

# Introduction à CUDA C

Amina Guermouche

Télécom SudParis









# Interrogation du GPU

- Quel est le nombre de GPUs disponibles
- Quelle est la taille de la mémoire disponible?
- Quelles sont les caractéristiques des GPUs?

```
    cudaDeviceProp prop;  
    int count;  
    cudaGetDeviceCount(&count);  
  
    for (int i = 0; i < count; i ++)  
    {  
        cudaGetDeviceProperties(&prop, i);  
        printf( "' Taille total de la memoire globale  
               %ld\n'", prop.totalGlobalMem);  
    }
```

- On peut même choisir le GPU qu'on veut selon des critères!!!!

```
    cudaChooseDevice(&dev, &prop)
```



# Hello, World!

```
int main (void) {  
  
    printf ( "Hello , World!\n" );  
    return 0;  
}
```

- Compilation avec nvcc (compilateur NVIDIA)
- nvcc ne se plaint pas s'il n'y a pas de code pour le device





# Hello, World!

## Code du device

```
__global__ void kernel (void)
```

- `__global__` indique que :
  - Le code s'exécute sur le device
  - Le code est appelé du host
- La partie device et interface est gérée par le compilateur nvidia
- La partie host par le compilateur C
- La syntaxe est obligatoire
- `__global__` ne retourne pas de valeur, JAMAIS

ok on utilise le GPU pour appeler la fonction kernel qui ne fait rien, super!!!!

























# Programmation parallèle en CUDA

- Chaque appel parallèle à `add(...)` est appelé **block**
- L'accès à un block donné se fait via `blockIdx.x`
- Chaque `blockIdx.x` référence un élément du tableau

```
__global__ void add (int *a, int *b, int *c){  
    c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];  
}
```

## Programmation parallèle en CUDA

- Chaque appel parallèle à `add(...)` est appelé **block**
- L'accès à un block donné se fait via `blockIdx.x`
- Chaque `blockIdx.x` référence un élément du tableau

```
__global__ void add (int *a, int *b, int *c){
    c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];
}
```

Block<sub>0</sub>

c[0] = a[0] + b[0]

Block<sub>1</sub>

c[1] = a[1] + b[1]

Block<sub>2</sub>

c[2] = a[2] + b[2]

Block<sub>3</sub>

c[3] = a[3] + b[3]





## Programmation parallèle en CUDA : le main

```
// Copie des données vers le Device
cudaMemcpy (gpu_a, a, size, cudaMemcpyHostToDevice);
cudaMemcpy (gpu_b, b, size, cudaMemcpyHostToDevice);

add <<< N, 1 >>> (gpu_a, gpu_b, gpu_c);

//Copie du resultat
cudaMemcpy(c, gpu_c, size, cudaMemcpyDeviceToHost);

free(a); free(b); free(c);
cudaFree (gpu_a);
cudaFree (gpu_b);
cudaFree (gpu_c);
return 0;
}
```

Et si on avait un vecteur à 2 dimensions (une matrice donc)?

- Le nombre de blocks lancés représentent une grille (*grid*)
- Le nombre de blocks par dimension est limité (`maxGridSize[3]`)
- `blockIdx.x`, `blockIdx.y`, `blockIdx.z`
- `dim3 grid(DIM, DIM)` initialise la variable `grid` de type `dim3` qui indique la dimension de la grille (2D)
- `gridDim.x`, `gridDim.y` donnent la dimension de la grille



## Addition de deux matrices

```
#define N 512 //taille d'une dimension de la
              matrice

__global__ void add (int *a, int *b, int *c){
    int x = blockIdx.x;
    int y = blockIdx.y;
    int indice = x + y * gridDim.x;
    c[indice] = a[indice] + b[indice];
}

int main (void){
    int *a, *b, *c;
    int *gpu_a, *gpu_b, *gpu_c;
    int size = N * sizeof(int);
    ...
    dim3 grid(N, N);
    add <<<grid, 1 >>> (dev_a, dev_b, dev_c);
    ...
}
```















