# Computación Paralela: CUDA

JHON STEEVEN CABANILLA ALVARADO

MIGUEL CHAVEINTE GARCÍA

G01

# Declaración de bloques e hilos

```
1   //*******************************************************//
2       //                  Blocks & Threads                    //
3       //*******************************************************//
4       // Declaración del numero de bloques y del numero de hilos por bloque
5       int num_blocks;
6       int num_blocks_reduction;
7       int threads_per_block_reduction = 1024;
8       int threads_per_block = 1024;
9
10      // Calculamos el numero de bloques en funcion del tamaño del problema
11      num_blocks = num_searchers / threads_per_block;
12      num_blocks_reduction = (rows * columns) / threads_per_block_reduction;
13
14      // Si el tamaño del vector no es multiplo del tamaño del bloque, lanzamos otro bloque completo (tendremos hilos sobrantes)
15      if (num_searchers % threads_per_block != 0)
16      {
17          num_blocks++;
18      }
19
20      if ((rows * columns) % threads_per_block_reduction != 0)
21      {
22          num_blocks_reduction++;
23      }
```

# Lanzamientos de Kernels

```
1   /* 3.2. Terrain initialization */
2       // Inicializamos las matrices del device
3       computeInitialitationMatrix<<<num_blocks_reduction, threads_per_block_reduction, 0, stream[0]>>>(dev_heights, dev_trails, rows * columns);
4       CHECK_CUDA_LAST();
5
6       /* 3.3. Searchers initialization */
7       searchersInitialization<<<num_blocks, threads_per_block, 0, stream[1]>>>(dev_searchers, dev_total_steps, num_searchers);
8       CHECK_CUDA_LAST();
9
10      /* 4. Compute searchers climbing trails */
11      computeSearchersTrails<<<num_blocks, threads_per_block, 0, stream[2]>>>(num_searchers, rows, columns, dev_searchers, dev_heights, dev_trails, x_min, x_max, y_min, y_max);
12      CHECK_CUDA_LAST();
13
14      /*4.1*/
15      computeFollows<<<num_blocks, threads_per_block, 0, stream[2]>>>(dev_searchers, dev_trails, num_searchers, columns);
16      CHECK_CUDA_LAST();
17
18      /* 5. Compute the leading follower of each searcher */
19      computeFollower<<<num_blocks, threads_per_block, 0, stream[2]>>>(dev_searchers, num_searchers);
20      CHECK_CUDA_LAST();
21
22      /* 6. Compute accumulated trail steps to each maximum */
23      computeAccumulatedTrail<<<num_blocks, threads_per_block, 0, stream[2]>>>(dev_searchers, dev_total_steps, num_searchers);
24      CHECK_CUDA_LAST();
```

