

---

# Calcul parallèle avancé

---

**ENSIMAG 3A - IF**

**17 février 2014**

## Chapitre 3 - Patterns de programmation

---

# Agenda

**1 Quelques exemples**

2 Premiers patterns

3 Ameliorations des patterns

## Exemple 1 : Calcul matriciel

- ▶ Tous les threads d'un bloc doivent participer : chaque thread charge un élément de  $M$  et un élément de  $N$  des blocs de la matrice.
- ▶ On associe les éléments chargés à chaque thread en faisant en sorte que les accès soient contigus.
- ▶ Il est nécessaire de synchroniser les threads pour être sûr que tous éléments sont bien chargés et que tous les éléments sont bien utilisés.
- ▶ `__syncthreads()` permet de synchroniser les threads d'un même bloc.

```
//Charger un bloc matriciel
int tx = threadIdx.x
int ty = threadIdx.y
// Accede au bloc 0 avec un indexage 2D :
...M[Row][tx]
...N[ty][Col]
// Accede au bloc 1 avec un indexage 2D :
...M[Row][1*TILE_WIDTH+tx]
...N[1*TILE_WIDTH+ty][Col]
// Accede aux elements avec un indexage 1D
...M[Row*Width + m*TILE_WIDTH + tx]
...N[(m*TILE_WIDTH+ty) * Width + Col]
```

## Exemple 1 : Multiplication par blocs

```
__global__ void MatrixMulKernel(float* d_M, float* d_N, float* d_P, int Width) 1
{
    __shared__ float ds_M[TILE_WIDTH][TILE_WIDTH]; 2
    __shared__ float ds_N[TILE_WIDTH][TILE_WIDTH]; 3
    int bx = blockIdx.x; int by = blockIdx.y; 4
    int tx = threadIdx.x; int ty = threadIdx.y; 5
    // Position de l'element de P sur lequel on travail 6
    int Col = bx * TILE_WIDTH + tx; 7
    float Pvalue = 0; 8
    // Boucle sur l'ensemble les blocs de M et N necessaire pour 9
    // calculer un element de 10
    for (int m = 0; m < Width/TILE_WIDTH; ++m) { 11
        // Chargement collaboratif en memoire partagee 12
        ds_M[ty][tx] = d_M[Row*Width + m*TILE_WIDTH+tx]; 13
        ds_N[ty][tx] = d_N[(m*TILE_WIDTH+ty)*Width+Col]; 14
        __syncthreads(); 15
        for (int k = 0; k < TILE_WIDTH; ++k) 16
            Pvalue += ds_M[ty][k] * ds_N[k][tx]; 17
        __syncthreads(); 18
    } 19
    d_P[Row*Width+Col] = Pvalue; 20
} 21
```

## Exemple 1 : Découpage

Afin de déterminer exactement le matériel, on peut récupérer les informations de façon dynamique

- ▶ Nombres d'accélérateurs dans le système

```
int dev_count; 1  
cudaGetDeviceCount(&dev_count); 2
```

---

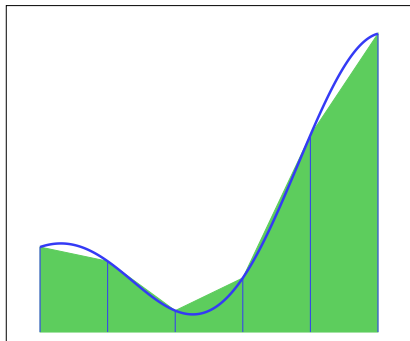
- ▶ Propriétés

```
cudaDeviceProp dev_prop; 1  
for (i = 0; i < dev_count; i++) { cudaGetDeviceProperties( &dev_prop, i) 2  
    ;  
// on peut decider si l'accélérateurs a des ressources suffisantes. 3  
} 4
```

