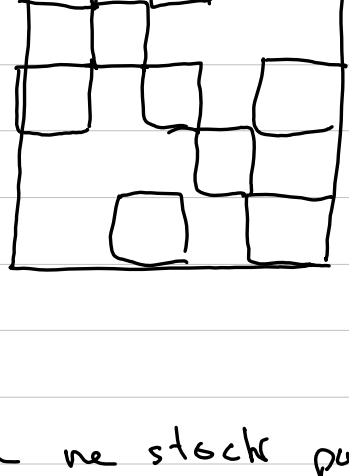


# Multiplication Multithread

## Matrices Sparse Symétriques par block



- Matrices composées majoritairement de 0
- Données en blocks (taille 2, 3, 4 et symétrique)

• On ne stocke pas les 0  
On stocke donc pour chaque colonne les block.

Forme classique:

Stockage requis

Stockage requis pour k blocks et une matrice de taille n

$$f(k, n) = (S\_block)^2 \times k + k \times 2$$

(émission, position block)

ou on peut aligner les blocks si on

Format classique:

$$f(k, n) = n^2$$

$$(S\_block)^2 \times k + k \times 2 \geq n^2$$

$$\Leftrightarrow k \geq \frac{n^2}{2 + (S\_block)^2}$$

Nombre de block max:  $\left(\frac{n}{S\_block}\right)^2$

$$\text{densité} = \frac{k}{\left(\frac{n}{S\_block}\right)^2}$$

$$\text{densité} \geq \frac{n^2}{n^2 \frac{(2 + (S\_block)^2)}{(S\_block)^2}} = \frac{1}{1 + \frac{2}{(S\_block)^2}}$$

En terme de stockage on est gagné par une densité plus petite que:

$$S\_block = 2 \Rightarrow \frac{2}{3}$$

$$S\_block = 3 \Rightarrow \frac{9}{11}$$

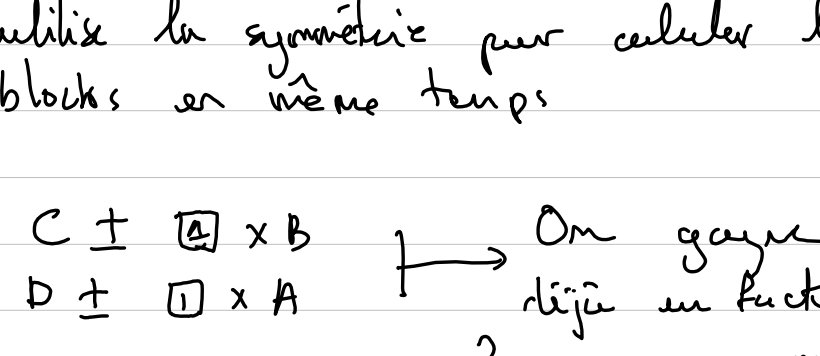
Ilème genre de densité en terme de calcul

Matrices sparse symétriques:

- Systèmes physiques
- bip de systèmes physiques

Multiplication vectorielle:

• Algo single thread: parcourir les blocks



Lorsque l'on parcourt on

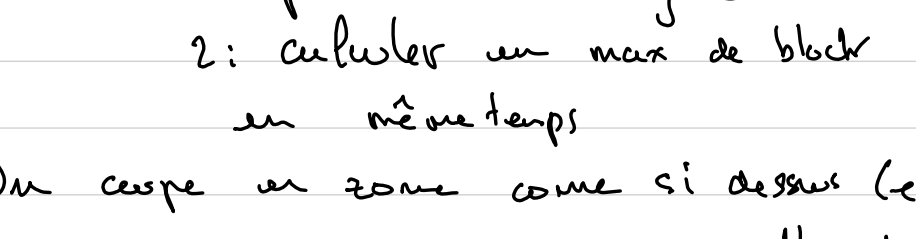
utilise la symétrie pour calculer les deux blocks en même temps

$$\begin{matrix} C \pm [0] \times B \\ D \pm [1] \times A \end{matrix} \rightarrow \begin{matrix} \text{On gagne} \\ \text{rien en facteur} \\ 2 \text{ car on ne} \\ \text{charge pas deux fois le} \\ \text{block dans le processeur} \end{matrix}$$

Les chargements mémoire sont le facteur limitant et l'arithmétique est moins coûteuse

