

Programmation multiGPU OpenMP *versus* MPI

Gabriel Noaje, Michaël Krajecki, Christophe Jaillet

gabriel.noaje@univ-reims.fr

**Équipe SysCom, Laboratoire CReSTIC
Université de Reims Champagne-Ardenne, France**

17/02/2011

Sommaire

I. Introduction (GPU & multiGPU)

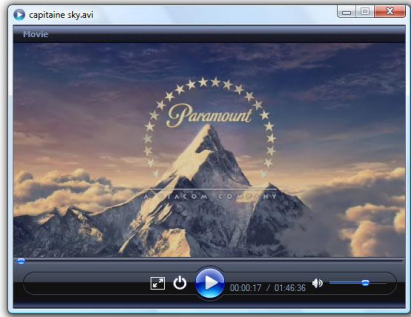
- *GPU*
- *GPU Computing*
- *Programmation de GPU*
- *Produit matriciel sur GPU (CUDA)*

II. Gestion des architectures multi-cartes

III. Résultats expérimentaux

IV. Conclusions

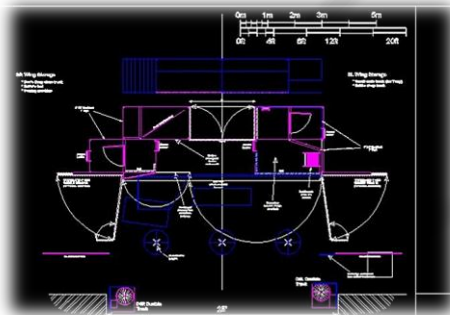
Carte graphique (Graphic Processing Unit - GPU)



lecture de films



Animations 3D



CAO (Conception Assistée par Ordinateur)



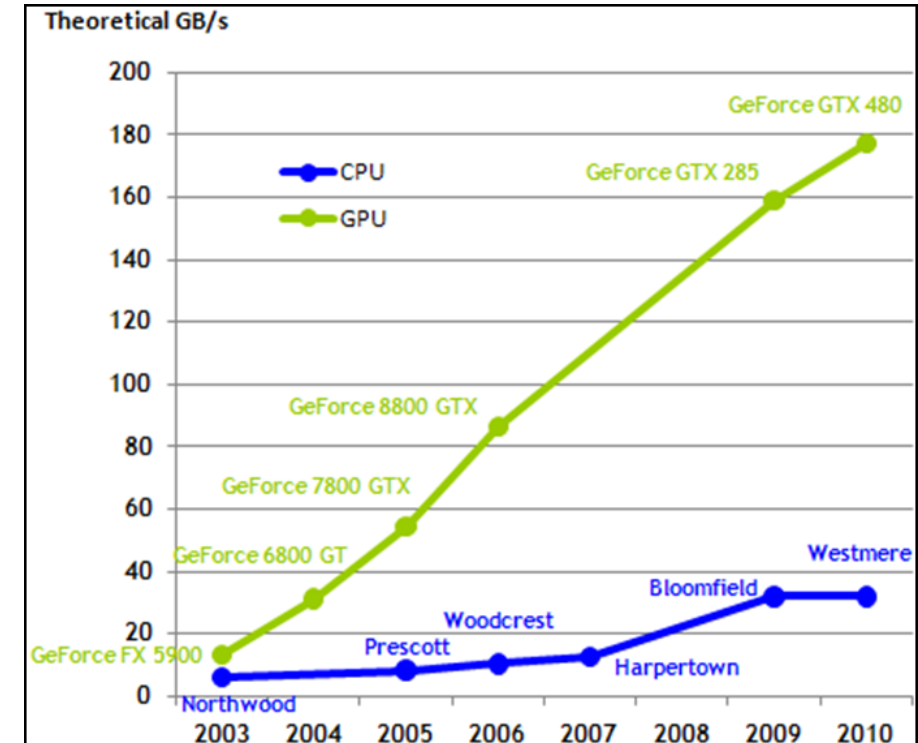
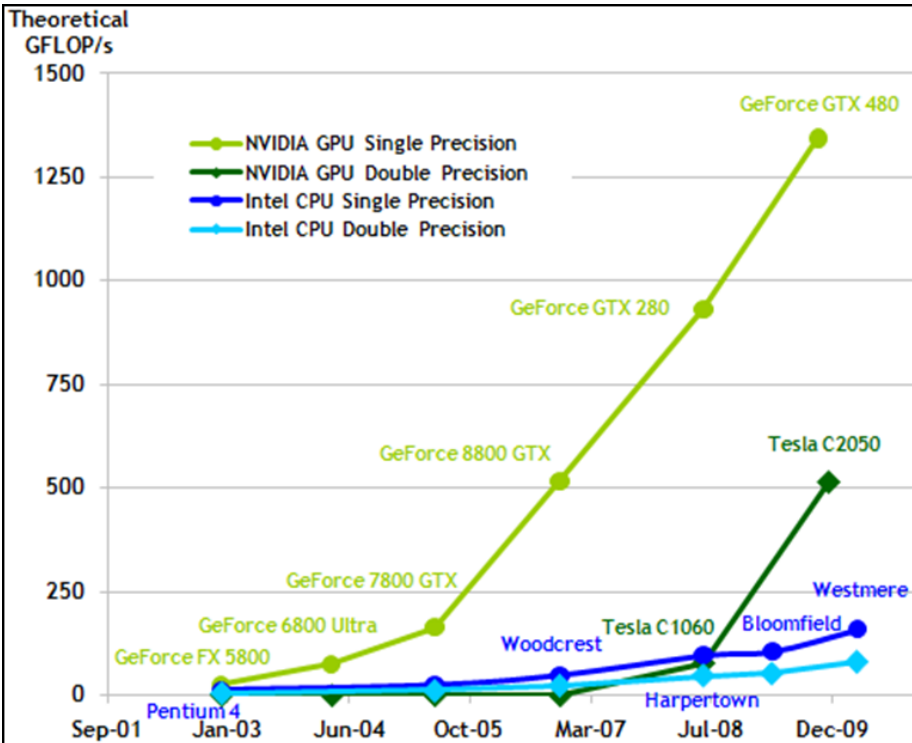
Jeux

