

Introduction à CUDA C

Amina Guermouche

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
()
```

Fonctionnement

- 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

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\n", prop.totalGlobalMem);
}
```

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

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

Exercise 1

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);
```

```
    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

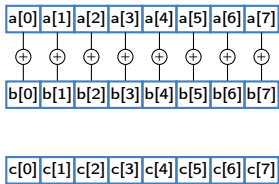
- Comment exécuter le code en parallèle?

On veut faire N fois add en parallèle

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

```
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];
}
```

Block₀
$$c[0] = a[0] + b[0]$$
Block₁
$$c[1] = a[1] + b[1]$$
Block₂
$$c[2] = a[2] + b[2]$$
Block₃
$$c[3] = a[3] + b[3]$$

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);

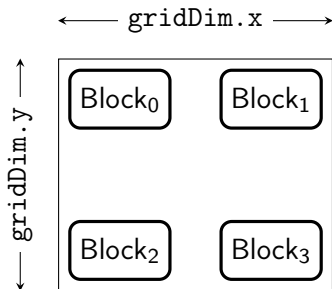
//Copie du resultat
cudaMemcpy(c, gpu_c, size , cudaMemcpyDeviceToHost);

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

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

Grille et blocks



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);
    ...
}
```


Threads

- 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;
}
```