vent mieux 32 op 1 load que 16 op 1 load

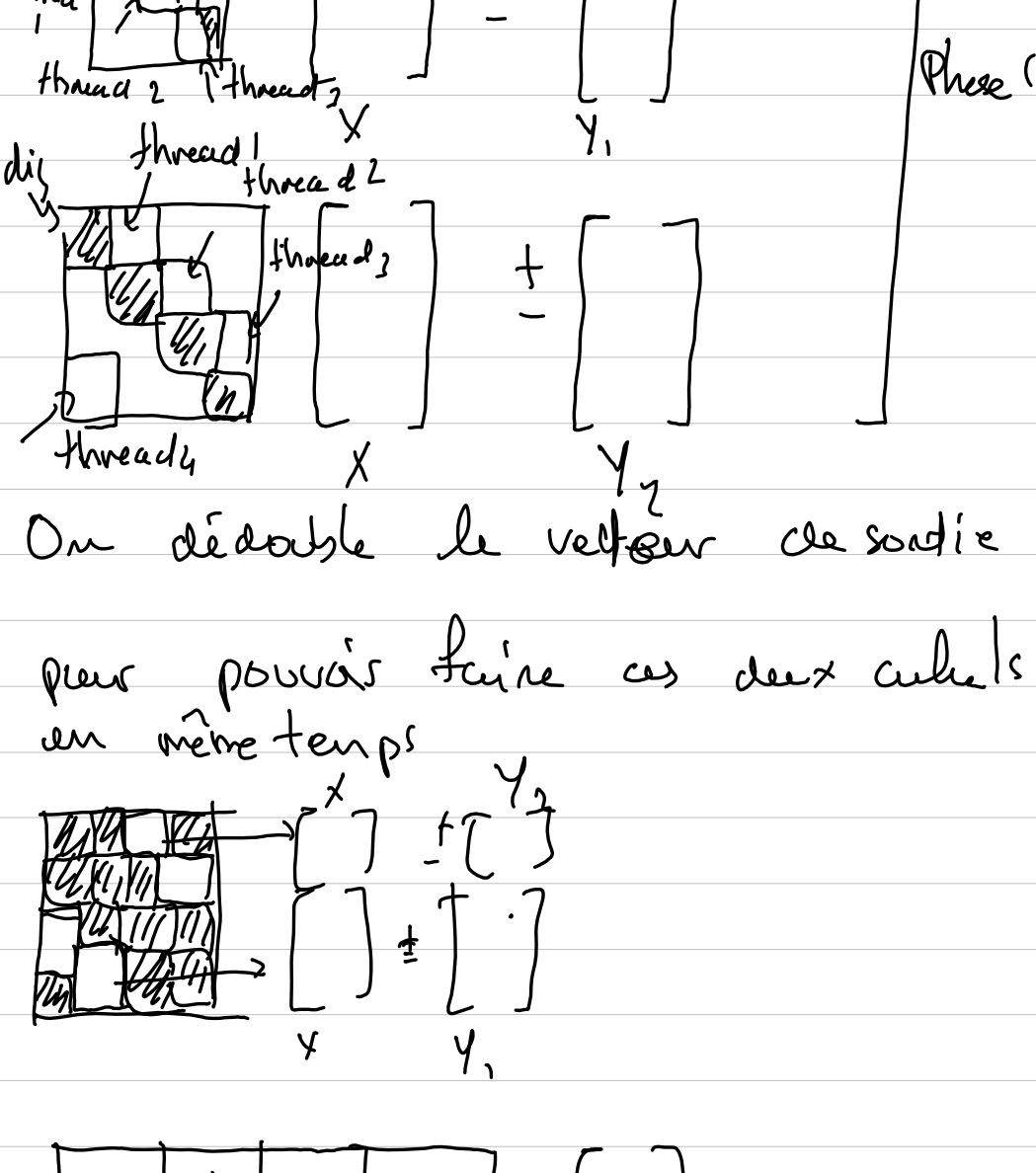
• Multithread sur contrôle d'accès mémoire thread-safe ex: 4 threads



On veut 1: profiter de la symétrie  
2: calculer un max de block en même temps

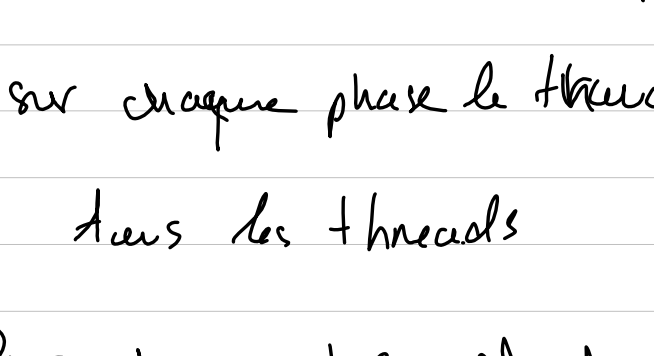
On coupe en zone comme si dessous (ex threads = 4)

On fait les calculs dans cet ordre



On dédouble le vecteur de sortie pour pouvoir faire ces deux calculs en même temps

$$\begin{bmatrix} x \\ y_1 \end{bmatrix} \pm \begin{bmatrix} y_2 \\ y_3 \end{bmatrix}$$



Avec cette méthode on peut régulariser sur chaque phase le travail permis dans les threads

Par chaque phase et chaque threads

On a une quantité de blocks à calculer

On stocke tout dans un tableau par phase

$$[16 | 16 | 16 | 16 | 16 | \dots] \text{ Phase 1}$$

thread 1 thread 2 thread 3 thread 4

afin d'optimiser le cache grâce à la localité et les positions des blocks

$$[2 | 2 | 2 | 2 | \dots]$$

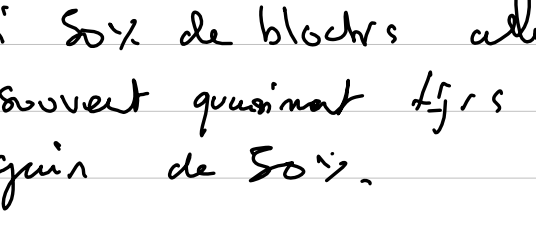
thread 1 thread 2 thread 3

(ix, ix)

On a également une multiplication optimisée pour calculer  $AX_1$  et  $AX_2$  en même temps pour minimiser les loads

• Améliorations possibles:

Profiter de l'alignement des blocks ex:



Ces deux blocks sont alignés et affecteront la même région de  $Y_1$ .  
Chaque modification de  $Y_1$  crée un load.  
On pourrait accumuler les résultats et ne modifier qu'une fois qu'un block est aligné à l'appareil.  
(Alignement sur X pour un block implique alignement selon Y pour son symétrique)  
Si 50% de blocks alignés. Ceci peut arriver souvent quasiment 1/3 de 50%.