

Introduction à CUDA C

Amina Guermouche

CUDA C

- Basé sur le standard C
- Extension du langage pour la programmation hétérogène
- API pour gérer les GPU, la mémoire, ...
- 3 niveaux d'abstraction : groupe de threads, mémoire partagée et une barrière de synchronisation

Fonctionnement

CPU (Host)



```
int main(int
argc, char
**argv)
```

GPU (**Device**)

```
__device__
ma_fonction_gpu
()
```

Fonctionnement

CPU (Host)



```
int main(int
argc, char
**argv)
```

GPU (**Device**)



```
__global__
ma_fonction_interface
()
```

```
__device__
ma_fonction_gpu
()
```

- 1 Host : Copier les données d'entrée de la mémoire du CPU à la mémoire du GPU
- 2 Device : Charger les instructions sur le GPU
- 3 Device : Copier les données vers la mémoire du CPU

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

Et si on faisait faire quelque chose au GPU (1/2)

```
__global__ void add (int *a, int *b, int *c){
    *c = *a + *b;
}

int main ( void ){
    int a, b, c;
    int *gpu_a, *gpu_b, *gpu_c;
    int size = sizeof(int);
```


Et si on faisait faire quelque chose au GPU (2/2)

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

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

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

return 0
}
```


Et si on faisait faire quelque chose au GPU (2/2)

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

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

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

//Liberation de l'espace alloue
cudaFree(gpu_a);
cudaFree( gpu_b);
cudaFree ( gpu_c);

return 0
}
```

Et si on faisait faire quelque chose au GPU (2/2)

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

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

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

//Liberation de l'espace alloue
cudaFree(gpu_a);
cudaFree( gpu_b);
cudaFree ( gpu_c);

return 0
}
```

addition de 2 entiers : super utilisation du parallélisme

- Comment exécuter le code en parallèle ?

On veut faire N fois add en parallèle

```
add<<<1, 1>>> (gpu_a, gpu_b, gpu_c)
```

addition de 2 entiers : super utilisation du parallélisme

- Comment exécuter le code en parallèle ?

On veut faire N fois add en parallèle

```
add<<<1, 1>>> (gpu_a, gpu_b, gpu_c)
```

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

addition de 2 entiers : super utilisation du parallélisme

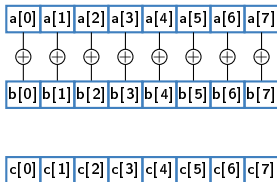
- Comment exécuter le code en parallèle ?

On veut faire N fois add en parallèle

`add<<<1, 1>>> (gpu_a, gpu_b, gpu_c)`

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

- Dans ce cas, autant faire un add sur un vecteur



- Comment sont exprimés les indices sur le GPU ?

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 : le main

```
#define N 512 //nombre d'elements du tableau
int main (void){
    int *a, *b, *c;
    int *gpu_a, *gpu_b, *gpu_c;
    int size = N * sizeof(int);
    // allocation de l'espace pour le device}
    cudaMalloc((void **)&gpu_a, size);
    cudaMalloc((void **)&gpu_b, size);
    cudaMalloc((void **)&gpu_c, size);

    a=(int*) malloc (size);
    b=(int*) malloc (size);
    random_ints(a, N);
    random_ints(b, N);
```


Programmation parallèle en CUDA : le main

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

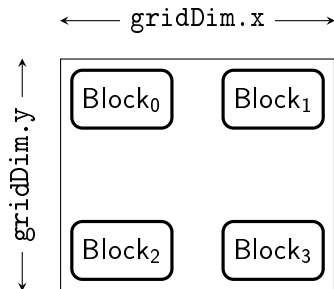
add <<< N, 1 >>> (gpu_a, gpu_b, gpu_c);

// 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;
}
```

- 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

Grille et blocks



- Un block peut être divisé en plusieurs threads parallèles
- CUDA définit un unique id par thread `threadIdx.x`
- On utilise `threadIdx.x` au lieu de `blockIdx.x`