---

- ▶ cudaDeviceProp est une structure C

## Example 2 : Calcul Integral



La formule du trapèze s'écrit  $\forall f \in C^2([a; b]), \exists \xi \in [a; b]$

$$I = h \left( f(a) + 2 \sum_{k=1}^{n-1} f(a + kh) + f(b) \right) - (b - a) \frac{h^2}{12} f^{(2)}(\xi)$$

## Example 2 : Pseudo-code

```
Input b, a, n      1
h = (b-a)/n        2
I = (f(a) + f(b))/2.0;  3
for (i = 0; i<= n-1; i++)  4
{                    5
    x_i = x_i + h    6
    I += f(x_i)      7
}                    8
I = h*I             9
```

## Example 2 : Pseudo-code parallèle

```
Find b, a, n                                1
h = (b-a)/n                                2
local_n = n/n_p                             3
local_a = a + id * local_n*h                4
local_b = local_a + local_n*h               5
local_I = Trap(local_a, local_b, local_n)    6
                                              7
If (id == 0)                                8
{                                             9
    I = local_I;                             10
    for (i = 1; i<= n_p; i++)               11
    {                                         12
        I += local_I;                       13
    }                                         14
    Display I;                               15
}                                             16
```



## Example 3 : Tri Bitonic

- ▶ Une suite bitonic est une suite de nombre  $a_0, a_1, \dots, a_{n-1}$  qui augmente de façon monotone en valeur, atteint un maximum et décroît de manière monotone.

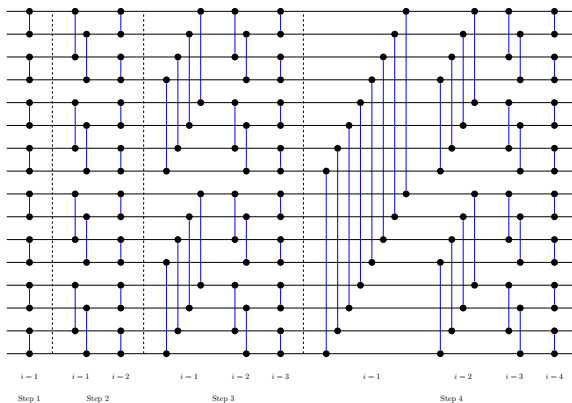
$$a_0 < a_1 < \dots < a_{i-1} < a_i > a_{i+1} > \dots > a_{n-2} > a_{n-1}$$

pour somme valeur de  $i$ . Une suite est également considérée comme bitonic si la relation précédente est atteinte par permutation circulaire des nombres.

- ▶ Une suite de 2 éléments est bitonic.
- ▶ 4 éléments :  $(a_0, a_1, a_2, a_3)$ 
  - ▶ Trier  $(a_0, a_1)$  tel que  $a_0 \leq a_1$
  - ▶ Trier  $(a_2, a_3)$  tel que  $a_2 \geq a_3$
- ▶ 8 éléments :  $(a_0, a_1, a_2, a_3, a_4, a_5, a_6, a_7)$ 
  - ▶ Réaliser un trie bitonic pour  $(a_0, a_1, a_2, a_3)$  et  $(a_4, a_5, a_6, a_7)$ .
  - ▶ Faire une séparation bitonic pour faire en sorte que  $(a_0, a_1, a_2, a_3)$  soit croissante.
  - ▶ Faire une séparation bitonic pour faire en sorte que  $(a_4, a_5, a_6, a_7)$  soit décroissante.

## Example 3 : Idée Séquentielle

- Pour les listes de longueur  $l = 2, 4, 8, \dots 2^m = n$
- Utiliser le tric bitonic afin de créer des groupes successifs de  $l$  éléments alternativement croissants et décroissants.