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;
}
```

Vecteurs de taille quelconque

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

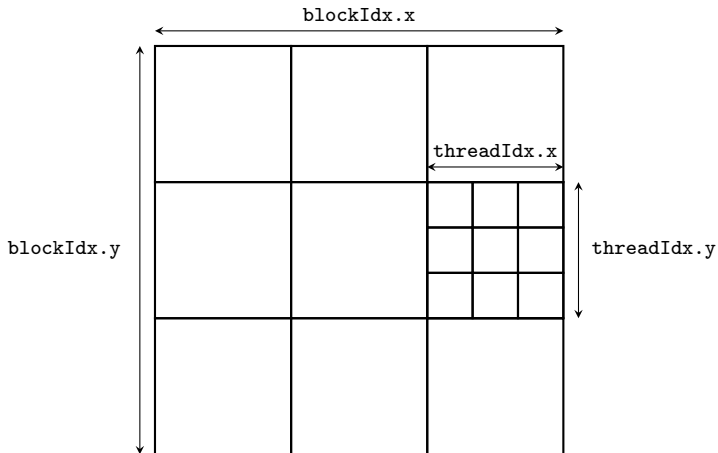
```
__global__ void add (int *a, int *b, int *c,
    int n){
    int indice = threadIdx.x + blockIdx.x *
        blockDim.x;
    if (indice < n)
        c[indice] = a[indice] + b[indice];
}
```

- Au niveau du main

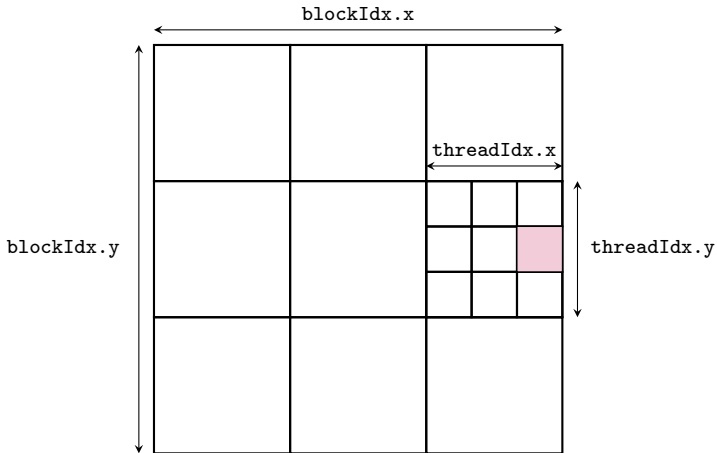
```
add <<< (N+THREAD_PER_BLOCK-1)/
    THREAD_PER_BLOCK, THREAD_PER_BLOCK >>> (
    gpu a, gpu b, gpu c, N);
```

Exercise 3

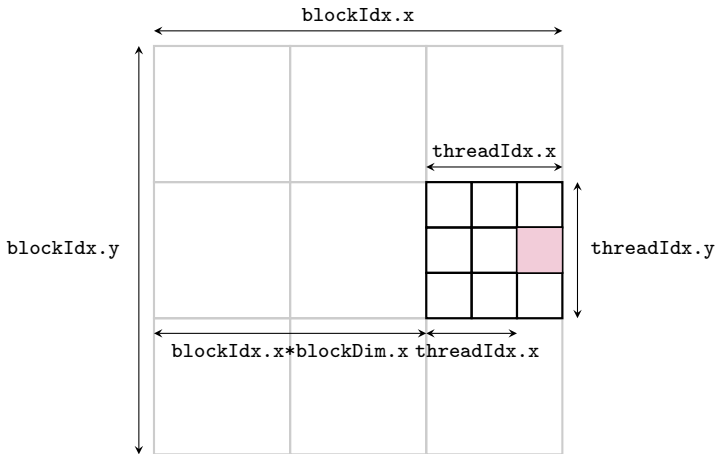
Des blocks et des threads



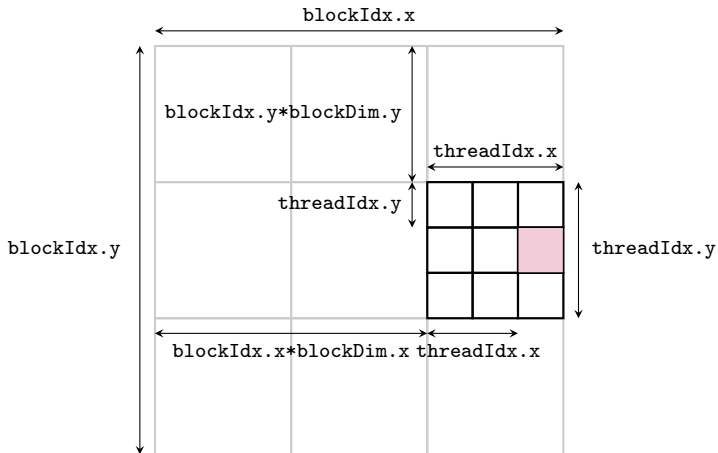
Des blocks et des threads



Des blocks et des threads

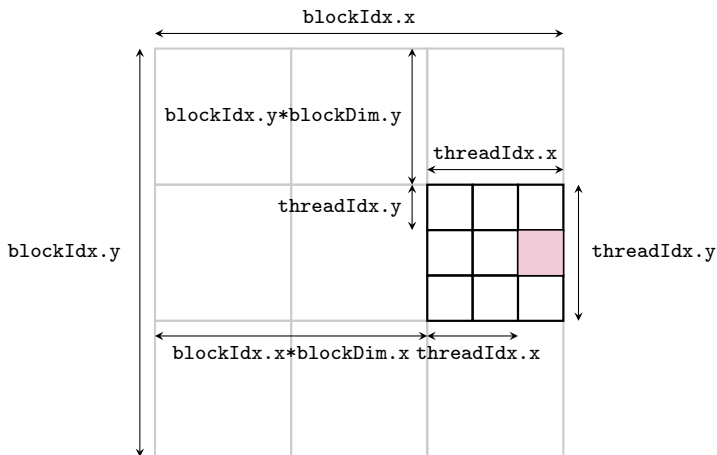


Des blocks et des threads



Addition de 2 matrices

- `colonne = blockIdx.x*blockDim.x+theadIdx.x`
- `ligne = blockIdx.y*blockDim.y+theadIdx.y`

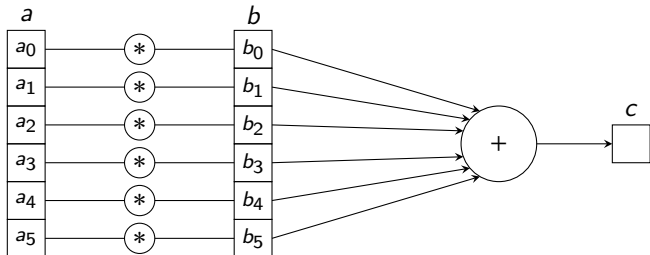


Pourquoi s'embêter avec des threads et des blocks?

- 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

Exercise 4

Produit scalaire (dot product)



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

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 tmp = 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 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

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
    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;
    }
}

```

Produit scalaire (dot product) : 1 seul block

```
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, sizeof(int));

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

Programmation parallèle en CUDA : le main

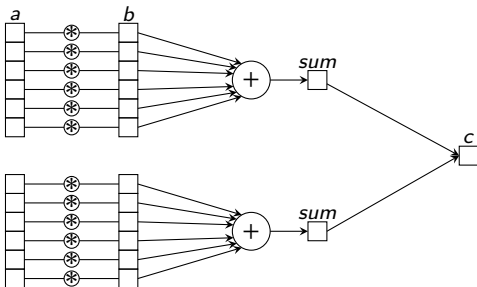
```
// Copie des donnees vers le Device
cudaMemcpy (gpu_a,a,size ,cudaMemcpyHostToDevice);
cudaMemcpy (gpu_b,b,size ,cudaMemcpyHostToDevice);

