

Prueba GPU - 1

Fabián Trigo

Universidad de Valparaíso, Valparaíso

(Dated: Septiembre 12, 2023)

Prueba de GPU, se comienza con unas preguntas sobre el código tratado

```
1 __global__ void vecAddKernel(float* A, float* B, float* C, int n){
2     int i = threadIdx.x + blockDim.x * blockIdx.x * 2;
3     if (i < n){ C[i] = A[i] + B[i] };
4     i += blockDim.x;
5     if (i < n){ C[i] = A[i] + B[i] };
6 }
7
8 int vecAdd(float* A, float* B, float* C, int n){
9     // "n" es el número de elementos en A, B y C
10    int size = n * sizeof(float);
11    float* A_d; float* B_d; float* C_d;
12    cudaMalloc((void **)&A_d, size);
13    cudaMalloc((void **)&B_d, size);
14    cudaMalloc((void **)&C_d, size);
15    cudaMemcpy(A_d, A, size, cudaMemcpyHostToDevice);
16    cudaMemcpy(B_d, B, size, cudaMemcpyHostToDevice);
17    vecAddKernel<<<ceil(n / 2048.0), 1024>>>>(A_d, B_d, C_d, n);
18    cudaMemcpy(C, C_d, size, cudaMemcpyDeviceToHost);
19 }
```

Observando el código a simple vista, tenemos un `vecAddKernel` que toma los arrays A,B,C.

I. PREGUNTA 1

1. Si el tamaño de los arrays A, B y C es 50.000 elementos cada uno, ¿cuántos bloques de threads están generados? Se generan en la línea: `vecAddKernel<<<ceil(n / 2048.0), 1024>>>>(A_d, B_d, C_d, n)`; Téngase en cuenta las cantidades `n=50000`; `ceil(n/2048) => 25` Osea que tendríamos asignados unos 25 bloques, cada bloque tendría unas 1024 threads por bloque.

II. PREGUNTA 2

2. Si el tamaño de los arrays A, B y C es 50.000 elementos cada uno, ¿cuántos warps hay en cada bloque de threads?

En la arquitectura de Nvidia, cada bloque posee warps de 32 threads que se ejecutan de manera paralela. Debido a que hay unos 1024 threads por bloque, tendríamos unos

$$\text{warp} = \frac{1024}{32} = 32$$

por bloque.

III. PREGUNTA 3

Si el tamaño de los arrays A, B y C es 50.000 elementos cada uno, ¿cuántos threads en total tendrá la grilla?

La grilla es este conjunto completo de bloques y threads, para ello bastará con que multipliquemos la cantidad de threads por bloque por la cantidad de bloques

$$\text{threads totales} = \frac{\text{threads}}{\text{por bloque}} \cdot \frac{\text{bloques}}{=} 1024 \cdot 25 = 25600 \text{ threads en grilla}$$

IV. PREGUNTA 4

¿Cuántos sumas calcula cada thread en el kernel?

Como podemos observar en las líneas 1 a 6 del código, al ingresar se encuentra el ‘threadIdx.x’ así como el ‘blockIdx.x’, esto posiciona al thread en un lugar específico de la memoria, es importante el hecho de que el ‘blockIdx.x’ está multiplicado por 2, quiere decir que es el tamaño de steps que damos en los bloques.

Cada thread ejecutará la función ‘voidAddKernel’, lo que quiere decir que tomara su index (threadId, blockDim) y el siguiente (threadId, blockDim+1), por tanto realizará 2 cálculos como mucho.

V. PREGUNTA 5

Escribir otra versión del código arriba tal que cada thread calcula 4 sumas y por lo tanto puede lanzar el kernel con la mitad del número de threads usado en el código arriba.

```

1  __global__ void vecAddKernel(float* A, float* B, float* C, int n){
2      int i = threadIdx.x + blockDim.x * blockIdx.x * 4;
3      if (i < n){ C[i] = A[i] + B[i] }; // (x,y)
4      i += blockDim.x;
5      if (i < n){ C[i] = A[i] + B[i] }; // (x,y+1)
6      i += blockDim.x;
7      if (i < n){ C[i] = A[i] + B[i] }; // (x,y+2)
8      i += blockDim.x;
9      if (i < n){ C[i] = A[i] + B[i] }; // (x,y+3)
10 }
11
12 int vecAdd(float* A, float* B, float* C, int n){
13     int size = n * sizeof(float);
14     float* A_d; float* B_d; float* C_d;
15     cudaMalloc((void **)&A_d, size);
16     cudaMalloc((void **)&B_d, size);
17     cudaMalloc((void **)&C_d, size);
18     cudaMemcpy(A_d, A, size, cudaMemcpyHostToDevice);
19     cudaMemcpy(B_d, B, size, cudaMemcpyHostToDevice);
20     // cada bloque se lanzaría con la mitad del número de threads como es pedido
21     // lo que equivaldría a tener 16 warps, con sus 32 threads corriendo en paralelo
22     vecAddKernel<<<ceil(n / 2048.0), 512>>>(A_d, B_d, C_d, n);
23
24     cudaMemcpy(C, C_d, size, cudaMemcpyDeviceToHost);
25 }

```

VI. PREGUNTA 6

El siguiente fragmento de código contiene un error, identifíquelo.

```

1  int main(){
2      dim3 bloques (1024, 1024); // 1024 es el tamaño máximo de bloques
3      dim3 threads (1024, 1024); // el problema va aquí
4      ...
5      kernel2D<<<bloques, threads>>>();
6      ...
7  }

```

Existe la limitacion de que por bloque no pueden haber mas de 1024 thread en total; osea que al multiplicar las dimensiones de los threads estas deben entregarnos un numero menor a 1024, en este caso la cantidad de threads por bloque

$$\text{threadsd por bloque} = 1024 * 1024 = 1048576$$

para arreglar este problema hemos de cambiar el numero total de threads

```
1 int main(){
2     dim3 bloques (1024, 1024); // 1024 es el tama o maximo de bloques
3     dim3 threads (32, 32); // el total daria 1024
4     ...
5     kernel2D<<<bloques, threads>>>();
6     ...
7 }
```