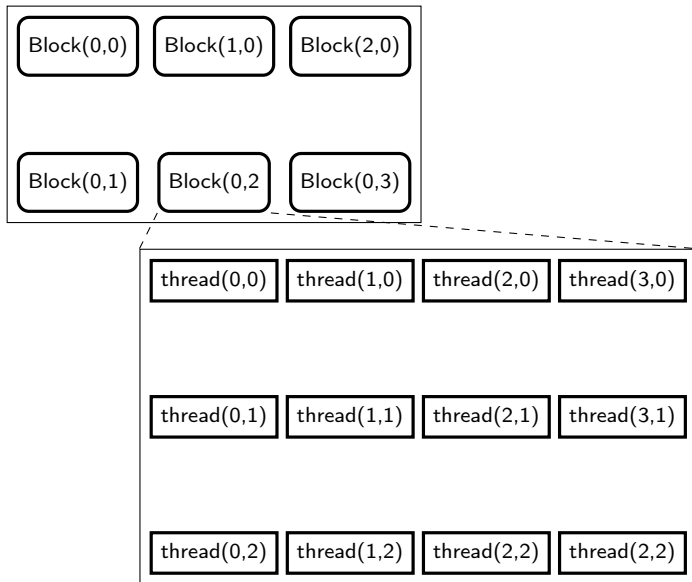




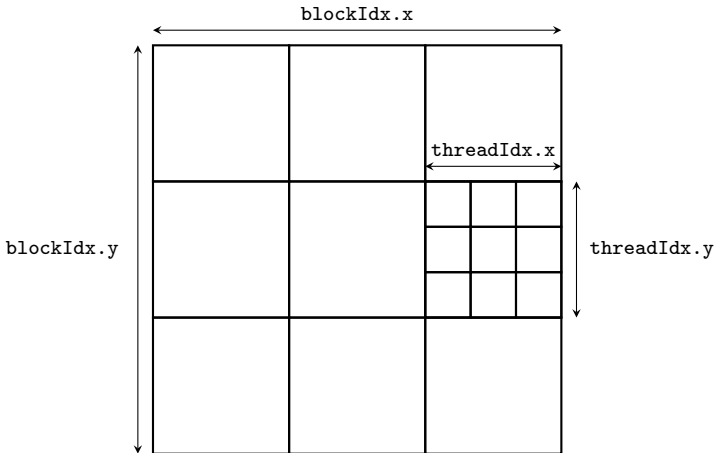




## Des blocks et des threads

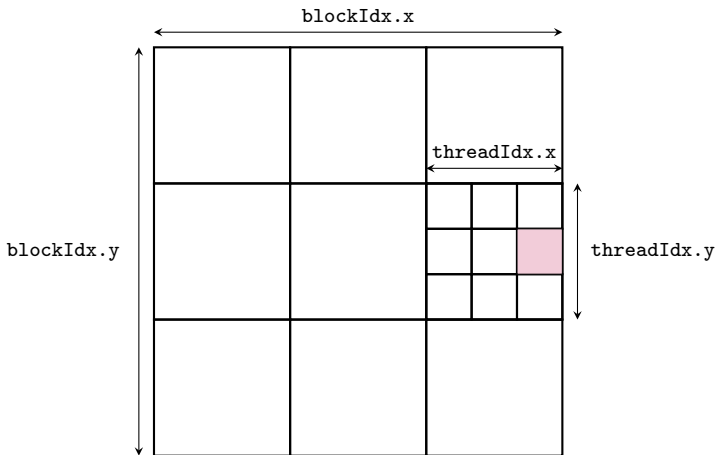


# Des blocks et des threads

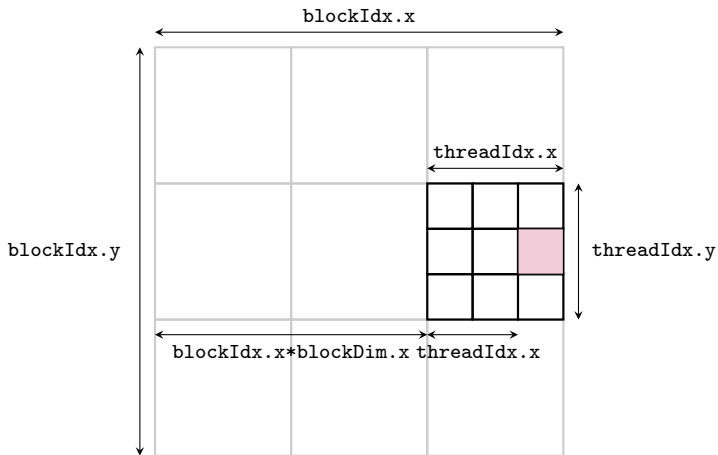




# Des blocks et des threads



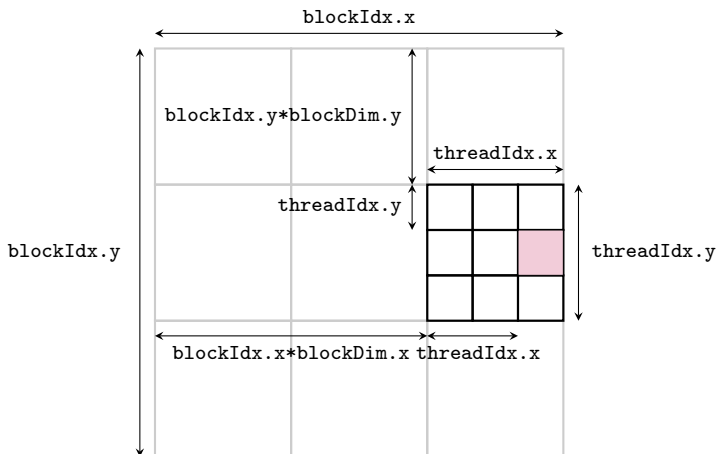
# Des blocks et des threads





## Addition de 2 matrices

- $\text{colonne} = \text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x}$
- $\text{ligne} = \text{blockIdx.y} * \text{blockDim.y} + \text{threadIdx.y}$



