

# Multi-dimensional grids/blocks and coalescence

Oguz Kaya

Assistant Professor  
Université Paris-Saclay, Gif-sur-Yvette, France

# Objectifs

- Utiliser l'indexage de grilles/blocs 2 ou 3 dims
- Réaliser accès "rapide" à la mémoire globale du GPU
- Premier TP sur multiplication des matrices

# Outline

- 1 Indexage multi-dimensionnel de blocs et threads
- 2 Coalescence
- 3 Lab assignment: Matrix multiplication

- 1 Indexage multi-dimensionnel de blocs et threads
- 2 Coalescence
- 3 Lab assignment: Matrix multiplication

# Multiplication d'un tableau 1D par blocs

Récapultatif de l'exemple de multiplication d'un tableau par blocs

```
#include <stdio>
#include "cuda.h"

#define N 1024
float A[N];
float c = 2.0;

__device__ float dA[N];

__global__ void multiplyArray(int n, float c)
{
    int i = blockIdx.x;
    dA[i] *= c;
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) { A[i] = i; }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    multiplyArray<<<N, 1>>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    return 0;
}
```

- Multiplier un tableau 1D par une grille de blocs.
- Chaque bloc multiplie 1 élément.
- Lancer le kernel avec  $N$  blocs.
- **Question:** Que ferait-on si le tableau était plutôt  $A[N][N]$ ?

# Multiplication d'un tableau 2D par blocs

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int i = blockIdx.x / n;
    int j = blockIdx.x % n;
    dA[i][j] *= c;
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    multiplyArray2D<<<N * N, 1>>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
    return 0;
}
```

- Chaque bloc multiplie 1 élément.
- Il faut lancer  $N^2$  blocs au total.
- Chaque  $N$  blocs consécutifs multiplient une ligne de  $A$ .
- Une division et un modulus pour trouver  $i$  et  $j$

# Multiplication d'un tableau 2D par blocs 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    dA[blockIdx.x][blockIdx.y] *= c;
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    dim3 dimGrid;
    dimGrid.x = N;
    dimGrid.y = N;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, 1>>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
    return 0;
}
```

- **dim3** définit une topologie 3D d'indexage (**dim3.{x,y,z}**).
- Mettre **dim3.z = 1** pour 2D.
- Le nombre total de blocs égale à la multiplication de dimensions.
- Pas de division ni modulus pour trouver  $i$  et  $j$

# Multiplication d'un tableau 2D par blocs 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int i = blockIdx.x;
    int j = blockIdx.y * blockDim.x + threadIdx.x;
    if (j < n) { dA[i][j] *= c; }
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    int blockSize = 1024;
    dim3 dimGrid;
    dimGrid.x = N;
    dimGrid.y = N / blockSize;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, blockSize>>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
    return 0;
}
```

- Chaque bloc multiplie `blockSize` éléments.
- Chaque thread multiplie 1 élément.
- Threads travaillent sur éléments consécutifs dans une **ligne**.
- Il faut lancer  $N^2 / \text{blockSize}$  blocs au total.
- Chaque blocs multiplie une sous-ligne de  $A$ .
- Il faut vérifier le dépassement du tableau.



# Multiplication d'un tableau 2D par blocs 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y;
    if (i < n) { dA[i][j] *= c; }
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    int blockSize = 1024;
    dim3 dimGrid;
    dimGrid.x = N / blockSize;
    dimGrid.y = N;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, blockSize>>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
    return 0;
}
```

- Chaque thread multiplie 1 élément.
- Il faut lancer  $N^2 / blockSize$  blocs au total.
- Threads travaillent sur éléments consécutifs dans une **colonne**.
- Chaque blocs multiplie une sous-colonne de  $A$ .
- Lequel est mieux (sous-ligne vs. sous-colonne)?
- Qu'est-ce qui se passe si  $A$  a peu de lignes/colonnes?