GPU, pourquoi ?

PERFORMANCE !

puissance (GFLOPs)

mémoire (GB/s)

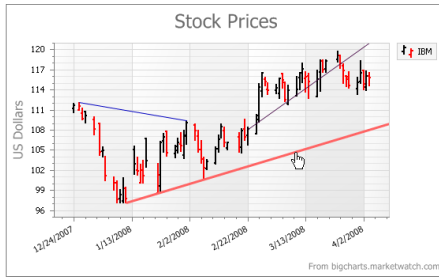


Credit: NVIDIA Programming Guide

- puissance de calcul supérieure aux processeurs classique (simple et double précision)
- accélérateurs matériels

- bande passante mémoire très importante (*mais juste à l'intérieur du GPU*)
- caractéristiques adaptées à **certain**s types de problèmes

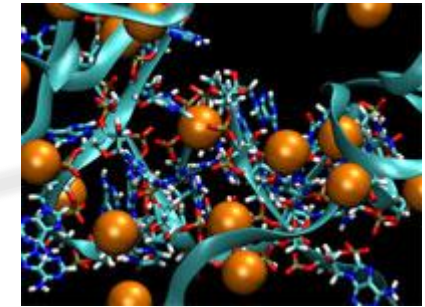
GPU Computing



Analyse financière



Simulation météo



Docking moléculaire

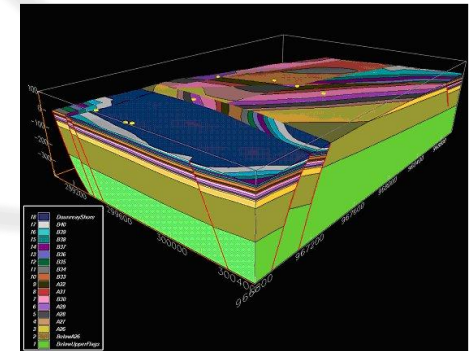
GPU



Imagerie médicale



Détection de logiciels malveillants



Modélisation géologique

Programmation de GPU

Initialement – utilisation d’appels graphiques :