# Cómputo estadístico

```
1   // Shared Memory
2       int shared_memory = threads_per_block * sizeof(int);
3       int global_reduction = threads_per_block_reduction * sizeof(int);
4
5
6       // Maximo accum steps
7       cudaMemsetAsync(max_steps, INT_MIN, sizeof(int));
8       cudaMalloc(&max_steps, sizeof(int)); CHECK_CUDA_LAST();
9       reductionMax<<<num_blocks, threads_per_block, shared_memory>>>(dev_total_steps, num_searchers, max_steps); CHECK_CUDA_LAST();
10      cudaMemcpyAsync(&max_accum_steps, max_steps, sizeof(int), cudaMemcpyDeviceToHost); CHECK_CUDA_LAST();
11
12      // Maxima altura ++ Total Heights
13      cudaMalloc(&max_altura, sizeof(int)); CHECK_CUDA_LAST();
14      cudaMalloc(&alturas, sizeof(unsigned long long int)); CHECK_CUDA_LAST();
15      reductionMax<<<num_blocks_reduction, threads_per_block_reduction, global_reduction>>>(dev_heights, rows * columns, max_altura); CHECK_CUDA_LAST();
16      reductionAddHeights<<<num_blocks_reduction, threads_per_block_reduction, global_reduction>>>(dev_heights, rows * columns, alturas); CHECK_CUDA_LAST();
17      cudaMemcpyAsync(&max_height, max_altura, sizeof(int), cudaMemcpyDeviceToHost); CHECK_CUDA_LAST();
18      cudaMemcpyAsync(&total_heights, alturas, sizeof(unsigned long long int), cudaMemcpyDeviceToHost); CHECK_CUDA_LAST();
19
20      // Num local max
21      cudaMalloc(&max_local, sizeof(int)); CHECK_CUDA_LAST();
22      computeLocalMax<<<num_blocks, threads_per_block, shared_memory>>>(dev_searchers, num_searchers, max_local); CHECK_CUDA_LAST();
23      cudaMemcpyAsync(&num_local_max, max_local, sizeof(int), cudaMemcpyDeviceToHost); CHECK_CUDA_LAST();
24
25      // Total tainted
26      cudaMalloc(&tainted_prueba, sizeof(int)); CHECK_CUDA_LAST();
27      reductionAdd<<<num_blocks_reduction, threads_per_block_reduction, global_reduction>>>(dev_trails, rows * columns, tainted_prueba); CHECK_CUDA_LAST();
28      cudaMemcpyAsync(&total_tainted, tainted_prueba, sizeof(int), cudaMemcpyDeviceToHost); CHECK_CUDA_LAST();
29
30      // Sincronizacion final de hilos
31      cudaDeviceSynchronize();
```

# Reducciones

```
1   // Load array values in the shared memory (0 if out of the array)
2       if (globalPos < size)
3       { // and array[globalPos] != INT_MIN
4           if (array[globalPos] != INT_MIN)
5           {
6               buffer[threadIdx.x] = array[globalPos];
7           }
8           else
9           {
10              buffer[threadIdx.x] = 0;
11          }
12      }
13      else
14          buffer[threadIdx.x] = 0;
15
16      // Wait for all the threads of the block to finish
17      __syncthreads();
18
```

```
1   // Reduction tree
2       for (int step = blockDim.x / 2; step >= 1; step /= 2)
3       {
4           if (threadIdx.x < step)
5               if (buffer[threadIdx.x] < buffer[threadIdx.x + step])
6                   buffer[threadIdx.x] = buffer[threadIdx.x + step];
7           __syncthreads();
8       }
```

# Condiciones de carrera

```
1   /* Stop if searcher finds another trail */
2       int check;
3       // check = atomicAdd(&accessMat( tainted, pos_row, pos_col ), 1);
4       check = atomicCAS(&accessMat(trails, pos_row, pos_col), -1, search);
5
6       if (check != -1)
7       {
8           searchers[search].follows = check;
9           search_flag = 1;
10      }
```

# Ejemplo de kernel

```cuda
__global__ void computeFollows(Searcher *searchers, int *trails, int num_searchers, int columns)
{

    // Calculamos la posicion global de cada hilo
    int gid = threadIdx.x + blockDim.x * blockIdx.x;

    // Podemos tener más hilos de los necesarios. Por ello, añadimos este IF para que los hilos ociosos no entren en esta parte del codigo
    if (gid < num_searchers)
    {
        int pos_row = searchers[gid].pos_row;
        int pos_col = searchers[gid].pos_col;

        searchers[gid].follows = accessMat(trails, pos_row, pos_col);
    }
}
```

# Optimizaciones

- Eliminar matriz tainted
- Eliminar el atributo id del struct Searchers
- Declaración de variables para ahorrar llamadas a funciones
- Uso de Streams

