

La Programmation GPU

Programmation séquentielle

Traditionnellement, les logiciels sont basés sur le calcul **séquentiel**:

- Un problème est découpé en instructions.
- Ces instructions sont exécutées séquentiellement les unes après les autres.
- Elles sont exécutées par un seul processeur.
- À un instant donné, une seule instruction est exécutée.
- La performance est déterminée principalement par la fréquence (Hz) du processeur.

programmation parallèle

La programmation parallèle permet l'utilisation de plusieurs ressources de calcul pour résoudre un problème donné:

- Un problème est découpé en parties qui peuvent être lancées simultanément.
- Chaque partie est découpée en instructions.
- Les instructions de chaque partie sont exécutées en parallèle en utilisant plusieurs processeurs.
- La performance est déterminée par: La fréquence et le nombre de processeurs, Le degré de parallélisation du problème

Architecture

CPU: "Central Processing Unit" :

Plusieurs unités d'exécutions (cœurs)

Plusieurs niveaux de mémoire (registres, L1, L2, L3, RAM)

Plusieurs ports d'exécution dans chaque cœur (ALUs, unités vectorielles)

Unités vectorielles , Exécution de quelques threads (1-4) simultanément Capable d'exploiter le parallélisme au niveau des instructions (micro-op buffer, renumérotation d'instructions, renommer les registres, ...)

Architecture

GPU: "Graphical Processing Unit", unité de calcul spécifique vectorielle consistant à:

Plusieurs unités d'exécutions (symmetric multiprocesseurs (SM)). Plusieurs niveaux de mémoire . Plusieurs unités vectorielles (2-4) larges (16-32 flottants) dans chaque SM. Exécution de milliers de threads simultanément. Grand tableau de registres (65K). Échange de threads très rapide. La plupart du circuit est consacrée aux unités vectorielles. Parallélisation prend l'effort

Programmation GPU

La prog GPU est adaptée au modèle d'exécution (single instruction multiple thread) (SIMT)

Un GPU consiste à plusieurs processeurs appelés "streaming multiprocessors" (SM).

`--global--` précise la définition d'un kernel GPU (sinon fonction CPU par défaut)

blockIdx.x est prédefini et donne l'identifiant d'un bloc dans un kernel GPU.

gridDim.x est prédefini et donne le nombre de blocs utilisés dans le kernel GPU en cours d'exécution.

Compilation: `nvcc programme.cu -o programme`

Exécution: `./programme`

Chaque appel de kernel est non-bloquant , mais on peut le rendre bloquant si on veut.

Threads identiques (exécutent le même code)

Threads organisés en blocs (de taille 32-1024)

Chaque bloc identique s'exécute sur un SM.

Blocs organisés en grilles et répartis sur tous les SMs

Note

Applications du calcul parallèle: Traitement des réseaux neurones. Graphiques (rendering, jeux vidéo, etc.). Simulations en physique.

Le transfert entre CPU et GPU de données se fait sur le bus PCI express (32Go/s débit chaque direction pour PCIe4).

Supercalculateur / Cluster : Un ensemble de machines (CPU+GPU) connecté. Connexion par un réseau avec une topologie particulière (anneau, grille, torus, clique, etc.)

omp.h: est la bibliothèque OpenMP qui fournit les fonctions nécessaires (e.g., pour obtenir `thid`, `numth`)

code

Listing 1: Hello World en OpenMP

```
1 #include <stdio>
2 #include "omp.h"
3 int main(int argc, char **argv)
4 {
5     #pragma omp parallel num_threads(3)
6     {
7         int thid = omp_get_thread_num();
8         int numth = omp_get_num_threads();
9         printf("Hello from thread %d / %d\n",
10              thid, OpenMP_numth);
11     }
12     return 0;
13 }
```

Listing 2: Hello World en CUDA

```
1 #include <stdio>
2 #include "cuda.h"
3
4 __global__ void helloWorld()
5 {
6     printf("Hello from block %d/%d\n",
7          blockIdx.x, gridDim.x);
8 }
9
10 int main(int argc, char **argv)
11 {
12     helloWorld<<<3, 1>>>>();
13     cudaDeviceSynchronize();
14     return 0;
15 }
```

