

DE LA RECHERCHE À L'INDUSTRIE



www.cea.fr

PTC-INCOME

*Intégration de lois de
Comportement Mécanique sur
GPU*

Salem Khellal

6 Novembre 2023

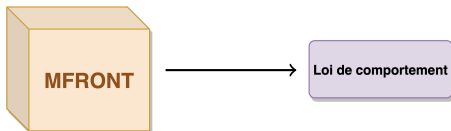


Maison
de la
Simulation

Loi de comportement

Loi de comportement

Relation mathématique qui décrit comment un matériau solide réagit aux charges, déformations et contraintes.



Loi de comportement

Relation mathématique qui décrit comment un matériau solide réagit aux charges, déformations et contraintes.

MFRONT

Générateur de code simplifiant l'écriture de propriétés matériaux, lois de comportement mécanique et modèles physico-chimiques.



Loi de comportement

Relation mathématique qui décrit comment un matériau solide réagit aux charges, déformations et contraintes.

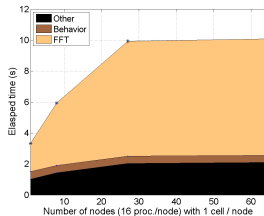
MFRONT

Générateur de code simplifiant l'écriture de propriétés matériaux, lois de comportement mécanique et modèles physico-chimiques.

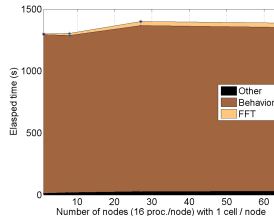
AMITEX

Solveur à base de Transformées de Fourier rapide (FFT).

Mise en contexte



Loi simple



Loi complexe



Objectifs

- ▶ Réaliser le portage sur GPU des lois de comportement dans MFRONT avec différents Programming Models (CUDA / Kokkos / SYCL)
- ▶ Proposer une solution flexible et performante sur GPU pour toutes les lois de comportement.

Points d'intégration

- ▶ Calculs indépendants entre chaque point d'intégration (embarassingly parallel) → Utilisation de functors
- ▶ Possibilité de déséquilibre de charge entre les points d'intégration

Structures de données avec MGIS/MFRONT

- ▶ Transferts mémoires CPU/GPU
- ▶ Gérer les structures propres à MGIS et à MFRONT dans le portage GPU
- ▶ Passage AOS (Array of Structs) vers SOA (Struct of Arrays) → rendre les accès mémoires coalescents

La loi de Hooke

Calcul du tenseur des contraintes σ à partir du tenseur des déformations ϵ dans chaque point d'intégration avec la relation suivante (cas linéaire):

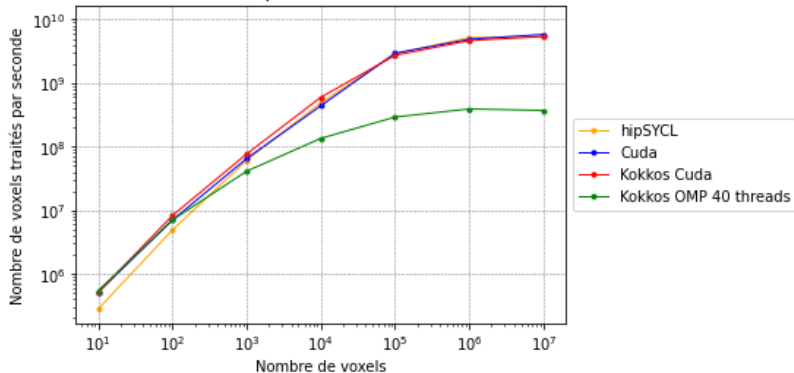
$$\underline{\sigma} = \lambda \text{trace}(\underline{\epsilon}) \mathbf{I} + 2\mu \underline{\epsilon}$$

- ▶ λ est le coefficient de Lamé
- ▶ μ est le module de cisaillement

Contributions

Performances des kernels pour la loi élasticit 

Performances des diff rents kernels pour le cas  lasticit  (lin aire) sur un GPU A100

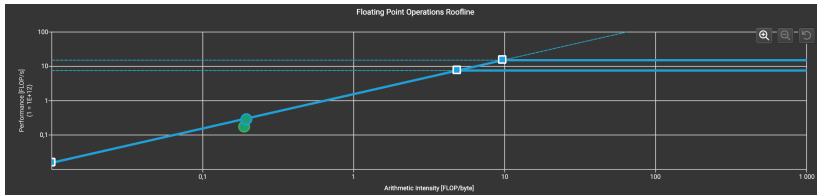


	Kokkos OMP (CPU)	Kokkos CUDA (GPU)	CUDA	hipSYCL
Kernel	0.03s	0.002s	0.002s	0.002s
H \rightarrow D	-	0.16s	0.16s	0.21s
D \rightarrow H	-	0.11s	0.15s	0.15s

- ▶ x10 speedup pour l'exécution du kernel,
- ▶ mais le transfert memoire est pénalisant sur GPU.

Contributions

Comparaison sur Roofline accès mémoire coalescents vs non coalescents



Avec des données coalescents, on est beaucoup plus proche du roof (performance maximale théorique).

Plasticité isotrope

Résolution d'une equation non-lineaire en utilisant un algorithme de Newton-Raphson scalaire en chaque point d'intégration:

$$\begin{cases} \underline{\epsilon}^{to} = \underline{\epsilon}^{el} + \underline{\epsilon}^p \\ \underline{\sigma} = \underline{D} : \underline{\epsilon}^{el} \\ \underline{\dot{\epsilon}}^p = \dot{p} \underline{n} \\ f(\sigma_{eq}, p) \dot{p} = 0, \text{ avec } f(\sigma_{eq}, p) \leq 0 \end{cases}$$

- ▶ $\underline{\epsilon}^{to}, \underline{\epsilon}^{el}, \underline{\epsilon}^p$ sont respectivement les déformations totale, élastique et plastique
- ▶ $\underline{n} = \frac{3}{2} \frac{\underline{s}}{\sigma_{eq}}$ est la direction d'écoulement
- ▶ \underline{s} est le déviateur des contraintes
- ▶ σ_{eq} est la norme de Von Mises.