# Multiplication d'un tableau 2D par blocs 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int blockDimSqrt = (int)sqrt((float)blockDim.x);
    int i = blockIdx.x * blockDimSqrt + threadIdx.x / blockDimSqrt;
    int j = blockIdx.y * blockDimSqrt + threadIdx.y % blockDimSqrt;
    if (i < n && j < n) { dA[i][j] *= c; }
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    int blockSize = 1024;
    dim3 dimGrid;
    dimGrid.x = N / 32;
    dimGrid.y = N / 32;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, blockSize>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
    return 0;
}
```

- Répartir 1024 threads en 2D.
- Chaque block s'occupe d'une sous-matrice  $32 \times 32$ .
- Threads consécutifs travaillent sur une sous-ligne.
- Trouver  $i$  et  $j$  nécessite division et modulus.

# Multiplication d'un tableau 2D par blocs 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int blockDimSqrt = (int)sqrt((float)blockDim.x);
    int i = blockIdx.x * blockDimSqrt + threadIdx.x % blockDimSqrt;
    int j = blockIdx.y * blockDimSqrt + threadIdx.y / blockDimSqrt;
    if (i < n && j < n) { dA[i][j] *= c; }
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    int blockSize = 1024;
    dim3 dimGrid;
    dimGrid.x = N / 32;
    dimGrid.y = N / 32;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, blockSize>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
    return 0;
}
```

- Répartir 1024 threads en 2D.
- Chaque block s'occupe d'une sous-matrice  $32 \times 32$ .
- Threads consécutifs travaille sur une **sous-colonne**.
- Trouver  $i$  et  $j$  nécessite division et modulus.
- Lequel est mieux (accès sous-ligne vs. sous-colonne)?

# Multiplication d'un tableau 2D par blocs 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;
    if (i < n && j < n) { dA[i][j] *= c; }
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    dim3 dimBlock;
    dimBlock.x = 32;
    dimBlock.y = 32;
    dimBlock.z = 1;
    dim3 dimGrid;
    dimGrid.x = N / 32;
    dimGrid.y = N / 32;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, dimBlock>>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
    return 0;
}
```

- Utiliser un **dim3** pour l'indexage 2D de threads.
- Threads consécutifs en **threadIdx.x** sont exécutés dans un **warp** (puis en **threadIdx.y**, puis en **threadIdx.z**).
- Donc chaque warp accède à une **sous-colonne** de  $A$ .

# Multiplication d'un tableau 2D par blocs 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int i = blockIdx.x * blockDim.y + threadIdx.y;
    int j = blockIdx.y * blockDim.x + threadIdx.x;
    if (i < n && j < n) { dA[i][j] *= c; }
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    dim3 dimBlock;
    dimBlock.x = 32;
    dimBlock.y = 32;
    dimBlock.z = 1;
    dim3 dimGrid;
    dimGrid.x = N / 32;
    dimGrid.y = N / 32;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, dimBlock>>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%i\n", A[1][2]);
    return 0;
}
```

- Utiliser un **dim3** pour l'indexage 2D de threads.
- Threads consécutifs en **threadIdx.x** sont exécutés dans un **warp** (puis en **threadIdx.y**, puis en **threadIdx.z**).
- Donc chaque warp accède à une **sous-ligne** de  $A$ .
- Lequel est mieux?

# Limites de dimensions de grille et de bloc

Pour une grille, il faut utiliser

- $\text{dim3.x} \leq 2^{31} - 1$
- $\text{dim3.y} \leq 65535$
- $\text{dim3.z} \leq 65535$

Pour un bloc, il faut utiliser

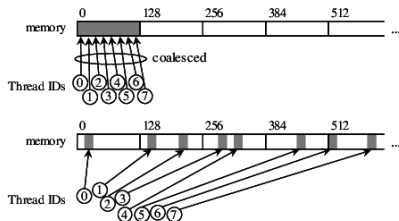
- $\text{dim3.x} \leq 1024$
- $\text{dim3.y} \leq 1024$
- $\text{dim3.z} \leq 64$
- Nombre total de threads  $\leq 1024$

# Outline

- 1 Indexage multi-dimensionnel de blocs et threads
- 2 Coalescence
- 3 Lab assignment: Matrix multiplication

# Coalescence

Il s'agit d'accès à la mémoire globale des threads dans un warp.