Listing 3: Multiplier un tableau en CUDA

```

1 #include <stdio>
2 #include "cuda.h"
3
4 #define N 1024
5 float A[N];
6 float c = 2.0;
7
8 __device__ float dA[N];
9
10 __global__ void multiplyArray(int n,
11                               float c)
12 {
13     int elemParBlock = n / gridDim.x;
14     int begin = blockIdx.x * elemParBlock;
15     int end;
16     if (blockIdx.x < gridDim.x - 1) {
17         end = (blockIdx.x + 1) * elemParBlock;
18     } else {
19         end = n;
20     }
21     for (int i = begin; i < end; i++) { dA[i] *= c; }
22 }
23
24 int main(int argc, char **argv)
25 {
26     // Initialisation
27     for (int i = 0; i < N; i++) { A[i] = i; }
28     // Copier le tableau vers le GPU
29     cudaMemcpyToSymbol(dA, A, N * sizeof(
30         float), 0,
31         cudaMemcpyHostToDevice);
32     multiplyArray<<<4, 1>>>>(N, c);
33     // Recopier le tableau multiplie vers
34     // le CPU
35     cudaMemcpyFromSymbol(A, dA, N * sizeof(
36         float), 0,
37         cudaMemcpyDeviceToHost);
38     printf("%lf\n", A[2]);
39     return 0;
40 }

```

Listing 4: Multiplier un tableau en CUDA

```

1 #include <stdio>
2 #include <iostream>
3 #include "cuda.h"
4 #define N 513
5 #define BSXY 32
6 float A[N][N], B[N][N], C[N][N];
7 __device__ float dA[N][N], dB[N][N], dC[N][N];
8 __global__ void multiplyMatrixGPUByBlocks(
9     int n)
10 {
11     int i = blockIdx.x; int j = blockIdx.y;
12     float c = 0.0;
13     for (int k = 0; k < n; k++) { c += dA[i][k] * dB[k][j]; }
14     dC[i][j] = c; }
15 //calculer avec blockDim.x threads par bloc.
16 __global__ void multiplyMatrixGPUByBlocksThreads1D(
17     int n)
18 {
19     int i = blockIdx.x;
20     int j = threadIdx.x + blockIdx.y * blockDim.x;
21     float c = 0.0;
22     for (int k = 0; k < n; k++) { c += dA[i][k] * dB[k][j]; }
23     dC[i][j] = c; }
24 // Creer un bloc pour le calcul de
25 // blockDim.x elements de C, calculer avec blockDim.x threads par bloc.
26 __global__ void multiplyMatrixGPUByBlocksThreads1DNonMultiple(
27     int n)
28 {
29     int i = blockIdx.x; element C[ i ][ j ]
30     , c a l c u l e r avec 1
31     int j = threadIdx.x + blockIdx.y * blockDim.x;
32     if (j < n) {
33         float c = 0.0;
34         for (int k = 0; k < n; k++) { c += dA[i][k] * dB[k][j]; }
35         dC[i][j] = c; } }
36 // Creer un bloc pour le calcul de
37 // blockDim.x * blockDim.y elements de C , calculer avec blockDim.x * blockDim.y threads par bloc.
38 __global__ void multiplyMatrixGPUByBlocksThreads2D(
39     int n)
40 {
41     int i = threadIdx.y + blockIdx.x * blockDim.y;
42     int j = threadIdx.x + blockIdx.y * blockDim.x;
43     float c = 0.0;
44     for (int k = 0; k < n; k++) { c += dA[i][k] * dB[k][j]; }
45     dC[i][j] = c; }