## Addition de deux matrices

```

#define N 2048 //taille de la matrice
#define THREAD_PER_BLOCK 512 //nombre de threads

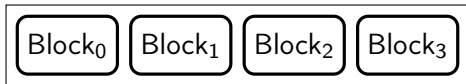
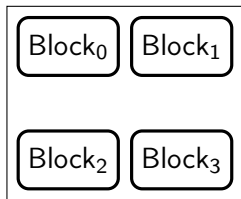
__global__ void add (int *a, int *b, int *c){
    int colonne = blockIdx.x*blockDim.x+threadIdx.x;
    int ligne = blockIdx.y*blockDim.y+threadIdx.y;
    int indice = ligne * N + colonne;
    // N = blockDim.x * blockDim.x;
    c[indice] = a[indice] + b[indice];
}

int main (int argc, char** argv)
{
    ...
    dim3 block(THREAD_PER_BLOCK, THREAD_PER_BLOCK);
    dim3 grid(N/blockDim.x, N/blockDim.y);
    add <<<grid, block >>> (dev_a, dev_b, dev_c);
    ...
}

```

## Pourquoi s'embêter avec des threads et des blocks ?

- 1 thread = 1 coeur
- 1 block = 1SM
- 1 block est exécuté sur 1SM
- Blocks :
  - Les blocks sont exécutés dans n'importe quel ordre, séquentiellement ou en parallèle
  - L'avantage est que ça scale automatiquement avec le nombre de SM





## Les paramètres du kernel

- Les blocks :
  - Le nombre de blocks doit être supérieur au nombre de SM (pour que tous travaillent)
  - Il devrait y avoir plusieurs blocks par SM, afin que d'autres blocks s'exécutent pendant une synchronisation
  - Si une synchronisation est utilisée, il vaut mieux utiliser plusieurs petits blocks qu'un grand
- Les threads :
  - Les SM ordonnancent les threads par groupe SIMD de 32 (warp) sur Quadro 620
    - Les threads d'un warp sont synchronisés
  - Les threads dans un block sont exécutés par groupe de 32
    - Un SM peut exécuter plusieurs blocks de manière concurrente



# Les paramètres du kernel

- Les blocks :
  - Le nombre de blocks doit être supérieur au nombre de SM (pour que tous travaillent)
  - Il devrait y avoir plusieurs blocks par SM, afin que d'autres blocks s'exécutent pendant une synchronisation

→ Si une synchronisation est utilisée, il vaut mieux utiliser plusieurs petits blocks qu'un grand
- Les threads :
  - Les SM ordonnancent les threads par groupe SIMD de 32 (warp) sur Quadro 620
    - Les threads d'un warp sont synchronisés

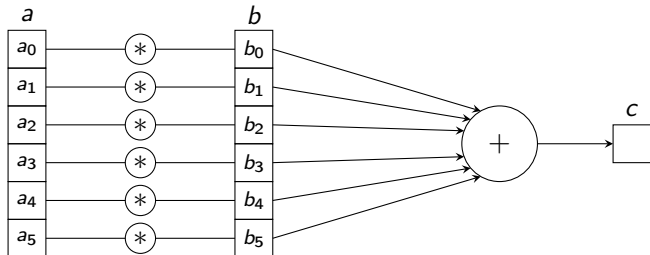
→ Les threads dans un block sont exécutés par groupe de 32

  - Un SM peut exécuter plusieurs blocks de manière concurrente
  - Utiliser des blocks de taille multiple de la taille du warp



# Execice 4

# Produit scalaire (dot product)



$$c = (a_0, a_1, a_2, a_3, a_4, a_5) \cdot (b_0, b_1, b_2, b_3, b_4, b_5)$$
$$= a_0b_0 + a_1b_1 + a_2b_2 + a_3b_3 + a_4b_4 + a_5b_5$$

## Produit scalaire (dot product) : 1 seul block

```
__global__ void dot (int *a, int *b, int *c){
    // chaque thread calcule le produit d'une paire
    int temp = a[threadIdx.x] * b[threadIdx.x];
}
```