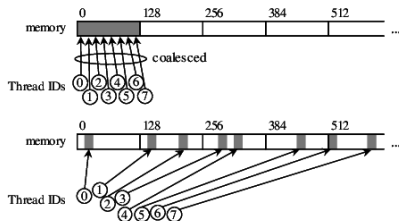


- Threads dans un warp exécute les instructions de manière **synchrone**.
- Chaque accès mémoire est également traité de manière **synchrone**.
- Si les threads accèdent à des éléments consécutifs dans la mémoire, ceci nécessite la lecture d'une ligne (donc 1 accès).



## Coalescence (cont.)

Il s'agit d'accès à la mémoire globale des threads dans un warp.



- S'il y a des sauts, chaque ligne touchée sera lue.
  - À la limite, 32 lignes lues pour un accès.
  - La plupart d'éléments dans la ligne sont inutilisés.
- **Règle:** Concevoir le kernel tel que les accès sont contigus en **threadidx.x**.

# Exemple: Multiplication d'un tableau 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int i = blockIdx.x;
    int j = blockIdx.y * blockDim.x + threadIdx.x;
    if (j < n) { dA[i][j] *= c; }
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    int blockSize = 1024;
    dim3 dimGrid;
    dimGrid.x = N;
    dimGrid.y = N / blockSize;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, blockSize>>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
    return 0;
}
```

- Est-il **coalescent**?
- Oui! La matrice est stockée par lignes, **threadIdx.x** aligné sur lignes.

# Exemple: Multiplication d'un tableau 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y;
    if (i < n) { dA[i][j] *= c; }
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    int blockSize = 1024;
    dim3 dimGrid;
    dimGrid.x = N / blockSize;
    dimGrid.y = N;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, blockSize>>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
    return 0;
}
```

- Est-il **coalescent**?
- Non! La matrice est stockée par lignes, **threadIdx.x** aligné sur colonnes.
- 32 lectures mémoires nécessaires pour chaque accès mémoire d'un warp.

# Exemple: Multiplication d'un tableau 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;
    if (i < n && j < n) { dA[i][j] *= c; }
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    dim3 dimBlock;
    dimBlock.x = 32;
    dimBlock.y = 32;
    dimBlock.z = 1;
    dim3 dimGrid;
    dimGrid.x = N / 32;
    dimGrid.y = N / 32;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, dimBlock>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
}
```

- Est-il **coalescent**?
- Non! La matrice est stockée par lignes, **threadIdx.x** aligné sur colonnes.
- 32 lectures mémoires nécessaires pour chaque accès mémoire d'un warp.

# Exemple: Multiplication d'un tableau 2D

Multiplier chaque élément d'un tableau  $A[N][N]$  par un scalaire  $c$

```
#include <stdio>
#include "cuda.h"

#define N 2048
float A[N][N];
float c = 2.0;

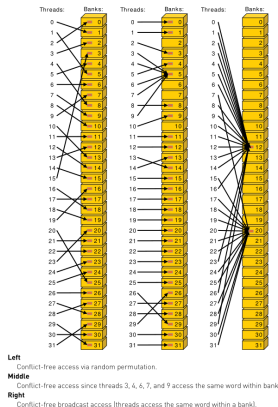
__device__ float dA[N][N];

__global__ void multiplyArray2D(int n, float c)
{
    int i = blockIdx.x * blockDim.y + threadIdx.y;
    int j = blockIdx.y * blockDim.x + threadIdx.x;
    if (i < n && j < n) { dA[i][j] *= c; }
}