```

Listing 5: Multiplier un tableau en CUDA

```

1 __global__ void multiplyMatrixGPU(
2     ByBlocksThreads2D(int n)
3 {
4     int i = threadIdx.y + blockIdx.x * blockDim.y;
5     int j = threadIdx.x + blockIdx.y * blockDim.x;
6     float c = 0.0;
7     for (int k = 0; k < n; k++) { c += dA[i][k] * dB[k][j]; }
8     dC[i][j] = c; }
9 __global__ void multiplyMatrixGPUByBlocksThreads2DNonMultiple(
10     int n)
11 {
12     int i = threadIdx.y + blockIdx.x * blockDim.y;
13     int j = threadIdx.x + blockIdx.y * blockDim.x;
14     if (i < n && j < n) { float c = 0.0;
15         for (int k = 0; k < n; k++) { c += dA[i][k] * dB[k][j]; }
16         dC[i][j] = c; } }
17 __global__ void multiplyMatrixGPUByBlocksThreads2DNonMultipleSharedMemory(
18     int n)
19 {
20     int row = threadIdx.y + blockIdx.x * blockDim.y;
21     int col = threadIdx.x + blockIdx.y * blockDim.x;
22     __shared__ float shA[BSXY][BSXY];
23     __shared__ float shB[BSXY][BSXY];
24     float c = 0.0;
25     const int nsteps = (n - 1) / BSXY + 1;
26     for (int step = 0; step < nsteps; step++) {
27         int offset = step * BSXY;
28         int stepRowB = offset + threadIdx.y;
29         int stepColA = offset + threadIdx.x;
30         if (row < n && stepColA < n) {
31             shA[threadIdx.y][threadIdx.x] = dA[stepRowB][stepColA];
32         } else {
33             shA[threadIdx.y][threadIdx.x] = 0.0; }
34         if (stepRowB < n && col < n) {
35             shB[threadIdx.y][threadIdx.x] = dB[stepRowB][col];
36         } else {
37             shB[threadIdx.y][threadIdx.x] = 0.0; }
38         __syncthreads();
39         for (int k = 0; k < BSXY; k++) {
40             c += shA[threadIdx.y][k] * shB[k][threadIdx.x]; }
41         __syncthreads();
42         if (row < n && col < n) { dC[row][col] = c; } }

```

Listing 6: Multiplier un tableau en CUDA

```

1 void multiplyMatrixCPU()
2 {for (int i = 0; i < N; i++) {
3     for (int j = 0; j < N; j++) {
4         C[i][j] = 0.0f;
5         for (int k = 0; k < N; k++) {
6             C[i][j] += A[i][k] * B[k][j];}}}}
7 void verifyResults()
8 {for (int i = 0; i < N; i++) {
9     for (int j = 0; j < N; j++) {
10        float c = 0.0f;
11        for (int k = 0; k < N; k++) {
12            c += A[i][k] * B[k][j];}
13        if (std::abs(C[i][j] - c) > 1e-6) {
14            std::cout << "Multiplication is
                incorrect for the element C["
                << i << "][" << j << "]" <<
                std::endl;
15            return;}}}}
16 std::cout << "Multiplication is correct
    !" << std::endl;
17 }
18 int main(int argc, char **argv)
19 {for (int i = 0; i < N; i++) {
20     for (int j = 0; j < N; j++) {
21         A[i][j] = i + j;
22         B[i][j] = i - j;}}
23     cudaMemcpyToSymbol(dA, A, N * N *
        sizeof(float), 0,
24         cudaMemcpyHostToDevice);
25     cudaMemcpyToSymbol(dB, B, N * N *
        sizeof(float), 0,
26         cudaMemcpyHostToDevice);
27     dim3 dimGrid;
28     dimGrid.x = (N - 1) / 32 + 1;
29     dimGrid.y = (N - 1) / 32 + 1;
30     dimGrid.z = 1;dim3 dimBlock;
31     dimBlock.x = 32;dimBlock.y = 32;
32     dimBlock.z = 1;
33     // multiplyMatrixGPUBlocks<<<dimGrid,
        dimBlock>>>(N);
34     // multiplyMatrixGPUBlocksThreads1D
        <<<dimGrid, dimBlock>>>(N);
35     // multiplyMatrixGPUBlocksThreads2D
        <<<dimGrid, dimBlock>>>(N);
36     // Recopier le tableau dC vers le CPU
37     cudaMemcpyFromSymbol(C, dC, N * N *
        sizeof(float), 0,
38         cudaMemcpyDeviceToHost);
39     // multiplyMatrixCPU();
40     verifyResults();
41     return 0;
42 }

```