//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



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)

```
#defin 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);
        }
    }
}
```

Produit scalaire (dot product)

```

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);
    c=(int*) malloc (sizeof(int));
    random_ints(a, N);
    random_ints(b, N);

```

51 / 84

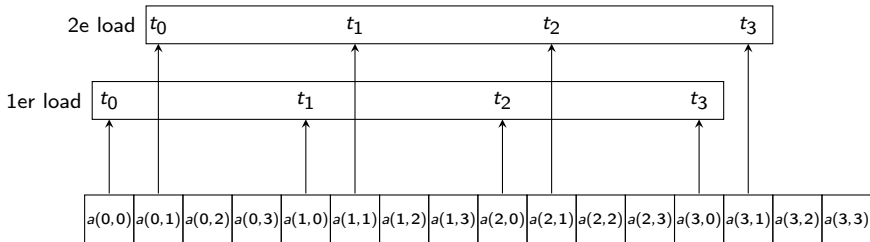
Exercise 5

Throughput de la mémoire

- 1 Minimiser les transferts de faible BW
 - Minimiser les transferts Host \leftrightarrow Device
- 2 Minimiser les transferts mémoire globale \leftrightarrow Device
 - Favoriser la mémoire partagée et les caches
 - La mémoire partagée est équivalente à un cache géré par l'utilisateur

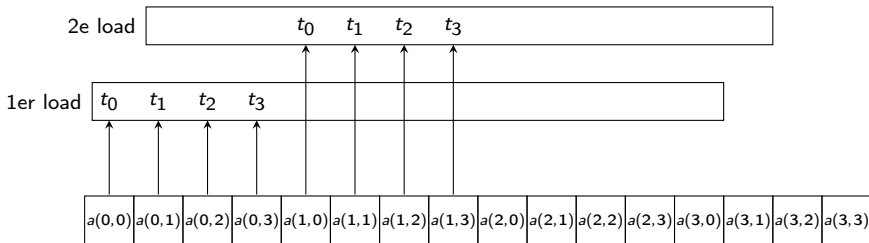
La mémoire globale

- Lorsque tous les threads font un load, le hardware détecte si les threads accèdent à un espace mémoire contigu
- Dans ce cas, le hardware groupe (*coalesces*) les accès en un seul accès à différentes location de la DRAM



La mémoire globale

- Lorsque tous les threads font un load, le hardware détecte si les threads accèdent à un espace mémoire contigu
- Dans ce cas, le hardware groupe (*coalesces*) les accès en un seul accès à différentes location de la DRAM



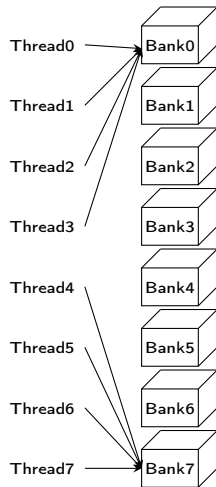
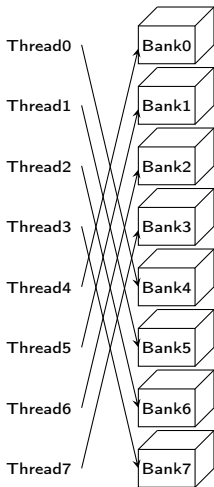
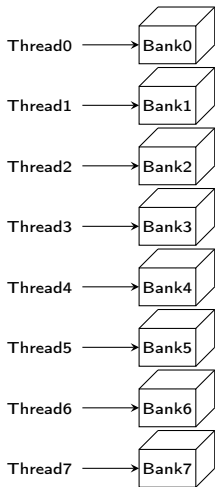
La mémoire partagée

- À utiliser pour éviter les accès non alignés
 - La mémoire est partagée en modules de taille égale, appelés *bank*
 - Il y a autant de *bank* que de threads dans un *warp*
 - L'accès aux *bank* est simultané
- Toute lecture/écriture de n adresses dans n bank différents est simultané

Bank conflict

- Si deux accès sont dans différentes adresses du même *bank*, il y a conflit
- ☹ L'accès en cas de conflit est sérialisé
- La mémoire partagée est rapide **tant qu'il n'y a pas de *bank conflict***
- Le hardware divise un accès mémoire avec conflit en autant d'accès nécessaire pour ne plus avoir de conflit
- ☹ La BW est réduite par un facteur égal au nombre d'accès créés pour éviter les conflits

