

Multiplicación de matrices con CUDA

Rafael Alonso David Peñalva

rafael.david@ucsp.edu.pe

1. Introducción

Se utiliza CUDA C/C++ para implementar dos kernels de multiplicación de matrices:

- Multiplicación de matrices (convencional)
- Multiplicación de matrices + Memoria compartida

Empezaremos detallando el input y el output de los kernels. Más adelante se explicará cada kernel por separado. Finalmente se verá la configuración de lanzamiento.

El código fuente se puede encontrar en: [CUDA-Matrix-Multiplication](#).

2. I/O

2.1. Input

El input para ambos kernels es una matriz generada aleatoriamente, se utiliza la función `generateRandomMatrix`:

```
1 void generateRandomMatrix(vector<float> &matrix, int N)
2 {
3     random_device rd;
4     mt19937 gen(rd());
5     uniform_real_distribution<float> dis(0.0, 9.0);
6
7     for (int i = 0; i < N; i++)
8         for (int j = 0; j < N; j++)
9             matrix[i * N + j] = round(dis(gen));
10 }
```

Listing 1: Implementación de `generateRandomMatrix`

Esta función genera números de coma flotante del 0 al 9 y se redondean. Esto se hizo para verificar las matrices resultantes con más facilidad, sin embargo, puede ser alterado para emitir flotantes no redondeados y en cualquier rango.

2.2. Output

El output de los kernels es una matriz producto de la multiplicación, se puede imprimir añadiendo el argumento `-DEBUG` al ejecutar.

```
1 matrixMulKernel 4 -DEBUG
```

Listing 2: Imprimir las matrices al ejecutar

3. Kernels

3.1. Multiplicación de matrices (convencional)

Este kernel realiza una multiplicación de matrices en paralelo.

```
1 __global__
2 void MatrixMulKernel(float *M, float *N, float *P, int Width)
3 {
4     // Se calculan las coordenadas globales de cada hilo en la matriz P
5     int Row = blockIdx.y * blockDim.y + threadIdx.y;
6     int Col = blockIdx.x * blockDim.x + threadIdx.x;
7     // Se verifica que el hilo este dentro de los limites v lidos de la
8     // matriz P
9     if ((Row < Width) && (Col < Width))
10    {
11        // Pvalue es una variable temporal para almacenar el resultado
12        // parcial de la multiplicacion realizada por cada hilo
13        float Pvalue = 0;
14        // Itera a traves del width de las matrices para realizar la
15        // multiplicacion
16        for (int k = 0; k < Width; ++k)
17        {
18            // Se realiza la multiplicacion de matrices y se acumulan los
19            // resultados
20            Pvalue += M[Row * Width + k] * N[k * Width + Col];
21        }
22        // Se actualiza la matriz de salida P con el resultado parcial
23        // calculado por el hilo
24        P[Row * Width + Col] = Pvalue;
25    }
26 }
```

Listing 3: Implementación de MatrixMulKernel

3.2. Multiplicación de matrices + Memoria compartida

Este kernel realiza una multiplicación de matrices en paralelo utilizando la memoria compartida de cada bloque.

```
1 __global__
2 void MatrixMulKernelShared(float *M, float *N, float *P, int width)
3 {
4     // Se definen dos matrices en memoria compartida del tama o del
5     // bloque (TILE_WIDTH)
6     __shared__ float Mds[TILE_WIDTH][TILE_WIDTH];
7     __shared__ float Nds[TILE_WIDTH][TILE_WIDTH];
8
9     // Obtenemos los ndices de bloque y de hilo
10    int bx = blockIdx.x;
11    int by = blockIdx.y;
12    int tx = threadIdx.x;
13    int ty = threadIdx.y;
14
15    // Se calculan las coordenadas globales de la matriz P
16    // correspondientes al hilo actual
17    int Row = by * TILE_WIDTH + ty;
```

```

16  int Col = bx * TILE_WIDTH + tx;
17
18  // Pvalue es una variable temporal para almacenar el resultado
   // parcial de la multiplicacion realizada por cada hilo
19
20  float Pvalue = 0;
21
22  // Se inicia un bucle para dividir la matriz en bloques mas peque os
23  for (int ph = 0; ph < width / TILE_WIDTH; ph++)
24  {
25      // Se cargan los bloques de las matriz M y N en memoria compartida
   Mds Nds
26      Mds[ty][tx] = M[Row * width + ph * TILE_WIDTH + tx];
27      Nds[ty][tx] = N[(ph * TILE_WIDTH + ty) * width + Col];
28
29      // Se sincronizan los hilos del bloque antes de continuar, caso
   // contrario no habrian datos para multiplicar
30      __syncthreads();
31
32      // Bucle para realizar la multiplicacion de matrices en bloques
   // cargados en memoria compartida
33      for (int k = 0; k < TILE_WIDTH; ++k)
34      {
35          // Se realiza la multiplicaci n de matrices y se acumulan los
   // resultados
36          Pvalue += Mds[ty][k] * Nds[k][tx];
37      }
38      // Se sincronizan los hilos del bloque para asegurar que todos los
   // hilos hayan terminado de multiplicar antes de iterar nuevamente
39      __syncthreads();
40  }
41  // Se actualiza la matriz de salida P con el resultado parcial
   // calculado por el hijo
42  P[Row * width + Col] = Pvalue;
43 }