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

Programmation parallèle en CUDA : le main

```
#define N 512 //nombre d'elements du tableau
int main (void){
    int *a, *b, *c;
    int *gpu_a, *gpu_b, *gpu_c;
    int size = N * sizeof(int);
    // allocation de l'espace pour le device}
    cudaMalloc((void **)&gpu_a, size);
    cudaMalloc((void **)&gpu_b, size);
    cudaMalloc((void **)&gpu_c, size);

    a=(int*) malloc (size);
    b=(int*) malloc (size);
    random_ints(a, N);
    random_ints(b, N);
```

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

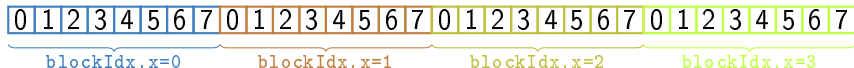
//Lancement de l'operation avec N threads
add <<< 1, N >>> (gpu_a, gpu_b, gpu_c);

//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;
}
```


Des blocks et des threads

- Les threads sont numérotés de 0 à `nb_thread` par block
- `threadIdx.x` est par block



- Si on a M threads/block, l'indice dans un vecteur est calculé par :
$$\text{indice} = \text{threadIdx.x} + \text{blockIdx.x} * M$$
- Le nombre de threads par block est donné par la variable `blockDim.x`
$$\text{indice} = \text{threadIdx.x} + \text{blockIdx.x} * \text{blockDim.x}$$

add avec blocks et threads

```
#define N (2048 * 2048) //taille du tableau
#define THREAD_PER_BLOCK 512 //nombre de threads

__global__ void add (int *a, int *b, int *c){
    int indice = threadIdx.x + blockIdx.x * blockDim.
        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);
    // allocation de l'espace pour le device}
    cudaMalloc((void **)&gpu_a, size);
    cudaMalloc((void **)&gpu_b, size);
    cudaMalloc((void **)&gpu_c, size);
```

Programmation parallèle en CUDA : le main

```
a=(int*) malloc (size);
b=(int*) malloc (size);
random_ints(a, N);
random_ints(b, N);

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

//Lancement de l'operation avec THREAD_PER_BLOCK
par block
add <<< N/THREAD_PER_BLOCK,THREAD_PER_BLOCK >>>
    (gpu_a, gpu_b, gpu_c);

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

Programmation parallèle en CUDA : le main

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

Exercise 3

Vecteurs de taille quelconque

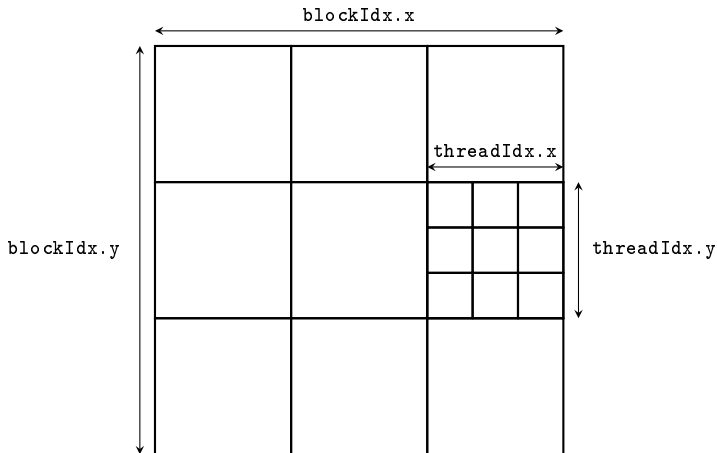
- Vecteur de taille non multiple de `blockDim.x`

```
__global__ void add (int *a, int *b, int *c,
    int n){
    int indice = threadIdx.x + blockIdx.x *
        blockDim.x;
    if (indice < n)
        c[indice] = a[indice] + b[indice];
}
```