# Agenda

1 Quelques exemples

2 Premiers patterns

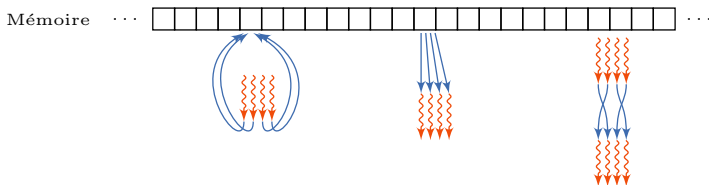
3 Ameliorations des patterns

# Principe

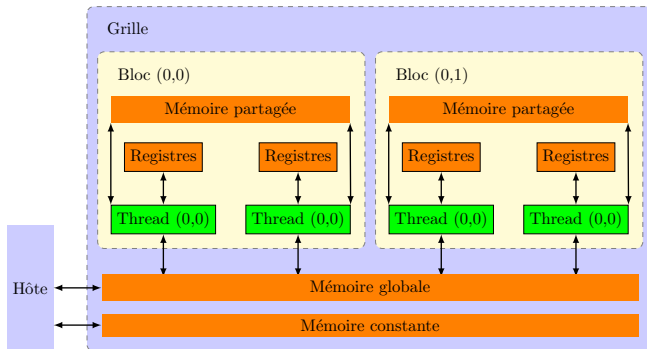
Idée :

Un accélérateur est efficace si un très grand nombre de threads résolvent un problème en travaillant ensemble.

Le point clé est la communication entre les threads.



## Les différents niveaux



- ▶ Registres : accès très haut débit en lecture et écriture (1 cycle) par les threads.
- ▶ Mémoire partagée : accès haut débit en lecture et écriture (5 cycles) par les blocs
- ▶ Mémoire globale : accès bas débit en lecture et écriture (500 cycles) par la grille
- ▶ Mémoire constante : accès haut débit en lecture seule (5 cycles)

## CUDA : les garanties

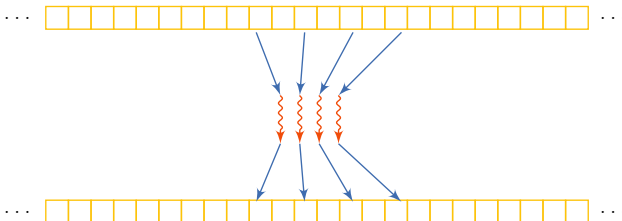
- ▶ Pas d'hypothèses sur l'affectation des blocs sur les unités de calculs.
- ▶ Pas de communication entre les blocs.
- ▶ Les threads, les blocs doivent se terminer avant de pouvoir exécuter les noyaux suivants.
- ▶ Tous les threads d'un bloc s'exécutent sur la même unité de calcul au même moment.
- ▶ Tous les blocs d'un noyau doivent être terminés avant que n'importe quel bloc du noyau suivant puisse être exécuté.

## Identification des patterns de base

```
#include <stdio.h> 1
2
#define NUM_BLOCKS 16 3
#define BLOCK_WIDTH 1 4
5
__global__ void identifierBloc() 6
{ 7
    printf("Je suis un thread dans le bloc %d\n", blockIdx.x); 8
} 9
10
int main(int argc, char **argv) 11
{ 12
    // Lancement du noyau 13
    identifierBloc<<<NUM_BLOCKS, BLOCK_WIDTH>>>(); 14
15
    // Force printf a afficher le resultat 16
    cudaDeviceSynchronize(); 17
18
    printf("Terminer\n"); 19
20
    return 0; 21
} 22
```

## Map pattern

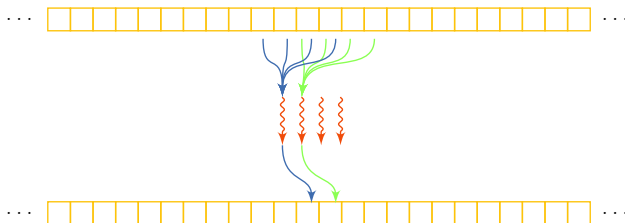
Les threads lisent et écrivent les données à partir d'emplacements mémoires spécifiques distincts.



