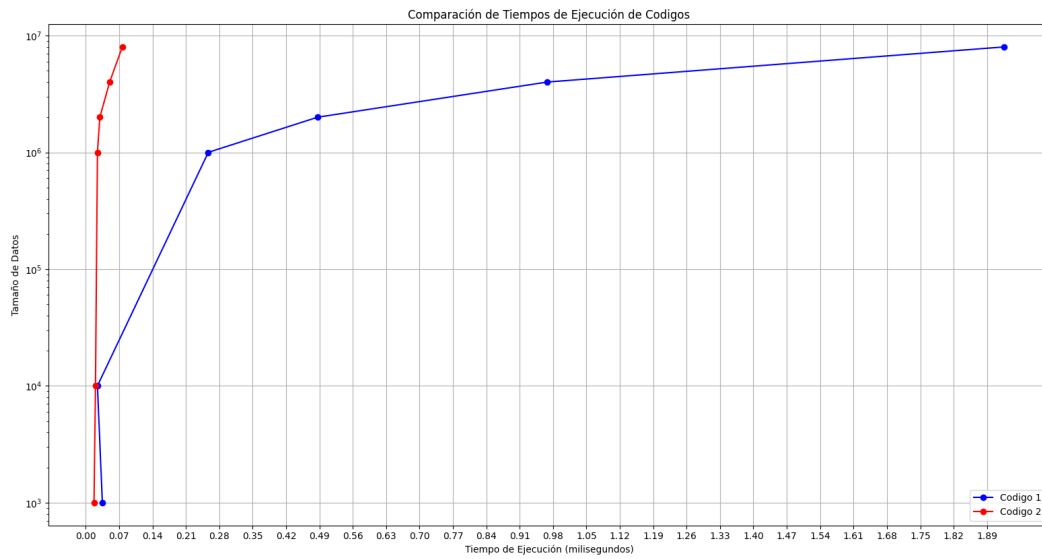


Figura 3.2: Tiempo máximo de ejecución

Tamaño de datos	Tiempo Código 1	Tiempo Código 2
1000	0.034912	0.017408
10000	0.024576	0.02048
1000000	0.257024	0.024576
2000000	0.4864	0.029664
4000000	0.96768	0.050176
8000000	1.92614	0.076736

Cuadro 3.2: Tiempo máximo de ejecución

3.3. Tiempos mínimo de ejecución

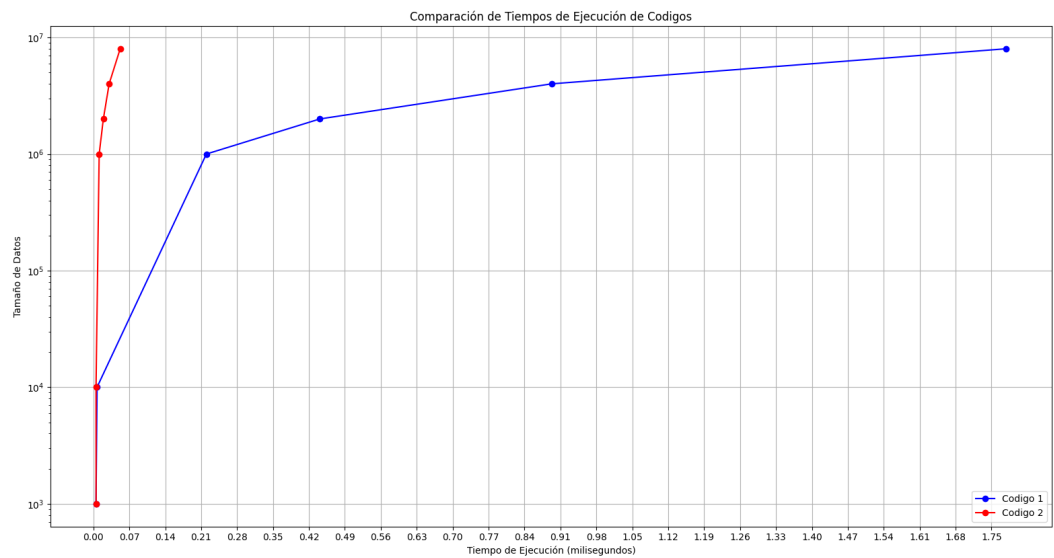


Figura 3.3: Tiempo mínimo de ejecución

Tamaño de datos	Tiempo Código 1	Tiempo Código 2
1000	0.00512	0.00528
10000	0.007168	0.005152
1000000	0.22016	0.011264
2000000	0.441344	0.019424
4000000	0.892928	0.03072
8000000	1.77869	0.052224

Cuadro 3.3: Tiempo mínimo de ejecución