int main(int argc, char **argv)
{
    // Initialisation
    for (int i = 0; i < N; i++) {
        for (int j = 0; j < N; j++) { A[i][j] = i + j; }
    }
    // Copier le tableau vers le GPU
    cudaMemcpyToSymbol(dA, A, N * N * sizeof(float), 0,
        cudaMemcpyHostToDevice);
    dim3 dimBlock;
    dimBlock.x = 32;
    dimBlock.y = 32;
    dimBlock.z = 1;
    dim3 dimGrid;
    dimGrid.x = N / 32;
    dimGrid.y = N / 32;
    dimGrid.z = 1;
    multiplyArray2D<<<dimGrid, dimBlock>>>(N, c);
    // Recopier le tableau multiplié vers le CPU
    cudaMemcpyFromSymbol(A, dA, N * N * sizeof(float), 0,
        cudaMemcpyDeviceToHost);
    printf("%f\n", A[1][2]);
}
```

- Est-il **coalescent**?
- Oui! La matrice est stockée par lignes, **threadIdx.x** aligné sur lignes.

# Coalescence rules



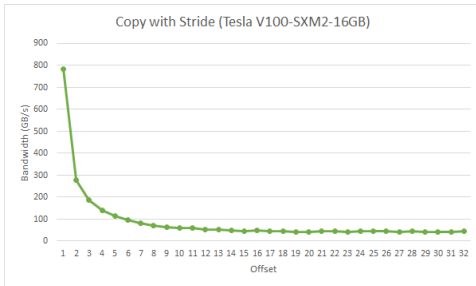
- Threads in a warp accessing the same memory slot = good performance (however, bandwidth is still potentially wasted).
- Threads in a warp accessing the same **memory line** in a **random** order = still good performance in new architectures (Volta and later).
- If coalescent access is difficult to do, **shared memory** might be useful (we will see soon).

# Example: Strided memory access to an array

```
--device.. float dA[N];  
  