Bank conflict



Mesure du temps

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

```

cudaEvent_t start, stop;
cudaEventCreate(&start);
cudaEventCreate(&stop);
...
cudaEventRecord(start); // le temps est sauvegarde
                          sur le GPU
saxpy<<<(N+255)/256, 256>>>(N, 2.0f, d_x, d_y);
cudaEventRecord(stop);

cudaEventSynchronize(stop); // Garantit que l'
                              evenement s'est execute : le CPU se bloque en
                              attendant le record de l'evenement

...
float milliseconds = 0;
cudaEventElapsedTime(&milliseconds, start, stop);

```


Exercise 6

Profiling d'une application CUDA

① nvprof (*deprecated*)

- `nvprof ./saxpy`
- `nvprof -o trace.prof ./saxpy`
- Ensuite ouverture de la trace avec : `nvvp -import trace.prof`

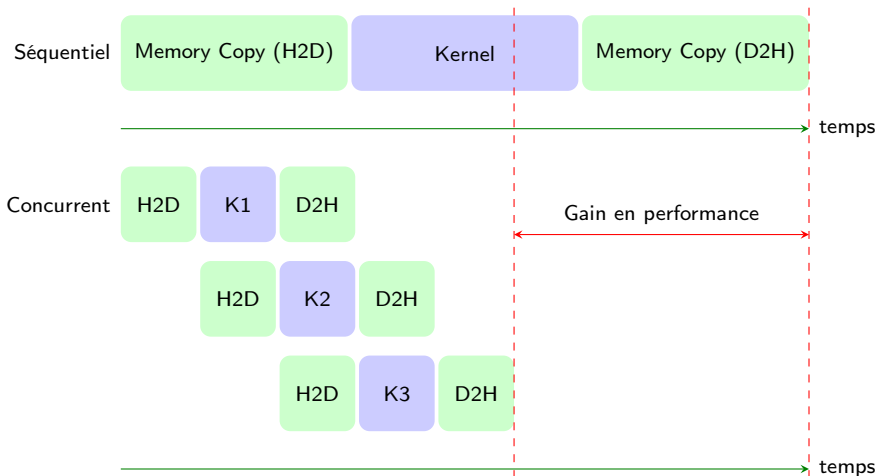
② nsys

- `nsys profile ./saxpy`
- Visualisation de la trace avec : `nsys-ui fichier.nsys-rep`
- Statistiques sur une exécution : `nsys stats fichier.nsys-rep`

Mesure de la puissance consommée

- `nvidia-smi` (System Management Interface) permet de questionner l'état du GPU
- Pour afficher la mesure toutes les secondes

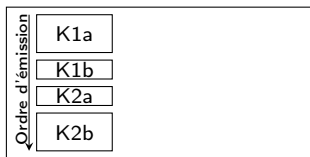
```
nvidia-smi -query-gpu=power.draw -format=csv -l 1
```

69 / 84

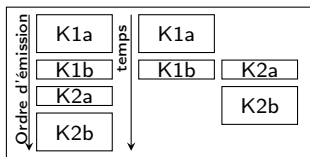
L'ordre d'émission des opérations est important

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



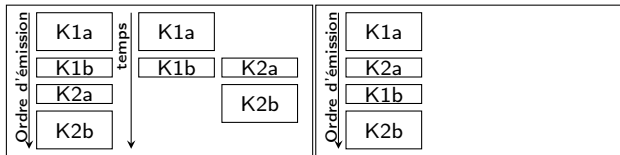
L'ordre d'émission des opérations est important

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



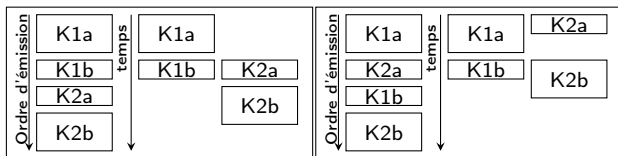
L'ordre d'émission des opérations est important

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



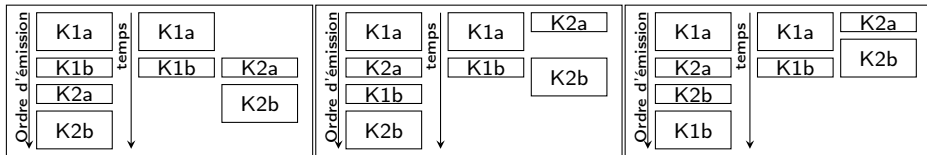
L'ordre d'émission des opérations est important

