

Simulação e Análise de Modelos de Difusão de Contaminantes em Água usando CUDA

Gabriel Amaral¹, Lucas Kuroki¹

¹Instituto de Ciência e Tecnologia – Universidade Federal de São Paulo (UNIFESP)
São José dos Campos – SP – Brazil

kenzo.kuroki@unifesp.br

Resumo. Este trabalho visa modelar a difusão de concentração em uma grade bidimensional. Foram desenvolvidas implementações de threads paralelas, utilizando o CUDA e o notebook Collab. Os resultados mostram melhorias consideráveis em relação ao sequencial e ao openMp.

1. Introdução

Para a modelagem de uma difusão de contaminantes em água, foram desenvolvidos códigos no formato sequencial e paralelizado através do OpenMP. Visando utilizar GPUs para calcular a equação de difusão, foi implementado o modelo da difusão em CUDA no ambiente Collab.

1.1. Metodologia

Foi necessário prepara-lo instalando ferramentas do CUDA para sua utilização, os comandos podem ser vistos abaixo.

```
1 // Verificar a presença de uma GPU
2 !nvidia-smi
3
4 // Instalar o compilador CUDA
5 !apt-get install -q nvidia-cuda-toolkit g++ freeglut3-dev libx11
   -dev libxmu-dev libxi-dev libglu1-mesa-dev
```

Listing 1. Comandos para preparar o ambiente

No código que implementa o modelo com CUDA, foram definidos os mesmos parametros de grade das versões anteriores. Em seguida temos a função da difusão que usa uma grade de threads, com cada thread atualizando uma célula da matriz(linha 12 e 13). Uma função auxiliar atomicAddDouble, que realiza somas atômicas, e um kernel computeDifference, que faz os cálculos e armazena em na global dif, nele, cada thread é responsável por calcular a diferença absoluta entre as matrizes(linha 47) e o vetor é usado para acumular as diferenças localmente dentro de cada bloco, realizando um reduction através da thread 0, que acumula os valores da global dif.

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <math.h>
4
5 #define N 2000          // Tamanho da grade
6 #define T 1000          // numero de iteracoes no tempo
```

```

7 #define D 0.1          // Coeficiente de difusao
8 #define DELTA_T 0.01
9 #define DELTA_X 1.0
10
11 __global__ void diff_eq(double *C, double *C_new) {
12     int i = blockIdx.x * blockDim.x + threadIdx.x;
13     int j = blockIdx.y * blockDim.y + threadIdx.y;
14
15     int index = i * N + j;
16
17     if (i > 0 && i < N-1 && j > 0 && j < N-1) {
18         C_new[index] = C[index] + D * DELTA_T * (
19             (C[index + N] + C[index - N] + C[index + 1] + C[
20                 index - 1] - 4 * C[index]) / (DELTA_X * DELTA_X)
21         );
22     }
23 }
24
25 static __device__ double atomicAddDouble(double* address, double
26     val) {
27     unsigned long long int* address_as_ull = (unsigned long long
28         int*)address;
29     unsigned long long int old = *address_as_ull, assumed;
30
31     do {
32         assumed = old;
33         old = atomicCAS(address_as_ull, assumed,
34             __double_as_longlong(val + __longlong_as_double(
35                 assumed)));
36     } while (assumed != old);
37
38     return __longlong_as_double(old);
39 }
40
41 __global__ void compute_difference(double *C, double *C_new,
42     double *dif) {
43     int i = blockIdx.x * blockDim.x + threadIdx.x;
44     int j = blockIdx.y * blockDim.y + threadIdx.y;
45
46     int index = i * N + j;
47
48     __shared__ double dif_per_block[256];
49     int t_idx = threadIdx.x + threadIdx.y * blockDim.x;
50
51     if (i < N && j < N) {
52         dif_per_block[t_idx] = fabs(C_new[index] - C[index]);
53     } else {
54         dif_per_block[t_idx] = 0.0;
55     }
56 }

```

```

50     }
51
52     __syncthreads();
53
54     for (int s = blockDim.x * blockDim.y / 2; s > 0; s >>= 1) {
55         if (t_idx < s) {
56             dif_per_block[t_idx] += dif_per_block[t_idx + s];
57         }
58         __syncthreads();
59     }
60
61     if (t_idx == 0) {
62         atomicAddDouble(dif, dif_per_block[0]);
63     }
64 }

```

Listing 2. código usando CUDA, parte 1

Na função principal, iniciamos as matrizes, a memória e alocamos memória na GPU, por meio do `cudaMalloc`. Após reservar espaço fizemos uma cópia do host para a GPU nas linhas 20 e 21, em seguida configuramos a quantidade de grids e blocos, no caso, $16 \times 16 = 256$ threads e o grid se ajusta conforme a variável global `N`.

No bloco de iterações no tempo, atualizamos a nova matriz de concentração, e a cada 100 iterações zeramos o `ddiff`, chamamos um novo cálculo por meio da função `computeDifference`, pegamos a diferença acumulada para o host e exibimos a diferença média entre iterações. E no final, colocamos a matriz atualizada no host e imprimimos no terminal a concentração no centro.