- Cg de NVIDIA
- HLSL de Microsoft

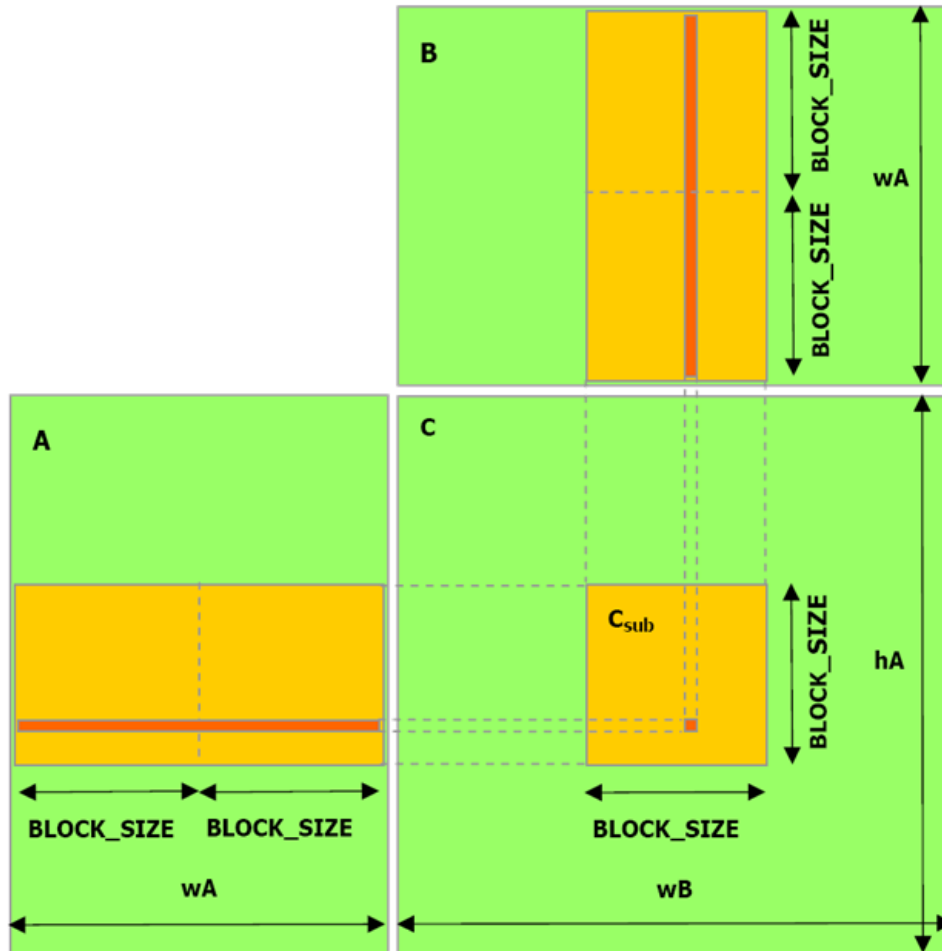
Aujourd’hui – écosystème de programmation de GPU :

- ATI Stream de AMD
- CUDA de NVIDIA
- OpenCL de Khronos Group
- Direct Computing de Microsoft

CUDA (bloqué sur le matériel NVIDIA = *aucune portabilité !*)

- = “**C**ompute **U**nified **D**evice **A**rchitecture”
- extension du standard ANSI C
- plus facile à apprendre (comparé à Cg, HLSL)
- ouvre l’architecture et donne à l’utilisateur un accès complet aux ressources

Produit matriciel classique sur GPU (CUDA)