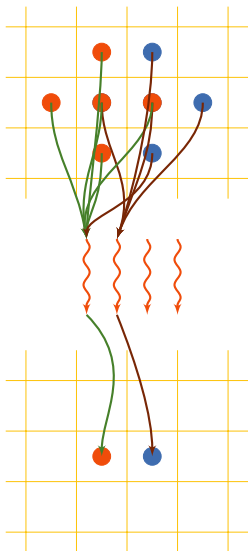


## Gather pattern

Les threads lisent les données à partir d'emplacements mémoires spécifiques et distincts, réalisent une opération sur ces données et écrivent le résultat sur un unique emplacement.

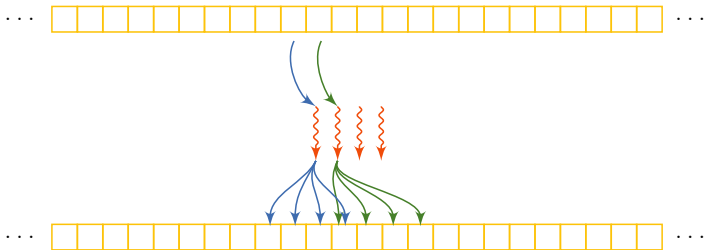


## Gather pattern

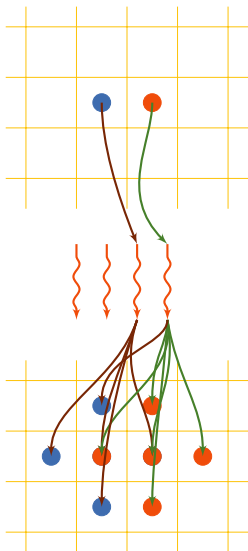


## Scatter pattern

Les threads calculent l'emplacement mémoire dans lequel la résultat sera écrit.

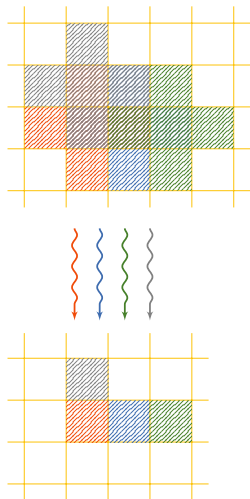


## Scatter pattern (2)



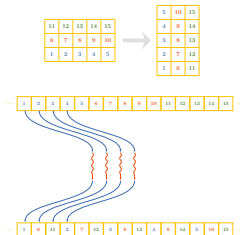
## Motif pattern

Les threads mettent à chaque élément d'un tableau en utilisant les éléments voisins en utilisant un même motif.



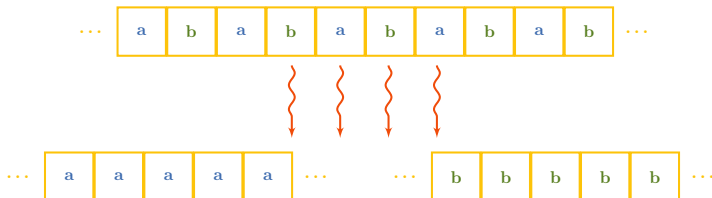
## Transposition pattern

Chaque thread lis un élément du tableau et le réécrit de manière épars, selon la taille de la colonne.



## Transposition pattern (2)

Le concept de transposition peut également à des tableaux de structures : on obtient une opération qui transforme un tableau de structures en structure de tableaux.