## Produit scalaire (dot product) : 1 seul block

```
__global__ void dot (int *a, int *b, int *c){  
    // chaque thread calcule le produit d'une paire  
    int temp = a[threadIdx.x] * b[threadIdx.x];  
}
```

- Le calcul est local au processus
- Les variables temp ne sont pas accessibles aux autres processus
- Mais il faut partager les données pour faire la somme finale

## Partager les données entre les threads (d'un même block)

- Les threads d'un block partagent une zone mémoire appelée *shared memory*
- Caractéristiques
  - Extrêmement rapide
  - *on-chip*
  - Déclarée avec `__shared__`
- Des blocks sur le même SM partagent la même mémoire partagée globale. Donc pour une mémoire de 48KB avec N blocks sur le même SM, chaque block possédera  $48/N$  de mémoire partagée

## Produit scalaire (dot product) : 1 seul block

```

#define N 512 //taille du tableau
__global__ void dot (int *a, int *b, int *c){
    __shared__ int tmp[N]
    // chaque thread calcule le produit d'une paire
    int temp[threadIdx.x] = a[threadIdx.x] * b[
        threadIdx.x];

    //Le thread 0 effectue la somme
    if (0 == threadIdx.x){
        int sum = 0;
        for (int i = 0; i < N; i++)
            sum = sum + temp[i];
        *c = sum;
    }
}

```



## Synchronisation des threads d'un même block

- Grâce à la fonction `__syncthreads()`
- Synchronise uniquement les threads d'un même block





## Programmation parallèle en CUDA : le main

```

// Copie des donnees vers le Device
cudaMemcpy (gpu_a, a, size, cudaMemcpyHostToDevice);
cudaMemcpy (gpu_b, b, size, cudaMemcpyHostToDevice);

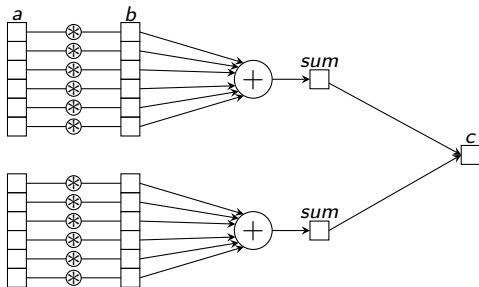
//Lancement de l'operation avec N threads et un
    seul block
dot <<< 1, N >>> (gpu_a, gpu_b, gpu_c);

//Copie du resultat
cudaMemcpy(&c, gpu_c, sizeof(int),
           cudaMemcpyDeviceToHost);

free(a); free(b);
cudaFree(gpu_a);
cudaFree(gpu_b);
cudaFree(gpu_c);
return 0;
}

```

# Plus de parallélisme : plusieurs blocks



## Produit scalaire (dot product)

```

#define N (2048 * 2048)
#define THREAD_PER_BLOCK 512 //taille du tableau
__global__ void dot (int *a, int *b, int *c){
    __shared__ int tmp[THREADS_PER_BLOCK]
    // chaque thread calcule le produit d'une paire
    int indice = threadIdx.x + blockIdx.x * blockDim.
        x;
    temp[indice] = a[indice] * b[indice];

    //Synchronisation (dans le block)
    __syncthreads();

    //Le thread 0 effectue la somme
    if (0 == threadIdx.x){
        int sum = 0;
        for (int i = 0; i < THREADS_PER_BLOCK; i++)
            sum = sum + temp[i];
        *c = sum;
    }
}

```

# Race condition

- c est dans la mémoire globale
- plusieurs thread 0 peuvent y accéder en même temps

## Race condition

- c est dans la mémoire globale
- plusieurs thread 0 peuvent y accéder en même temps
- **Les opérations atomiques :**
  - Les opération de lecture, modification et écriture sont ininterrompibles
  - Plusieurs opérations atomiques possibles avec CUDA :

- `atomicAdd()`

- `atomicSub()`

- `atomicMin()`

- `atomicMax()`

- `atomicInc()`

- `atomicDec()`

- `atomicExch()`

- `atomicCAS()`



## Produit scalaire (dot product)

```

#define N (2048 * 2048)
#define THREAD_PER_BLOCK 512 //taille du tableau
__global__ void dot (int *a, int *b, int *c){
    __shared__ int tmp[THREADS_PER_BLOCK]
    // chaque thread calcule le produit d'une paire
    int indice = threadIdx.x + blockIdx.x * blockDim.
        x;
    temp[indice] = a[indice] * b[indice];

    //Synchronisation (dans le block)
    __syncthreads();

    //Le thread 0 effectue la somme
    if (0 == threadIdx.x){
        int sum = 0;
        for (int i = 0; i < THREADS_PER_BLOCK; i++)
            sum = sum + temp[i];
        atomicAdd(*c, sum);
    }
}

```































































