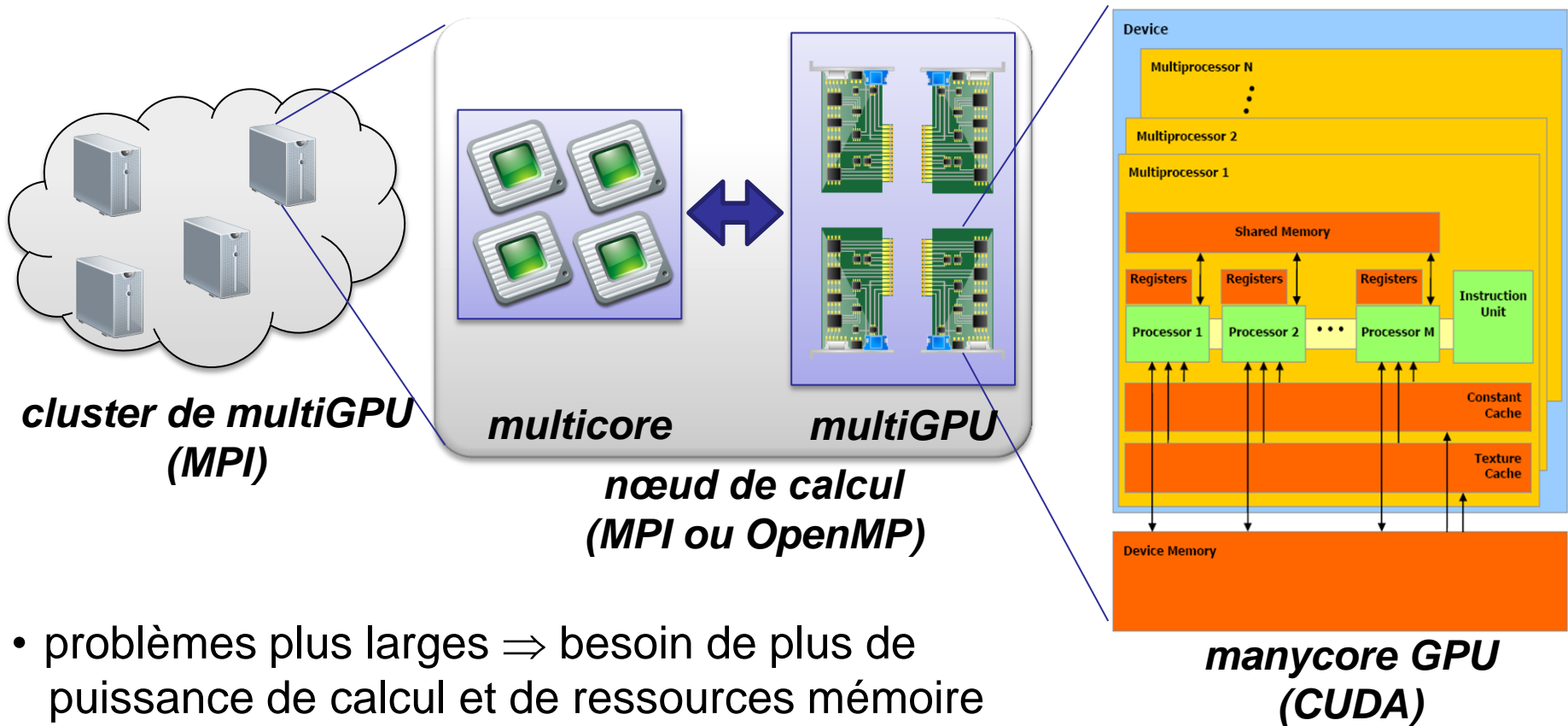
Produit matriciel $C = A \times B$

- correspond au modèle d'exécution ($C_{sub} \leftrightarrow$ block thread)
- utilise les différentes hiérarchies mémoire (mémoire globale, mémoire partagée)
- problème standard de calcul intensif \Rightarrow notre étude

Sommaire

- I. Introduction (GPU & multiGPU)
- II. Gestion des architectures multi-cartes
 - *Environnement multiGPU*
 - *Contrôle des GPU avec OpenMP ou MPI*
 - *Produit matriciel (strategie MPI versus OpenMP)*
- III. Résultats expérimentaux
- IV. Conclusions

Environnement multiGPU



- problèmes plus larges \Rightarrow besoin de plus de puissance de calcul et de ressources mémoire
- applications “simples” | multithread
multiprocesseur