--global.. void stridedAccess(int stride)  
{  
    float f = dA[threadIdx.x * stride];  
    // ...  
}
```

- How does performance evolve in terms of **stride**?
- For **stride** = 1, reading a single line of 128 bytes.
- For **stride** = 2, reading two lines of 128 bytes (half of which is unused).
- ...
- For **stride** = 32, reading 32 lines of 128 bytes (31/32 of which is unused).

# Example: Strided memory access to an array (cont.)



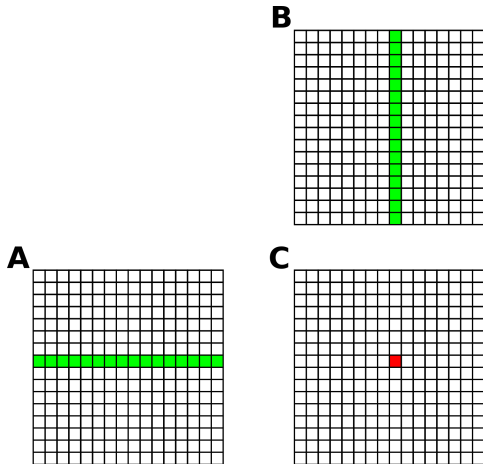
- Effective bandwidth falls rapidly.



- 1 Indexage multi-dimensionnel de blocs et threads
- 2 Coalescence
- 3 Lab assignment: Matrix multiplication

# Matrix multiplication

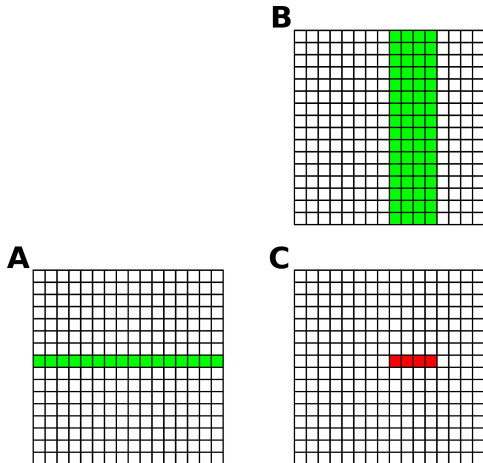
Let  $A, B, C$  be  $N \times N$  matrices.



- The multiplication  $C = AB$  corresponds to the computation  $C[i][j] = \sum_{k=0}^{N-1} A[i][k]B[k][j]$ .
- First kernel: Create one block/thread to compute each  $C[i][j]$ .

# Matrix multiplication

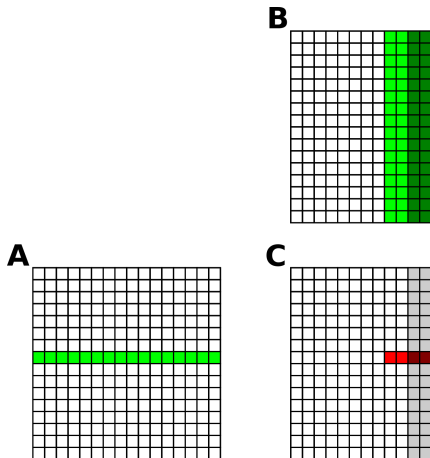
Let  $A, B, C$  be  $N \times N$  matrices.



- The multiplication  $C = AB$  corresponds to the computation  $C[i][j] = \sum_{k=0}^{N-1} A[i][k]B[k][j]$ .
- Second kernel: Use  $P$  threads per block, each block computing  $P$  consecutive elements of a row of  $C$  ( $P = 4$  in this figure).
- Suppose that  $N$  is divisible by  $P$ .

# Matrix multiplication

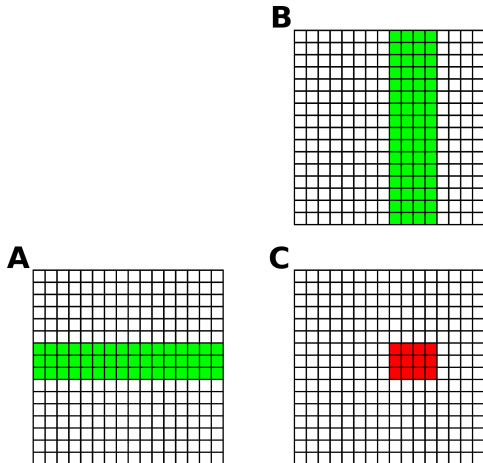
Let  $A, B, C$  be  $N \times N$  matrices.



- The multiplication  $C = AB$  corresponds to the computation  $C[i][j] = \sum_{k=0}^{N-1} A[i][k]B[k][j]$ .
- Second kernel: Use  $P$  threads per block, each block computing  $P$  consecutive elements of a row of  $C$  ( $P = 4$  in this figure).
- Suppose that  $N$  is **not** divisible by  $P$ .

# Matrix multiplication

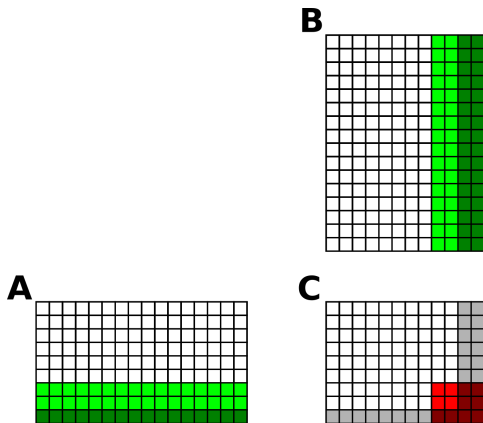
Let  $A, B, C$  be  $N \times N$  matrices.



- The multiplication  $C = AB$  corresponds to the computation  $C[i][j] = \sum_{k=0}^{N-1} A[i][k]B[k][j]$ .
- Fourth kernel: Use  $P \times Q$  threads per block, each block computing  $P \times Q$  consecutive elements in a tile of  $C$  ( $P = 4$  et  $Q = 3$  in this figure).
- Suppose that  $N$  is divisible by  $P$  and  $Q$ .

# Matrix multiplication

Let  $A, B, C$  be  $N \times N$  matrices.



- The multiplication  $C = AB$  corresponds to the computation  $C[i][j] = \sum_{k=0}^{N-1} A[i][k]B[k][j]$ .
- Fourth kernel: Use  $P \times Q$  threads per block, each block computing  $P \times Q$  consecutive elements in a tile of  $C$  ( $P = 4$  et  $Q = 3$  in this figure).
- Suppose that  $N$  is **not** divisible by  $P$  and  $Q$ .

## **Contact**

Oguz Kaya

Université Paris-Saclay and LRI, Paris, France

[oguz.kaya@lri.com](mailto:oguz.kaya@lri.com)

[www.oguzkaya.com](http://www.oguzkaya.com)