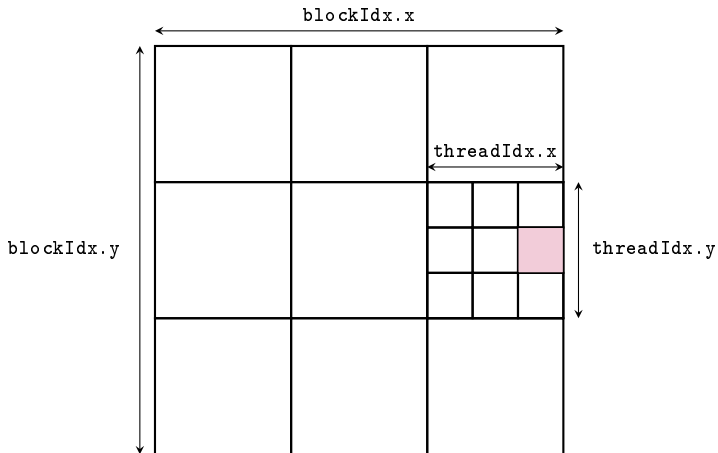
- Au niveau du main

```
add <<< (N+THREAD_PER_BLOCK-1)/  
    THREAD_PER_BLOCK, THREAD_PER_BLOCK >>> (  
    gpu a, gpu b, gpu c, N);
```

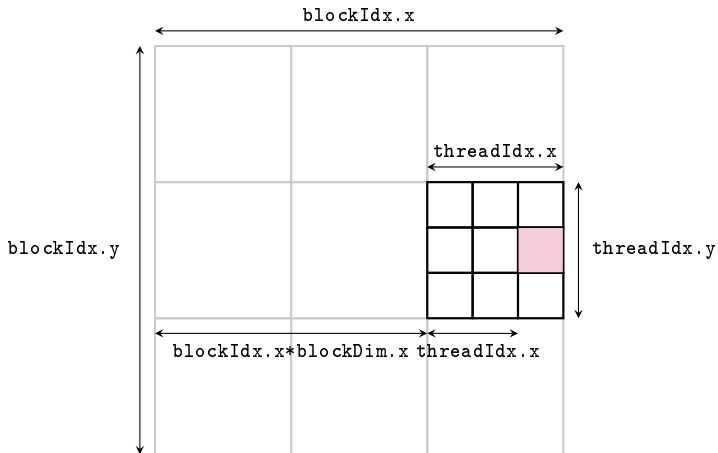

Des blocks et des threads



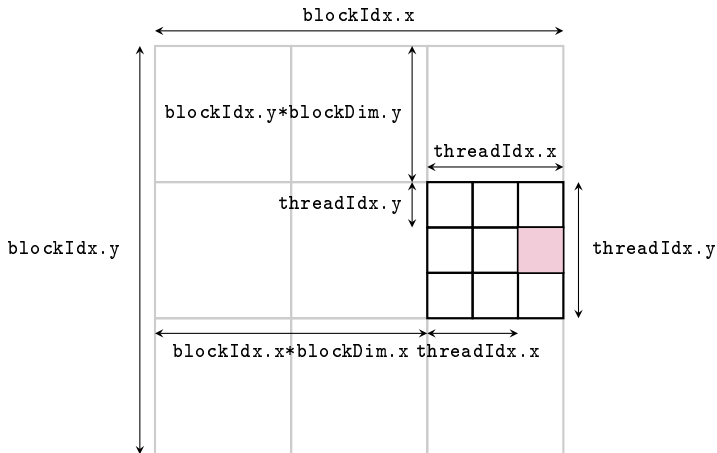
Des blocks et des threads



Des blocks et des threads

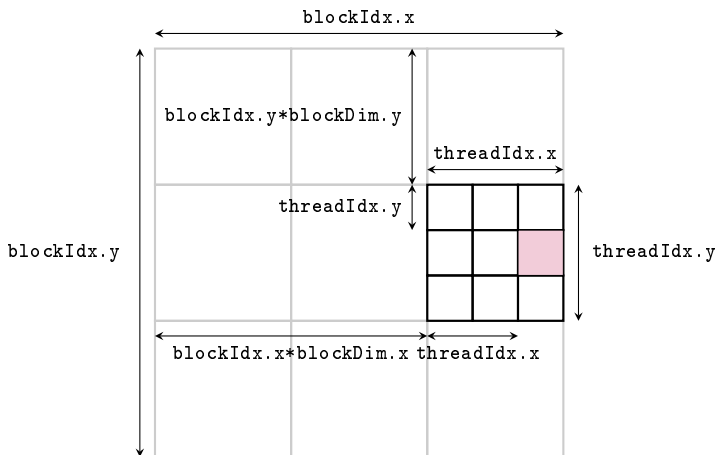


Des blocks et des threads



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

Pourquoi s'embêter avec des threads et des blocks ?