## Identification des patterns de base

```
float out[], in[]; 1
int i= threadId.x; 2
int j= threadId.y; 3

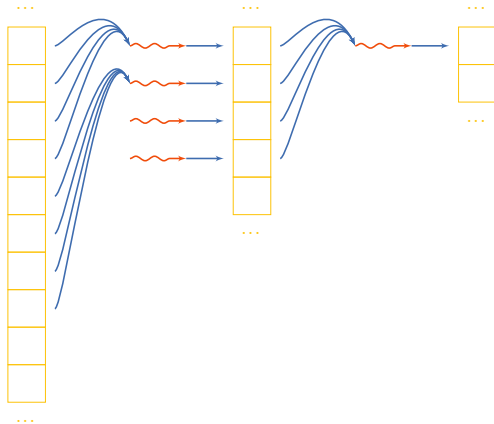
4
const float pi = 3.1415; 5

6
out[i] = pi * in[i]; 7
out[i + j*128] = in[j + i*128]; 8

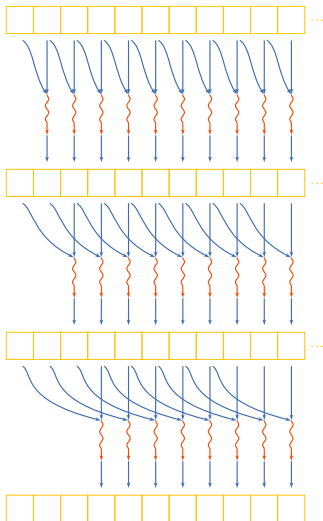
9
if (i %2 ) 10
{ 11
    out[i-1] += pi * in[i]; out[i+1] += pi*in[i]; 12
    out[i] = (in[i] + in[i-1] + in[i+1]) * pi/3.0f; 13
} 14
```



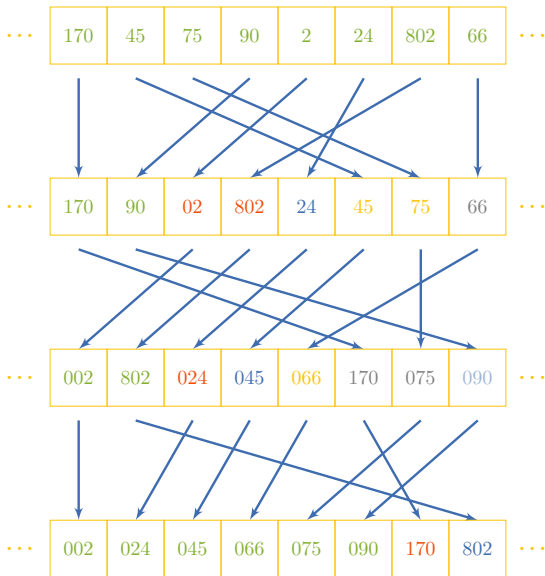
# Réduction



# Scan



# Trie



## Synchronisation

- ▶ Les threads peuvent accéder aux résultats des autres en utilisant la mémoire partagée ou la mémoire globale.
- ▶ Les threads peuvent donc travailler ensemble.
- ▶ A travers les patterns, il est nécessaire de faire particulièrement attention lorsqu'un thread a besoin de lire une donnée pendant qu'un autre l'écrit.
- ▶ Il est nécessaire de synchroniser les threads.

## Exemple de barrière

```
int idx = threadIdx.x; 1
__shared__ int array[128]; 2
3
array[idx] = threadIdx.x; 4
if(idx < 127) 5
    array[idx] = array[idx+1]; 6
```

## Barrière (2)

```
__global__ void foo() 1
{ 2
    __shared__ int s[1024]; 3
    int i = threadIdx.x; 4

    s[i] = s[i-1]; 5
    6
    if(i%2) s[i] = s[i-1]; 7
    8
    s[i] = (s[i-1] -2* s[i] + s[i+1]) /2.0; 9
    printf("s[%d] = %f \n",i,s[i]); 10
} 11
12
```

## Ecrire du code efficace

### 1. Maximiser l'intensité arithmétique

$$A_I = \frac{op}{memoire}$$

- ▶ On peut maximiser le nombre d'opérations par thread
- ▶ On peut minimiser le temps passer dans les accès mémoires : on déplace les données fréquemment utilisées dans les mémoires à accès rapides.
- ▶ On réalise des accès coalescer.

### 2. Eviter les divergences.

# Agenda

1 Quelques exemples

2 Premiers patterns

3 Améliorations des patterns



## Quelques exemples