Contributions

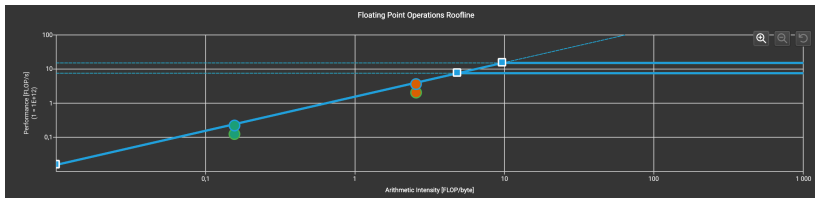
Comparaison de performances pour la loi plasticité avec 10^7 voxels

	Kokkos OMP (CPU)	Kokkos CUDA (GPU)	CUDA	hipSYCL
Kernel	0.09s	0.002s	0.002s	0.003s
H \rightarrow D		0.15s	0.18s	0.21s
D \rightarrow H		0.12s	0.17s	0.17s

► x45 speedup pour l'exécution du kernel.

Contributions

Comparaison sur Roofline accès mémoire coalescents vs non coalescents



Avec des données coalescents, on est beaucoup plus proche du roof comme dans le cas précédent .

La loi Norton

Résolution d'un système non-linéaire en utilisant un algorithme de Newton-Raphson vectoriel en chaque point d'intégration:

$$\begin{cases} \underline{\epsilon}^{to} = \underline{\epsilon}^{el} + \underline{\epsilon}^{vis} \\ \underline{\sigma} = \underline{D} : \underline{\epsilon}^{el} \\ \dot{\underline{\epsilon}}^{vis} = \dot{\rho} \underline{n} \\ \dot{\rho} = A \sigma_{eq}^m \end{cases}$$

- ▶ $\underline{\epsilon}^{to}, \underline{\epsilon}^{el}, \underline{\epsilon}^{vis}$ sont respectivement les déformations totale, élastique et visqueuse
- ▶ $\underline{n} = \frac{3}{2} \frac{\underline{s}}{\sigma_{eq}}$ est la direction d'écoulement
- ▶ \underline{s} est le déviateur des contraintes
- ▶ σ_{eq} est la norme de Von Mises.

Contributions

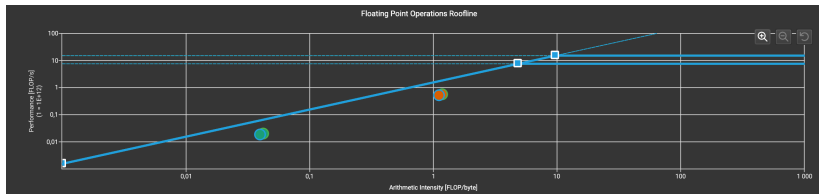
Comparaison de performances pour la loi Norton avec 10^7 voxels

	Kokkos OMP (CPU)	Kokkos CUDA (GPU)	CUDA	hipSYCL
Kernel	1.60s	0.29s	0.30s	0.33s
H → D		0.05s	0.05s	0.08s
D → H		0.05s	0.05s	0.07s

- ▶ x5 speedup pour l'exécution du kernel.
- ▶ Réduction de la quantité des données à copier à chaque itération → Overhead lié au transfert mémoire devient moins prononcé

Contributions

Comparaison sur Roofline accès mémoire coalescents vs non coalescents



Dans ce cas plus complexe, rendre les données coalescents ne permet pas de vraiment se rapprocher de la limite.

Newton vectoriel

Dans un algorithme de Newton vectoriel, on résout un système linéaire à chaque itération lorsque c'est possible:

$$F'(x_k)(x_{k+1} - x_k) = -F(x_k)$$

- ▶ $F'(x_k)$ la matrice Jacobienne à l'itération k
- ▶ x_k le vecteur contenant les variables internes à l'itération k

Hypothèse 1: Optimiser la décomposition LU ?

La Décomposition LU est l'étape qui prend le plus de temps. En testant un kernel qui réalise uniquement la décomposition LU, on a constaté avec une Roofline que ce n'est pas cette étape le problème.

Hypothèse 2: Grand nombre de registres par thread

Dans un algorithme de Newton vectoriel, on résout un système linéaire à chaque itération lorsque c'est possible:

$$F'(x_k)(x_{k+1} - x_k) = -F(x_k)$$

- ▶ $F'(x_k)$ la matrice Jacobienne à l'itération k nécessite le stockage de 12×12 double
- ▶ x_k le vecteur contenant les variables internes à l'itération k nécessite le stockage de 12 double

Dans l'algorithme actuel toutes ces données sont stockées dans la pile → Grand nombre de registres = problème sur GPU
Idée: déplacer le stockage de ces données vers la mémoire globale

Performances

- ▶ Loi simple (élasticité , plasticité): $\text{CPU} > \text{GPU}$
- ▶ Loi plus complexe (Norton): $\text{GPU} > \text{CPU}$

Perspectives

- ▶ Poursuivre l'optimisation du cas Norton (diminution du nombre de registres par thread)
- ▶ Tester les cas avec déséquilibre de charge en intégrant AMITEX
- ▶ Tester d'autres lois plus complexes

Merci pour votre contribution au PTC !

Thomas Helfer

Lionel Gélébart

Raphael Prat

Yushan Wang