



UNIVERSIDADE FEDERAL DE SÃO CARLOS
DEPARTAMENTO DE COMPUTAÇÃO

PROGRAMAÇÃO PARALELA E DISTRIBUÍDA

Exercício Programa - EP5

Professor: Hermes Senger

Aluno: Ivan Duarte Calvo, 790739

1 Metodologia

A parte mais importante a ser implementada foi o Kernel CUDA *matrixMulGPU*.

```
1 __global__ void matrixMulGPU( int * a, int * b, int * c )
2 {
3     int row = blockIdx.y*blockDim.y+threadIdx.y;
4     int col = blockIdx.x*blockDim.x+threadIdx.x;
5
6     int val = 0;
7
8     if (row < N && col < N) {
9         for (int i = 0; i < N; i++) {
10             val += a[row * N + i] * b[i * N + col];
11         }
12     }
13     c[row * N + col] = val;
14 }
```

Inicialmente, nas linhas 3 e 4 são calculados os valores de linha(*row*) e coluna(*col*). O cálculo desses valores é importante para transformar uma matriz 2D em um vetor de uma única dimensão. Após isso, é verificado se os valores determinados para *row* e *col* são menores que N, para garantir que a multiplicação seja realizada em posições corretas na matriz. Finalmente é realizado o cálculo principal, na linha 10, e o valor é atribuído para a matriz de resposta.

2 Hardware

O código foi colocado em execução em um ambiente do *Google Colab*, para avaliar qual o hardware disponibilizado pelo ambiente para a realização das tarefas foram utilizados os comandos *nvidia-smi* e *lscpu* para obter informações sobre a GPU e CPU respectivamente.

NVIDIA-SMI 525.85.12										Driver Version: 525.85.12										CUDA Version: 12.0									
GPU		Name		Persistence-M				Bus-Id		Disp.A		Volatile Uncorr. ECC																	
Fan		Temp		Perf		Pwr:Usage/Cap						Memory-Usage		GPU-Util		Compute M.													
																MIG M.													
0		Tesla T4				Off				00000000:00:04.0		Off				0													
N/A		42C		P0		26W / 70W				0MiB / 15360MiB				0%		Default													
																N/A													
Processes:																													
GPU		GI		CI		PID		Type		Process name						GPU Memory Usage													
ID		ID		ID																									
No running processes found																													

Figura 1: *Output* do comando *nvidia-smi*

```

1 Architecture:                x86_64
2 CPU op-mode(s):              32-bit, 64-bit
3 Byte Order:                   Little Endian
4 Address sizes:                46 bits physical, 48 bits virtual
5 CPU(s):                       2
6 On-line CPU(s) list:         0,1
7 Thread(s) per core:          2
8 Core(s) per socket:          1
9 Socket(s):                    1
10 NUMA node(s):                1
11 Vendor ID:                    GenuineIntel
12 CPU family:                   6
13 Model:                        79
14 Model name:                   Intel(R) Xeon(R) CPU @ 2.20GHz
15 Stepping:                     0
16 CPU MHz:                      2199.998
17 BogoMIPS:                     4399.99
18 Hypervisor vendor:            KVM
19 Virtualization type:          full
20 L1d cache:                    32 KiB
21 L1i cache:                    32 KiB
22 L2 cache:                     256 KiB
23 L3 cache:                     55 MiB
24 NUMA node0 CPU(s):           0,1

```

3 Resultados

Para a execução foram utilizadas matrizes quadradas NxN, com o N tendo os tamanhos: 128, 256, 512, 1024 e 2048. A quantidade de *threads* por bloco foi uma matriz 2D de tamanho (16, 16) e o número de blocos foi calculado através da divisão de N pela quantidade de threads por bloco. Os resultados obtidos foram os seguintes:

É importante salientar que nas tabelas e gráficos com resultados são apresentados valores para *speedup*, no entanto, não o *speedup* calculado não é o mesmo usual, pois nesse caso a quantidade de blocos está aumentando de acordo com o tamanho das matrizes multiplicadas. Portanto o calculo foi feito da seguinte forma:

$$Speedup = \frac{T_{cpu}}{T_{gpu}}$$

Tabela 1: Tempos de execução

Tamanho da matriz	Tempo na CPU	Tempo na GPU	<i>Speedup</i>
128	0.006966	0.000565	12.334225
256	0.098065	0.001374	71.361000
512	0.495764	0.002654	186.781904
1024	6.105657	0.013107	465.825727
2048	116.277195	0.086975	1336.910414

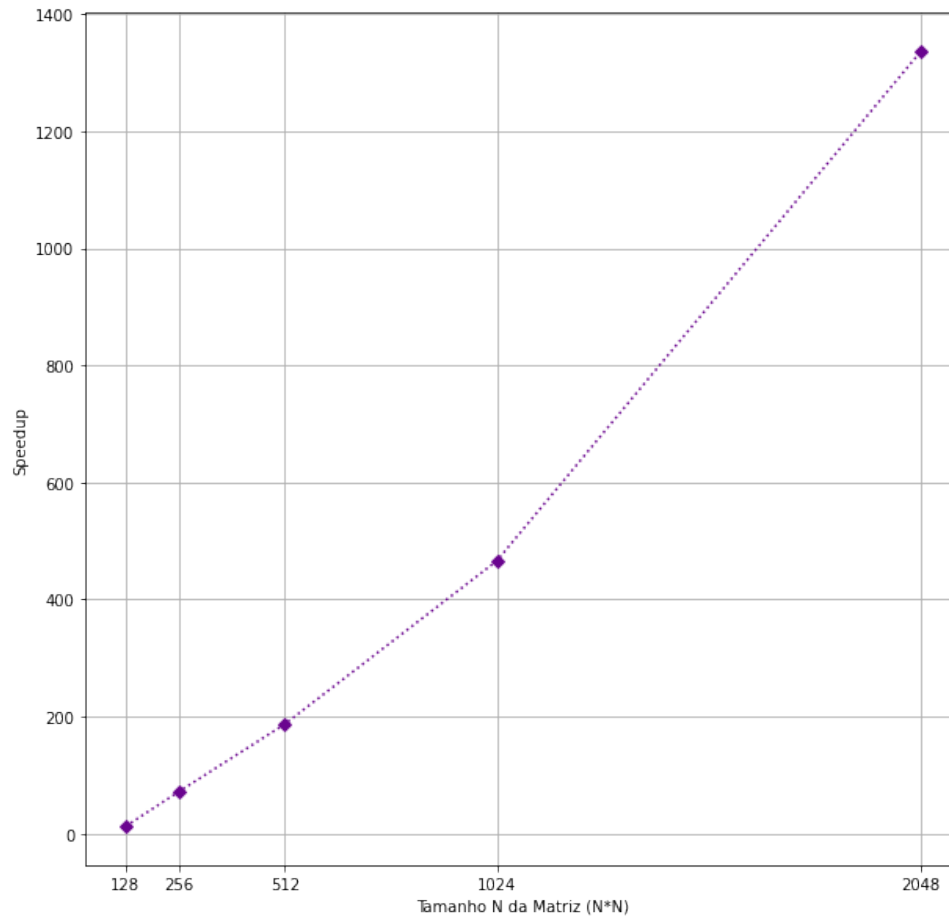


Figura 2: *Speedup* do tempo na GPU em comparação com a CPU

4 Conclusão

Os valores obtidos nos tempos de execução foram extremamente bem sucedidos, indicando a grande vantagem na implementação da paralelização através do CUDA visto que os ganhos nos tempos de execução são muito significativos.

A única possível contraindicação para a utilização seria por conta da alta complexidade para a implementação, contudo, ainda se mostra uma alternativa interessante por conta dos ótimos resultados de maneira que pode se considerar vantajoso superar a barreira inicial do aprendizado necessário.