- OK avec n'importe quelle API multithread ou avec passage de messages
- régions de calcul intensif déportées sur GPU
- multithreading \perp CUDA

Contrôle des GPU

	MPI	OpenMP
<i>environnement</i>	processus	thread
<i>mémoire</i>	distribuée	partagée
<i>échange des données</i>	explicite (passage des messages)	implicite (mémoire partagée)
<i>synchronisation</i>	explicite	explicite / implicite
<i>attachement à une carte graphique</i>	rang du processus	identifiant du thread

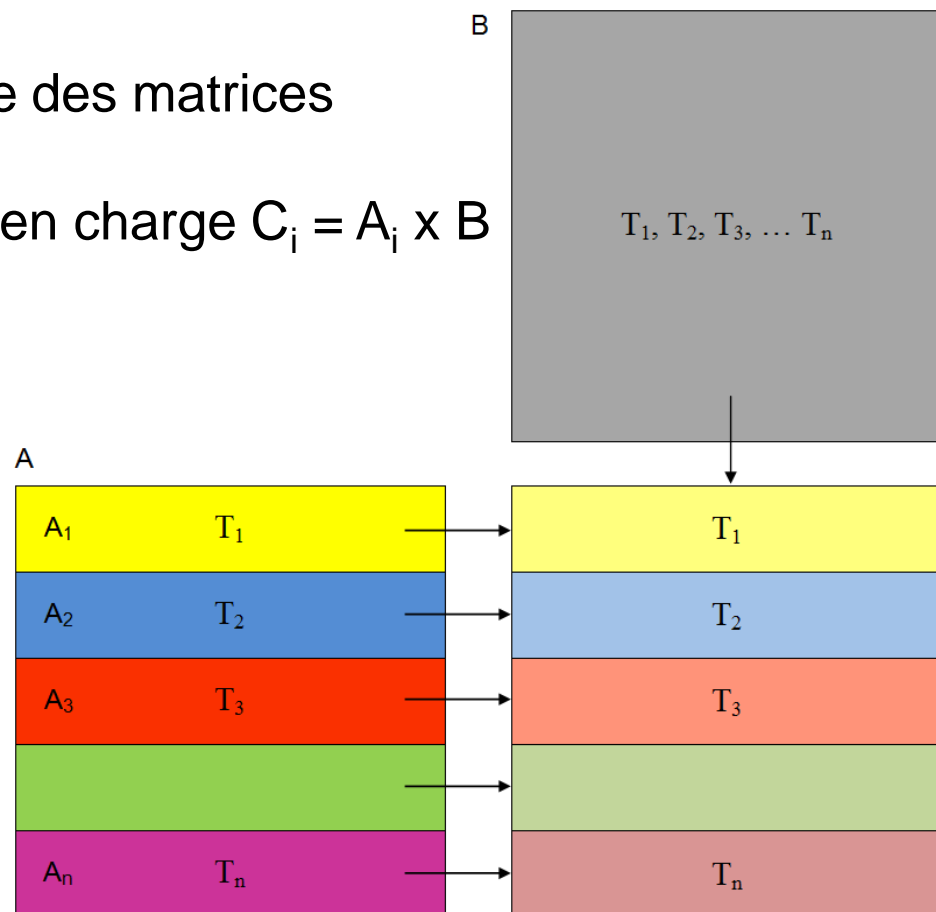
Produit matriciel (1)