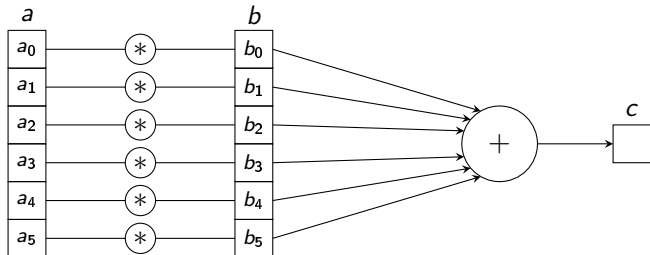
- 1 thread = 1 coeur
- 1 block = 1SM
- 1 block est exécuté sur 1SM
- Blocks :
 - Les blocks sont exécutés dans n'importe quel ordre, séquentiellement ou en parallèle
 - L'avantage est que ça scale automatiquement avec le nombre de SM
- Threads
 - Contrairement aux blocks, les threads peuvent
 - Communiquer
 - Se synchroniser
 - Ces opérations sont à l'intérieur d'un block

Les instructions de contrôles dans les warps

- Les instructions de contrôle (if, while, for, switch, do) affectent les performances, car les thread d'un warp vont diverger
- Les différents chemins sont sérialisés
- Lorsque toutes les exécutions sur les différents chemins sont finies, les threads convergent vers le même chemin
- Les conditions doivent minimiser les divergences
 - Par exemple une condition dépendant de $(threadIdx/warp_size)$
- Il faut utiliser plus de threads CUDA que de nombre de cœurs CUDA pour augmenter le parallélisme

Exercise 4

Produit scalaire (dot product)



$$\begin{aligned} c &= (a_0, a_1, a_2, a_3, a_4, a_5) \cdot (b_0, b_1, b_2, b_3, b_4, b_5) \\ &= a_0 b_0 + a_1 b_1 + a_2 b_2 + a_3 b_3 + a_4 b_4 + a_5 b_5 \end{aligned}$$


```
__global__ void dot (int *a, int *b, int *c){
    // chaque thread calcule le produit d'une paire
    int tmp = 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


```
#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
    tmp[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 + tmp[i];
        *c = sum;
    }
}
```


Synchronisation des threads d'un même block

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

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
    tmp[threadIdx.x] = a[threadIdx.x] * b[threadIdx.x]
};

//Synchronisation pour etre sur que tout les
//threads ont fini
__syncthreads();

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

— $\sqrt{1 - \frac{1}{n}}$

Programmation parallèle en CUDA : le main

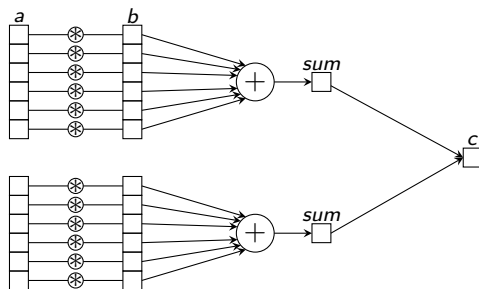
```
// Copie des donnees vers le Device
cudaMemcpy (gpu_a, a, size, cudaMemcpyHostToDevice);
cudaMemcpy (gpu_b, b, size, cudaMemcpyHostToDevice);

//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;
    tmp[threadIdx.x] = 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 + tmp[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 ininterruptibles
 - 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;
    tmp[threadIdx.x] = 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 + tmp[i];
            atomicAdd(*c, sum);
        }
    }
}

```


Mesure du temps

- Les transferts de données sont synchrones
- L'appel au kernel ne l'est pas

```
cudaMemcpy(d_x, x, N*sizeof(float),  
           cudaMemcpyHostToDevice);  
cudaMemcpy(d_y, y, N*sizeof(float),  
           cudaMemcpyHostToDevice);
```

```
t1 = myCPUTimer();  
saxpy<<<(N+255)/256, 256>>>(N, 2.0, d_x, d_y);  
cudaDeviceSynchronize();  
t2 = myCPUTimer();
```