```

Listing 4: Implementación de MatrixMulKernelShared

3.2.1. Lanzamiento

El lanzamiento debe configurado de esta manera:

```

1  int blockDim = TILE_WIDTH;
2  int gridDim = ceil(width / blockDim);
3
4  dim3 dimGrid(gridDim, gridDim, 1);
5  dim3 dimBlock(blockDim, blockDim, 1);
6
7  MatrixMulKernelShared<<<dimGrid, dimBlock>>>(d_M, d_N, d_P, width);

```

Listing 5: Configuración de lanzamiento

4. Comparación de rendimiento

Se realiza la comparación pasándole a ambas implementaciones n , operando matrices de n . En la gráfica 1 podemos observar en el eje X el tamaño de las matrices, y en el eje Y, los tiempos

en milisegundos.

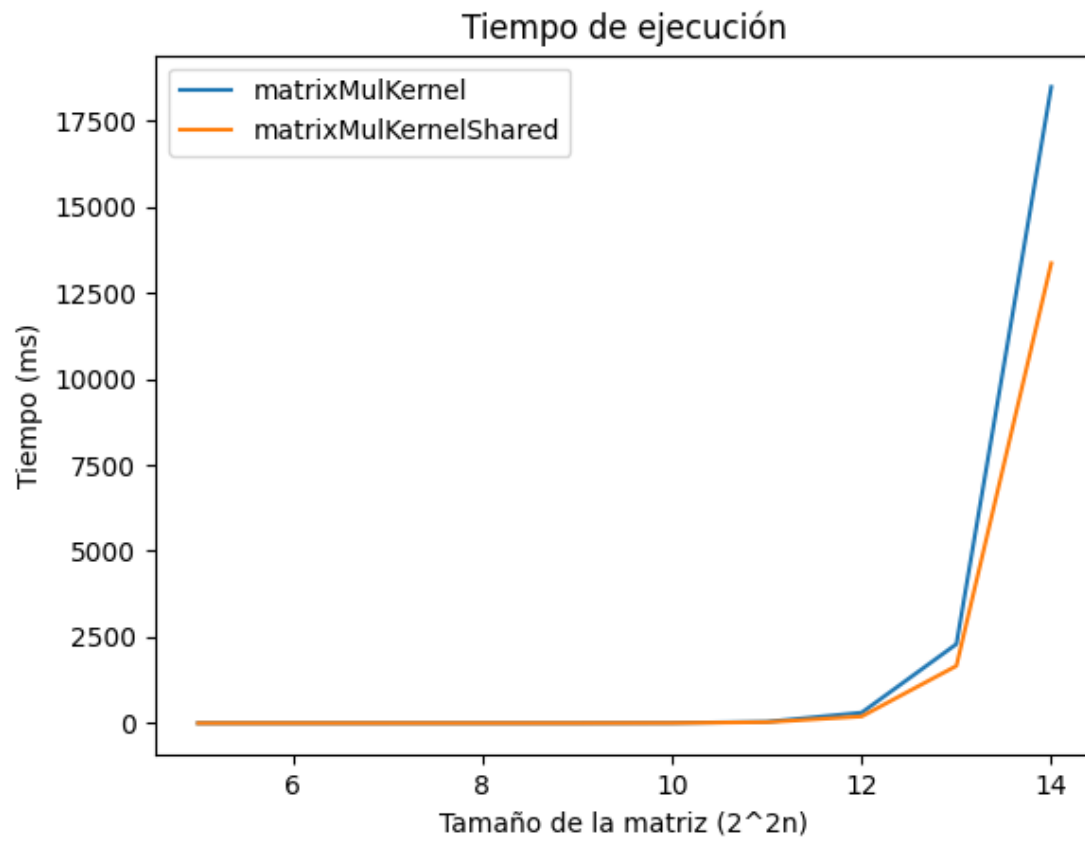


Figura 1: Comparación de tiempos