Produit matriciel $C = A \times B$

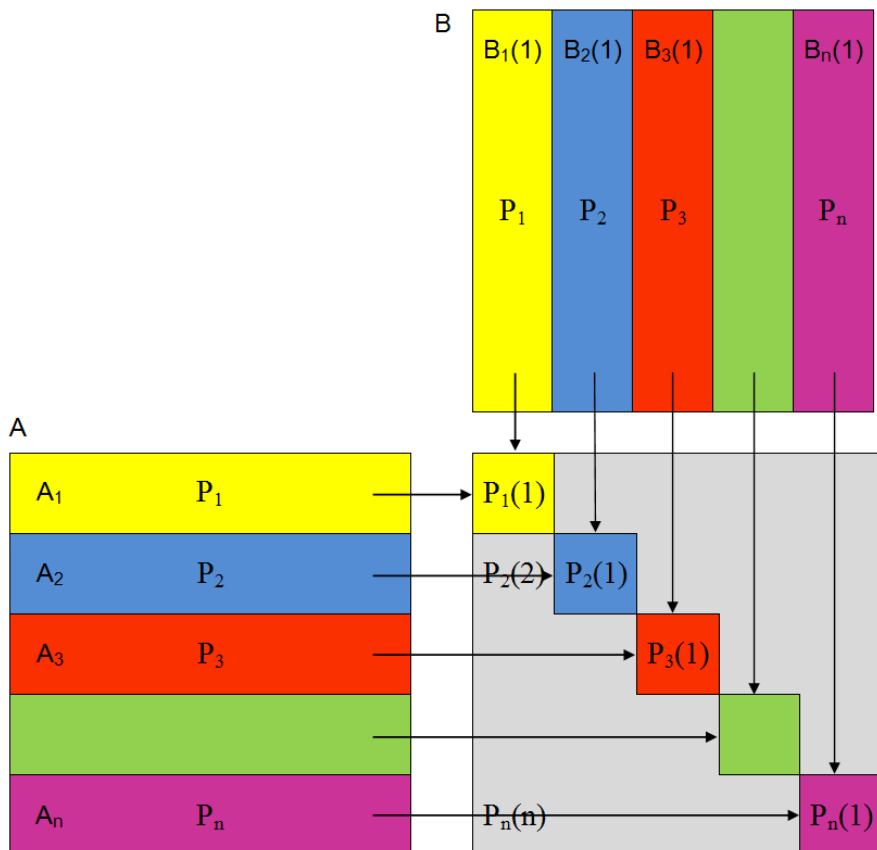
- matrices des grandes tailles
 - ⇒ diviser le travail ⇒ découpage des matrices
- partitionner la matrice A
 - ⇒ chaque unité de calcul prend en charge $C_i = A_i \times B$

Partitionnement OpenMP

- mémoire partagée
 - ⇒ plus de mémoire (nœud entier)
 - ⇒ matrice B non partitionnée
 - ⇒ pas d'effort de communication
- chaque thread s'attache à une carte
 - ⇒ allocation mémoire sur GPU
 - ⇒ transfert des blocs nécessaires (CPU → GPU)
 - ⇒ calcul (GPU)
 - ⇒ récupération des résultats (GPU → CPU)



Produit matriciel (2)



Partitionnement MPI

- mémoire distribuée
 - \Rightarrow matrice B découpée
 - \Rightarrow sous-matrices de B cyclées sur les processus
 - \Rightarrow celles de A restent fixes
- algorithme
 - \Rightarrow processus i attaché au GPU $_i$
 - \Rightarrow allocation mémoire (GPU) + transfert A_i
 - \Rightarrow [à chaque pas]
 - transfert B_j (CPU \rightarrow GPU)
 - $C_{ij} = A_i \times B_j$ (GPU)
 - transfert C_{ij} (GPU \rightarrow CPU)
- communications MPI
 - \Rightarrow buffering
 - stratégie du préchargement
 - recouvrement des communications MPI par des calculs GPU
 - coût de communication nul

Sommaire

I. Introduction (GPU & multiGPU)

II. Gestion des architectures multi-cartes

III. Résultats expérimentaux

- *Conditions expérimentales*
- *Comportement général (mesures de temps, irrégularités)*
- *Aspects spécifiques (pinned memory, recouvrement)*
- *Accélérations & limites des applications*

IV. Conclusions

Conditions expérimentales

Nombre du threads

- 1, 2 ou 4 threads/processus pour contrôler autant de cartes
4 threads/processus = utilisation complète du Tesla S1070
- 8 threads/processus \Rightarrow situations de compétition
deux threads/processus accèdent à une carte en même temps