```
cudaMemcpy(y, d_y, N*sizeof(float),  
           cudaMemcpyDeviceToHost);
```

Mesure du temps

- CUDA event API
- Les opérations sont séquentielles sur le GPU

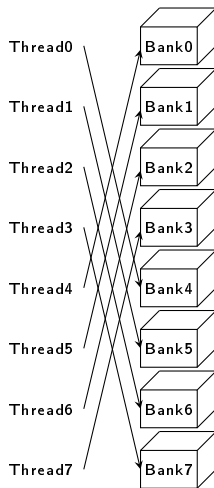
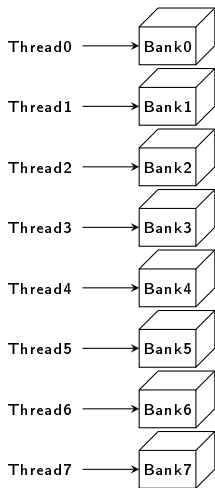
```

cudaEvent_t start, stop;
cudaEventCreate(&start);
cudaEventCreate(&stop);
...
cudaEventRecord(start);
saxpy<<<(N+255)/256, 256>>>(N, 2.0f, d_x, d_y);
cudaEventRecord(stop);

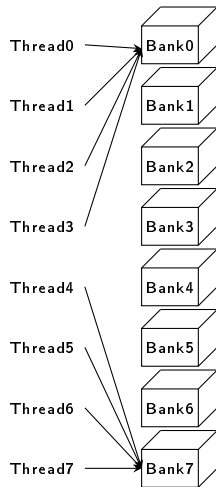
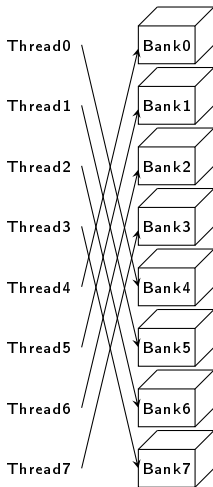
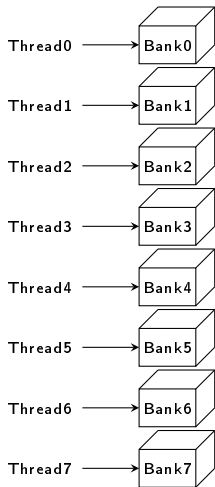
cudaEventSynchronize(stop); // Garantit que l'
    evenement s'est execute
...
float milliseconds = 0;
cudaEventElapsedTime(&milliseconds, start, stop);

```


La mémoire partagée



Bank conflict



Exercise 6

Nombre de threads dans les blocks pour le produit de matrices

- Chaque SM du Quadro K620 possède 48KB de mémoire partagée par block
- $TILE_WIDTH = 16$:
 - 2 vecteur de mémoire partagée de float d'une taille de
→ $2 \times 4B \times 16 \times 16 = 2KB$ pour un block
 - Donc jusqu'à 24 blocks pouvant s'exécuter sur le même SM
- $TILE_WIDTH = 32$
 - 2 vecteur de mémoire partagée de float d'une taille de
→ $2 \times 4B \times 32 \times 32 = 8KB$ pour un block
 - Donc jusqu'à 6 blocks pouvant s'exécuter sur le même SM
- Mais dans la vraie vie, on utilise des bibliothèques (cuBlas)

Accès à la mémoire

Exemple d'utilisation de la mémoire unifiée

```
int *a, *b;
int *gpu_a, *gpu_b;
cudaMalloc((void **)&gpu_a, N*sizeof(int));
cudaMalloc((void **)&gpu_b, N*sizeof(int));
a=(int*) malloc (N*sizeof(int));
b=(int*) malloc (N*sizeof(int));
for (int i = 0; i < N; i++) {
    x[i] = 1.0f;
    y[i] = 2.0f;
}
cudaMemcpy (gpu_a,a,N*sizeof(int),cudaMemcpyHostToDevice);
cudaMemcpy (gpu_b,b,N*sizeof(int),cudaMemcpyHostToDevice);
add <<< N, 1 >>> (gpu_a, gpu_b);
cudaMemcpy(b, gpu_b, N*sizeof(int), cudaMemcpyDeviceToHost);
free(a); free(b);
cudaFree (gpu_a);
cudaFree (gpu_b);
```

Exemple d'utilisation de la mémoire unifiée