# Streams

```
1   /* 3.2. Terrain initialization */
2       // Inicializamos las matrices del device
3       computeInitialitationMatrix<<<num_blocks_reduction, threads_per_block_reduction, 0, stream[0]>>>(dev_heights, dev_trails, rows * columns);
4       CHECK_CUDA_LAST();
5
6       /* 3.3. Searchers initialization */
7       searchersInitialization<<<num_blocks, threads_per_block, 0, stream[1]>>>(dev_searchers, dev_total_steps, num_searchers);
8       CHECK_CUDA_LAST();
9
10      /* 4. Compute searchers climbing trails */
11      computeSearchersTrails<<<num_blocks, threads_per_block, 0, stream[2]>>>(num_searchers, rows, columns, dev_searchers, dev_heights, dev_trails, x_min, x_max, y_min, y_max);
12      CHECK_CUDA_LAST();
13
14      /*4.1*/
15      computeFollows<<<num_blocks, threads_per_block, 0, stream[2]>>>(dev_searchers, dev_trails, num_searchers, columns);
16      CHECK_CUDA_LAST();
17
18      /* 5. Compute the leading follower of each searcher */
19      computeFollower<<<num_blocks, threads_per_block, 0, stream[2]>>>(dev_searchers, num_searchers);
20      CHECK_CUDA_LAST();
21
22      /* 6. Compute accumulated trail steps to each maximum */
23      computeAccumulatedTrail<<<num_blocks, threads_per_block, 0, stream[2]>>>(dev_searchers, dev_total_steps, num_searchers);
24      CHECK_CUDA_LAST();
```

# Pinned Memory

```
1   //*********************************************************//
2      //                  Pinned Memory                        //
3      //*********************************************************//
4      /*
5      cudaMallocHost( (void **)&searchers, sizeof(Searcher) * num_searchers); CHECK_CUDA_LAST();
6      cudaMallocHost( (void **)&total_steps, sizeof(int) * num_searchers); CHECK_CUDA_LAST();
7
8      cudaMallocHost( (void **)&heights, sizeof(int) * (size_t)rows * (size_t)columns); CHECK_CUDA_LAST();
9      cudaMallocHost( (void **)&trails, sizeof(int) * (size_t)rows * (size_t)columns); CHECK_CUDA_LAST();
10     cudaMallocHost( (void **)&tainted, sizeof(int) * (size_t)rows * (size_t)columns); CHECK_CUDA_LAST();
11     */
```

```
1   // Copia de datos del Host al Device
2       //*********************************************************//
3       //                         HostToDevice                    //
4       //*********************************************************//
5
6       /*
7       cudaMemcpyAsync( dev_heights, heights, sizeof(int) * (size_t)rows * (size_t)columns, cudaMemcpyHostToDevice, stream[0]); CHECK_CUDA_LAST();
8       cudaMemcpyAsync( dev_trails, trails, sizeof(int) * (size_t)rows * (size_t)columns, cudaMemcpyHostToDevice, stream[1]); CHECK_CUDA_LAST();
9       cudaMemcpyAsync( dev_tainted, tainted, sizeof(int) * (size_t)rows * (size_t)columns, cudaMemcpyHostToDevice, stream[2]); CHECK_CUDA_LAST();
10      cudaMemcpyAsync( dev_searchers, searchers, sizeof(Searcher) * num_searchers, cudaMemcpyHostToDevice, stream[3]); CHECK_CUDA_LAST();
11      cudaMemcpyAsync( dev_total_steps, total_steps, sizeof(int) * num_searchers, cudaMemcpyHostToDevice, stream[4]); CHECK_CUDA_LAST();
12      */
13
14      //cudaMemcpyAsync(dev_heights, heights, sizeof(int) * (size_t)rows * (size_t)columns, cudaMemcpyHostToDevice);
15      //CHECK_CUDA_LAST();
16      //cudaMemcpyAsync(dev_trails, trails, sizeof(int) * (size_t)rows * (size_t)columns, cudaMemcpyHostToDevice);
17      //CHECK_CUDA_LAST();
18      cudaMemcpyAsync(dev_searchers, searchers, sizeof(Searcher) * num_searchers, cudaMemcpyHostToDevice);
19      CHECK_CUDA_LAST();
20      //cudaMemcpyAsync(dev_total_steps, total_steps, sizeof(int) * num_searchers, cudaMemcpyHostToDevice);
21      //CHECK_CUDA_LAST();
```