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



L'ordre d'émission des opérations est important

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



streams (1/2)

- Problèmes de performances simples à corriger (à vous de jouer : exercice 7)

streams (1/2)

- Problèmes de performances simples à corriger (à vous de jouer : exercice 7)
- Symptômes
 - Un stream ne s'overlap pas avec les autres

Étude de différents cas d'utilisation (et d'erreurs) avec les streams (1/2)

- Problèmes de performances simples à corriger (à vous de jouer : exercice 7)
- Symptômes
 - Un stream ne s'overlap pas avec les autres
⇒ Chercher le stream par défaut
- Solutions
 - Lancer bien un kernel par stream (on évite le stream par défaut)

streams (1/2)

- Problèmes de performances simples à corriger (à vous de jouer : exercice 7)
- Symptômes
 - Un stream ne s'overlap pas avec les autres
 - ⇒ Chercher le stream par défaut
 - ⇒ Cherchez `cudaEventRecord`
- Solutions
 - Lancer bien un kernel par stream (on évite le stream par défaut)
 - Associer le stream à `cudaEventRecord`

streams (1/2)

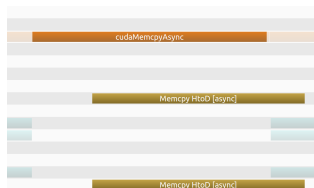
- Problèmes de performances simples à corriger (à vous de jouer : exercice 7)
- Symptômes
 - Un stream ne s'overlap pas avec les autres
 - ⇒ Chercher le stream par défaut
 - ⇒ Cherchez `cudaEventRecord`
 - Les copies mémoires ne s'overlappent pas
- Solutions
 - Lancer bien un kernel par stream (on évite le stream par défaut)
 - Associer le stream à `cudaEventRecord`
 - Utiliser la version asynchrone `cudaMemcpyAsync`
 - Associer un stream différent à `cudaMemcpyAsync`

Référence : <https://on-demand.gputechconf.com/gtc/2014/presentations/S4158-cuda-streams-best-practices-common-pitfalls.pdf>

Pour aller plus loin dans la gestion de la mémoire

- Beaucoup de temps est passé dans l'API pour la copie (`cudaMemcpy`)
 - Cuda indique que la mémoire est *pageable*
- ⇒ La mémoire est transférée via le host
- ⇒ La mémoire peut-être page-in et out par le système

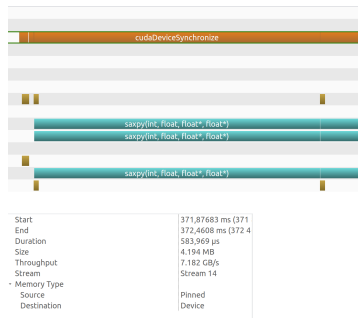
⚠ Afin que `cudaMemCpyAsync` soit réellement asynchrone (en plus des streams), il faut que la mémoire soit pinnée, sinon les opérations sur le GPU sont sérialisées.



Properties	
Memcpy HtoD [async]	
Start	401,46604 ms (401,46604 ms)
End	402,24226 ms (402,24226 ms)
Duration	776,225 µs
Size	4.194 MB
Throughput	5.403 GB/s
Stream	Stream 14
Memory Type	
Source	Pageable
Destination	Device

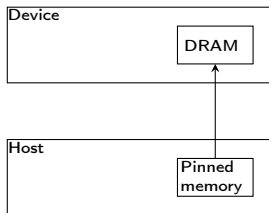
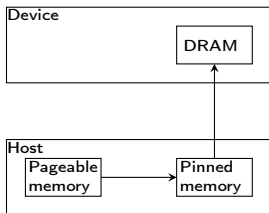
Pour aller plus loin dans la gestion de la mémoire

- Utiliser un mécanisme pour *pin* la mémoire
 - Elle devient accessible directement depuis le GPU
 - Il n'est pas possible pour le système de faire page-in ou page-out
- ⇒ La mémoire transférée via le moteur DMA
- `cudaHostMalloc` à la place de `malloc` (ou `cudaHostRegister`)
 - Bande passante plus élevée
 - Libère le CPU pour une exécution asynchrone