```
int *a, *b;
int *gpu_a, *gpu_b;
cudaMalloc((void **)&gpu_a, N*sizeof(int));
cudaMalloc((void **)&gpu_b, N*sizeof(int));
a=(int*) malloc (N*sizeof(int));
b=(int*) malloc (N*sizeof(int));
for (int i = 0; i < N; i++) {
    x[i] = 1.0f;
    y[i] = 2.0f;
}
cudaMemcpy (gpu_a,a,N*sizeof(int),cudaMemcpyHostToDevice);
cudaMemcpy (gpu_b,b,N*sizeof(int),cudaMemcpyHostToDevice);
add <<< N, 1 >>> (gpu_a,gpu_b);
cudaMemcpy(b,gpu_b,N*sizeof(int),cudaMemcpyDeviceToHost);
free(a);free(b);
cudaFree (gpu_a);
cudaFree (gpu_b);
```


[illegible]

Parallélisme de tâches avec des GPU : les streams

Les streams

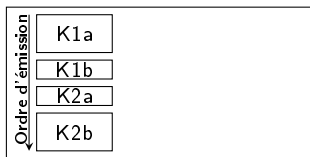
- Un stream : Une séquence d'opérations (comme une file d'attente pour le device)
- Stream par défaut : stream 0
 - Opérations de lecture et d'écriture synchrone (*HostToDevice* et *DeviceToHost*) avec le host
 - Les *kernels* sont asynchrones avec le host par défaut (des opérations CPU sont possibles)

Exemple de code avec streams

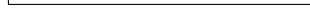
```
cudaStream_T stream1, stream2, stream3, stream4;
cudaStreamCreate(&stream1);
...
cudaMalloc(&data_dev1, size);
...
cudaMemcpyAsync(data_dev1, data_host1, size,
    cudaMemcpyHostToDevice, stream1);
kernel2 <<< grid, block, 0, stream2 >>> (... ,
    data_dev2, ...);
kernel3 <<< grid, block, 0, stream3 >>> (... ,