```

1  int main() {
2      // Aloca o de memoria no host
3      double *C = (double *)malloc(N * N * sizeof(double));
4      double *C_new = (double *)malloc(N * N * sizeof(double));
5      // Inicializa o da memoria
6      for (int i = 0; i < N; i++) {
7          for (int j = 0; j < N; j++) {
8              C[i * N + j] = 0.0;
9              C_new[i * N + j] = 0.0;
10         }
11     }
12     C[(N/2) * N + (N/2)] = 1.0; // Inicializar concentra o
        alta no centro
13     // Aloca o de memoria na GPU
14     double *d_C, *d_C_new, *d_dif;
15     cudaMalloc((void **)&d_C, N * N * sizeof(double));
16     cudaMalloc((void **)&d_C_new, N * N * sizeof(double));
17     cudaMalloc((void **)&d_dif, sizeof(double));
18
19     // Copiar dados do host para a GPU
20     cudaMemcpy(d_C, C, N * N * sizeof(double),
        cudaMemcpyHostToDevice);

```

```

21     cudaMemcpy(d_C_new, C_new, N * N * sizeof(double),
22               cudaMemcpyHostToDevice);
23
24     // Dimensões do bloco e da grade
25     dim3 threadsPerBlock(16, 16);
26     dim3 blocksPerGrid((N + threadsPerBlock.x - 1) /
27                       threadsPerBlock.x, (N + threadsPerBlock.y - 1) /
28                       threadsPerBlock.y);
29
30     // Iterações no tempo
31     for (int t = 0; t < T; t++) {
32         diff_eq<<<blocksPerGrid, threadsPerBlock>>>(d_C, d_C_new
33             );
34         if ((t % 100) == 0) {
35             cudaMemset(d_dif, 0, sizeof(double));
36             compute_difference<<<blocksPerGrid, threadsPerBlock
37                 >>>(d_C, d_C_new, d_dif);
38             double difmedio;
39             cudaMemcpy(&difmedio, d_dif, sizeof(double),
40                       cudaMemcpyDeviceToHost);
41             difmedio /= ((N-2)*(N-2));
42             printf("Iteração %d - Diferença média=%g\n", t,
43                   difmedio);
44         }
45
46         // Trocar os ponteiros
47         double *temp = d_C;
48         d_C = d_C_new;
49         d_C_new = temp;
50     }
51
52     // Copiar resultado final da GPU para o host
53     cudaMemcpy(C, d_C, N * N * sizeof(double),
54               cudaMemcpyDeviceToHost);
55
56     // Exibir resultado final no centro da grade
57     printf("Concentração final no centro: %f\n", C[(N/2) * N +
58           (N/2)]);
59
60     // Liberar memória
61     cudaFree(d_C);
62     cudaFree(d_C_new);
63     cudaFree(d_dif);
64     free(C);
65     free(C_new);
66     return 0;
67 }

```

Listing 3. código usando CUDA, parte 2

2. Resultados e Discussões

Para esta implementação em CUDA, utilizando a ferramenta time nos comandos de execução e obtivemos os seguintes valores de tempo:

```
1 nvcc -o diffusion diffusion.cu -lm
2 time ./diffusion
3 !nvcc -o diffusion diffusion.cu -lm
4 !time ./diffusion
5 Itera o 0 - Diferença m dia=2.00401e-09
6 Itera o 100 - Diferença m dia=1.23248e-09
7 Itera o 200 - Diferença m dia=7.81794e-10
8 Itera o 300 - Diferença m dia=5.11528e-10
9 Itera o 400 - Diferença m dia=4.21632e-10
10 Itera o 500 - Diferença m dia=3.62223e-10
11 Itera o 600 - Diferença m dia=3.05976e-10
12 Itera o 700 - Diferença m dia=2.57135e-10
13 Itera o 800 - Diferença m dia=2.21174e-10
14 Itera o 900 - Diferença m dia=2.00244e-10
15 Concentra o final no centro: 0.095045
16
17 real    0m1.536s
18 user    0m1.073s
19 sys 0m0.297s
```

Listing 4. Comandos e resultados da execução do código

Com o tempo real foram calculados os speedups das implementações antigas e dessa com CUDA, e comparadas na figura abaixo.

Numero de threads	Tempo (s)	Speedup
Serial	67,13	1
OpenMp(8 threads)	10,611	6,326
CUDA	1,536	43,704

Figure 1. Medidas de desempenho usando implementação sequencial, openMp e CUDA

Podemos observar que utilizando o OpenMp com 8 threads, obtivemos um speedup próximo de 6, já utilizando o CUDA, o ganho foi muito maior, com speedup de 43, isso revela a capacidade que GPUs possuem em realizar tarefas atômicas, como de redução e calculo de matrizes, com muita velocidade, sendo assim, uma ferramenta muito poderosa na área de pesquisa que exigem calculos simples mas em grande escala.