On remarquera que le `MemcpyAsync` a disparu

Pour aller plus loin dans la gestion de la mémoire



```
cudaHostRegister(y, N*sizeof(float),0);
```

```
for(int i = 0; i < 10; i++){
    cudaMemcpyAsync(gpu_y, y, N*sizeof(float),
        cudaMemcpyHostToDevice, stream2);
    saxpy<<<1, 1, 0, stream1>>>(N, 2.0f, gpu_x,
        gpu_y);
    cudaDeviceSynchronize();
}
cudaHostUnregister(y);
```


Accès à la mémoire

Différents types de mémoire

- La mémoire du GPU
 - Allouée via `cudaMalloc`
 - Ne peut pas être paged
- La mémoire du CPU
 - Allouée avec `malloc`(`callo`, `new`, ...)
 - Peut-être page-in et page-out par le système
- La mémoire pinned sur le CPU
 - Allouée via `cudaMallocHost`
 - Ne peut-être page-in et page-out par le système
- La mémoire mapped
 - Mémoire pinned mais mappée dans la mémoire du GPU
 - Les données ne sont pas sur le GPU : elles sont copiées à l'exécution
- Mémoire unifiée : une même zone mémoire accessible directement depuis le CPU et le GPU (le meilleur des deux mondes)

Référence : <https://on-demand.gputechconf.com/gtc/2014/presentations/S4158-cuda-streams-best-practices-common-pitfalls.pdf>

Différents accès à la mémoire

- Accès via des `memcpy` explicites
 - 😊 Bonnes performances (car les données sont disponibles)
 - 😞 Code assez lourd et erreurs fréquentes
 - 😞 Accès uniquement à la mémoire du GPU
- Accès via un mécanisme *Zero Copy* (non abordé dans ce cours) : les thread GPU accèdent directement à l'espace mémoire du CPU (via les fonctions `cudaHostAlloc` et `cudaHostGetDevicePointer`).
 - 😊 Accès à toute la mémoire
 - 😞 Vitesse limitée par le bus PCI et NVLink
 - 😞 Pas de localité des données : les données ne sont pas copiées sur le GPU. Le transfert se fait à l'exécution
- Mémoire unifiée : une même zone mémoire accessible directement depuis le CPU et le GPU (le meilleur des deux mondes)

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);

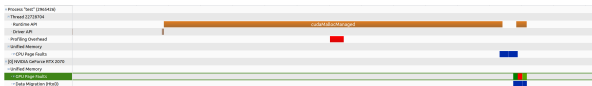
```

1000

1. *Journal of Management Studies*, 1990, 27, 1.

Mémoire unifiée : le prefetching

Sans Prefetch



CPU Page Faults

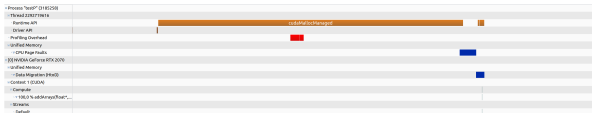
The segment mode is used for this timeline. In this mode the timeline is split into equal width segments and only aggregated data values for each time segment are shown.

Duration	Session
121.8045 ms (121)	1

The time taken to resolve CPU page faults within the segment:

Duration	Session
10 - 10 % [0 - 354.244 μs]	1
10 - 20 % [354.244 μs - 708.488 μs]	1
20 - 30 % [708.488 μs - 1,062.73 μs]	1
30 - 40 % [1,062.73 μs - 1,416.98 μs]	1
40 - 50 % [1,416.98 μs - 1,771.22 μs]	1
50 - 60 % [1,771.22 μs - 2,125.46 μs]	1
60 - 70 % [2,125.46 μs - 2,479.71 μs]	1
70 - 80 % [2,479.71 μs - 2,833.95 μs]	1

Avec Prefetch



CPU Page Faults

The segment mode is used for this timeline. In this mode the timeline is split into equal width segments and only aggregated data values for each time segment are shown.

Duration	Session
121.8045 ms (121)	1

Total CPU Page Faults: 24

The number of CPU page faults per second within the segment:

Duration	Session
10 - 10 % [0 - 354.244 μs]	1
10 - 20 % [354.244 μs - 708.488 μs]	1
20 - 30 % [708.488 μs - 1,062.73 μs]	1
30 - 40 % [1,062.73 μs - 1,416.98 μs]	1
40 - 50 % [1,416.98 μs - 1,771.22 μs]	1
50 - 60 % [1,771.22 μs - 2,125.46 μs]	1
60 - 70 % [2,125.46 μs - 2,479.71 μs]	1
70 - 80 % [2,479.71 μs - 2,833.95 μs]	1