    data_dev3, ...);
cudaMemcpyAsync(data_host4, data_dev4, size,
    cudaMemcpyDeviceToHost, stream4);
...
```

L'ordre d'émission des opérations est important

- Stream 1 : K1a, K1b
- Stream 2 : K2a, K2b

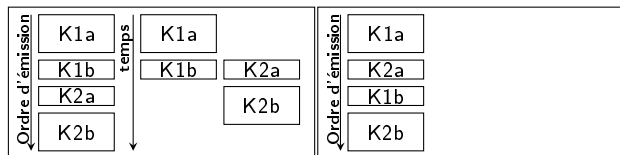


-



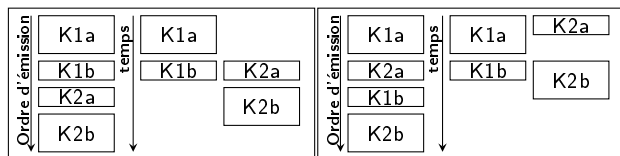
L'ordre d'émission des opérations est important

- Stream 1 : K1a, K1b
- Stream 2 : K2a, K2b



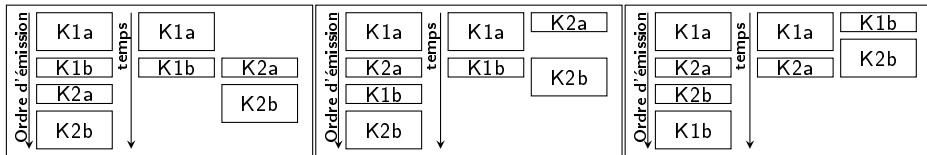
L'ordre d'émission des opérations est important

- Stream 1 : K1a, K1b
- Stream 2 : K2a, K2b



L'ordre d'émission des opérations est important

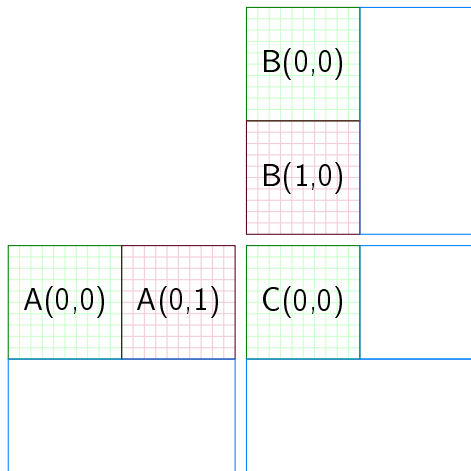
- Stream 1 : K1a, K1b
- Stream 2 : K2a, K2b



- MPI + CUDA :
<http://on-demand.gputechconf.com/gtc/2014/presentations/S4236-multi-gpu-programming-mpi.pdf>
- Mémoire unifiée CPU + GPU :
<http://www.drdobbs.com/parallel/unified-memory-in-cuda-6-a-brief-overview/240169095?pgno=1>

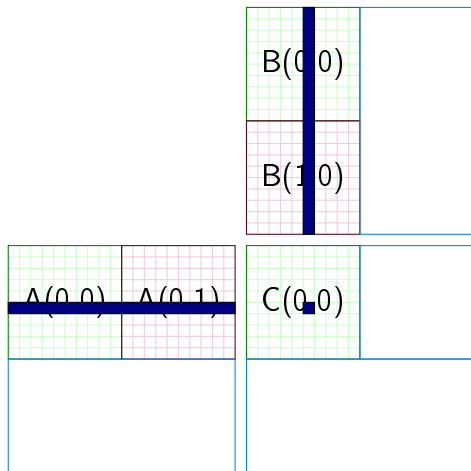
Multiplication de matrice

- $C(0,0) \rightarrow$
blockIdx.x=0,blockIdx.y=0
- $C(0,1) \rightarrow$
blockIdx.x=1,blockIdx.y=0
- $C(1,0) \rightarrow$
blockIdx.x=0,blockIdx.y=1
- $C(1,1) \rightarrow$
blockIdx.x=1,blockIdx.y=1



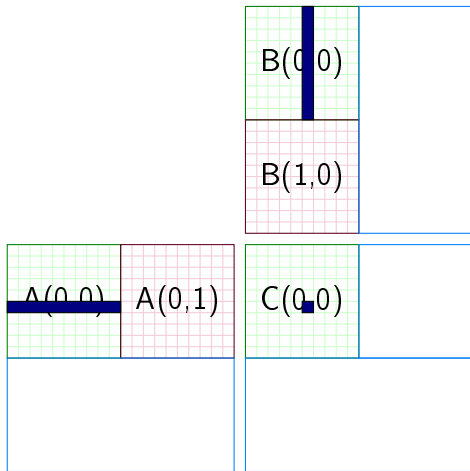
Multiplication de matrice

- 1 block calcule une tuile
- 1 thread calcule un élément de la tuile



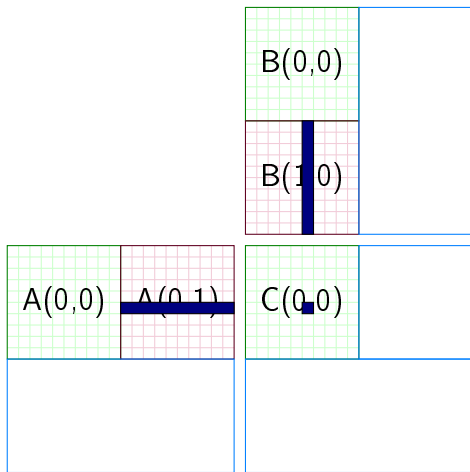
Multiplication de matrice

- Les blocks $A(0,0)$ et $B(0,0)$ sont chargés en mémoire partagée de façon collaborative par les threads (chacun charge un élément de chaque tuile)
- Les threads font le calcul partiel de C
- Chaque valeur est gardée dans le registre du thread



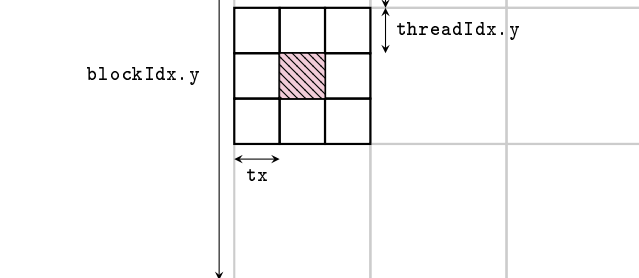
Multiplication de matrice

- Les blocks A(0,1) et B(1,0) sont chargés en mémoire partagée
- Les threads terminent le calcul de C
- À la fin, C est stockée en mémoire globale



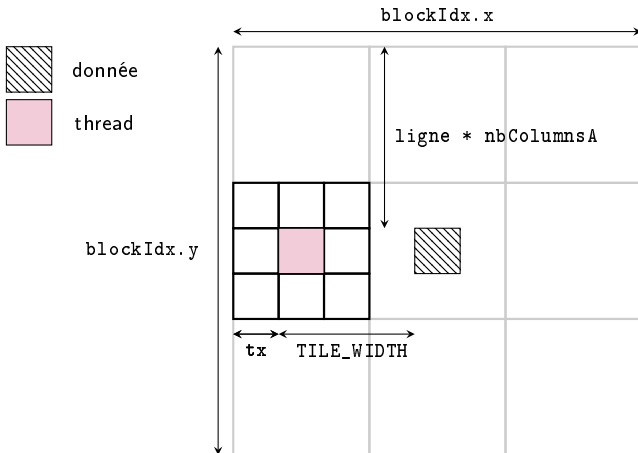
- `col = blockIdx.x*blockDim.x+threadIdx.x`

-



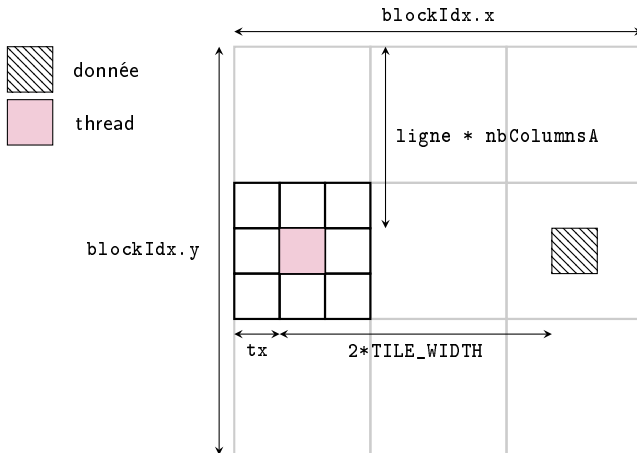
Les indices dans le produit de matrices par tuile : A

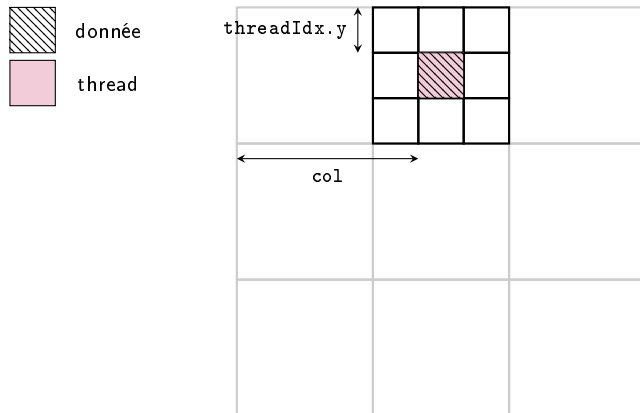
- `col = blockIdx.x*blockDim.x+threadIdx.x`
- `ligne = blockIdx.y*blockDim.y+threadIdx.y`
- `TILE WIDTH = blockDim.x`



Les indices dans le produit de matrices par tuile : A

- `col = blockIdx.x*blockDim.x+threadIdx.x`
- `ligne = blockIdx.y*blockDim.y+threadIdx.y`
- `TILE_WIDTH = blockDim.x`





80 / 82

80 / 82

80 / 82

Les indices dans le produit de matrices par tuile : B

- `col = blockIdx.x*blockDim.x+theadIdx.x`
- `ligne = blockIdx.y*blockDim.y+theadIdx.y`
- `TILE_WIDTH = blockDim.x`
- `indice = col + (nb_colB * TILE_WIDTH * tileId + ty * nb_col)`