Les temps mesurés

- **temps kernel** = temps nécessaire à calculer un bloc-ligne C_i
- **temps de transfert mémoire CUDA** = temps total nécessaire pour les transferts CPU-GPU / GPU-CPU
- **temps application** = temps total de l'application (initialisation de l'environnement d'exécution, initialisations des matrices, transferts mémoire, calculs)

Architecture matérielle



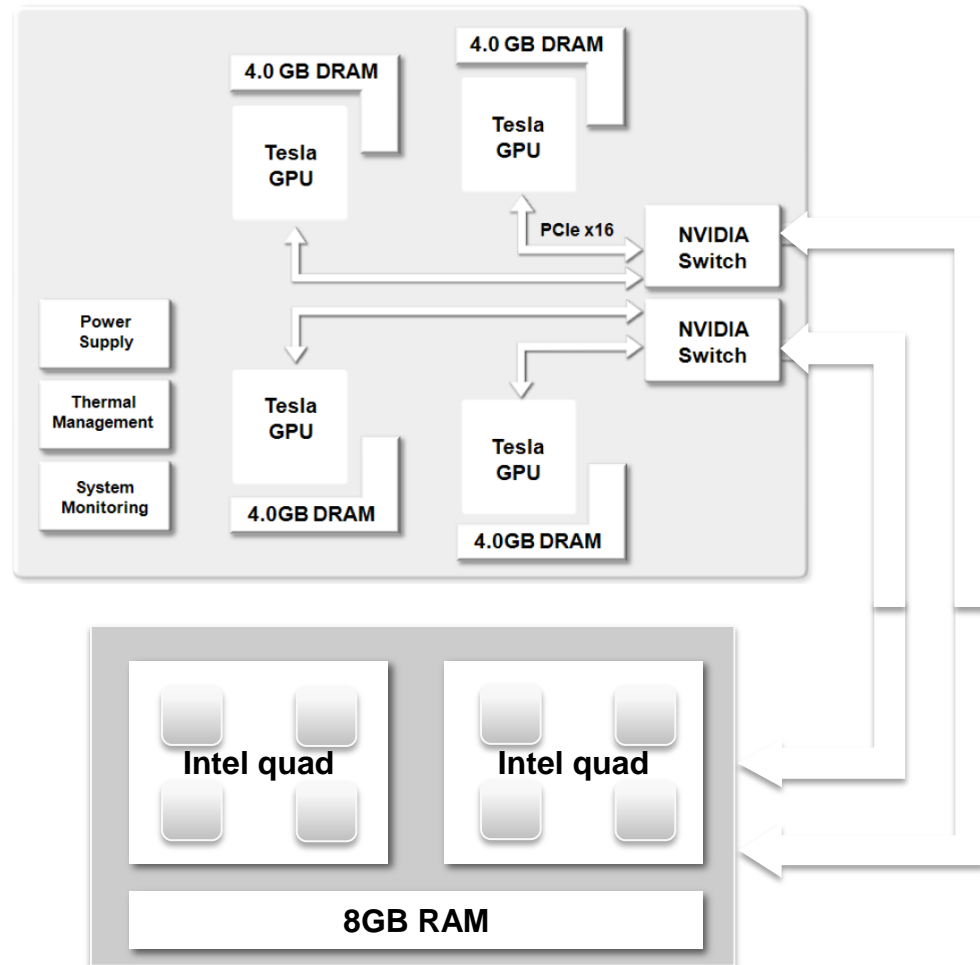
ROMEO – Centre de Calcul de Champagne-Ardenne

NVIDIA Tesla S1070

- 4 cartes graphiques “standard”
- 4Go DRAM / carte
- interconnectées par deux bridges PCI Express
- deux cartes d’interface

Serveur Bull NovaScale R425

- 2 quad-core Intel Xeon E5472
- 8Go mémoire vive DDR3



Comportement général

OpenMP (1, 2, 4 threads)

Temps kernel :

constamment divisé par 2
car le nombre de GPU double
(dim. sous-matrices / 2)

Temps application :

surcoût constant

- initialisation des environnements d'exécution
- allocations & initialisations de la mémoire hôte

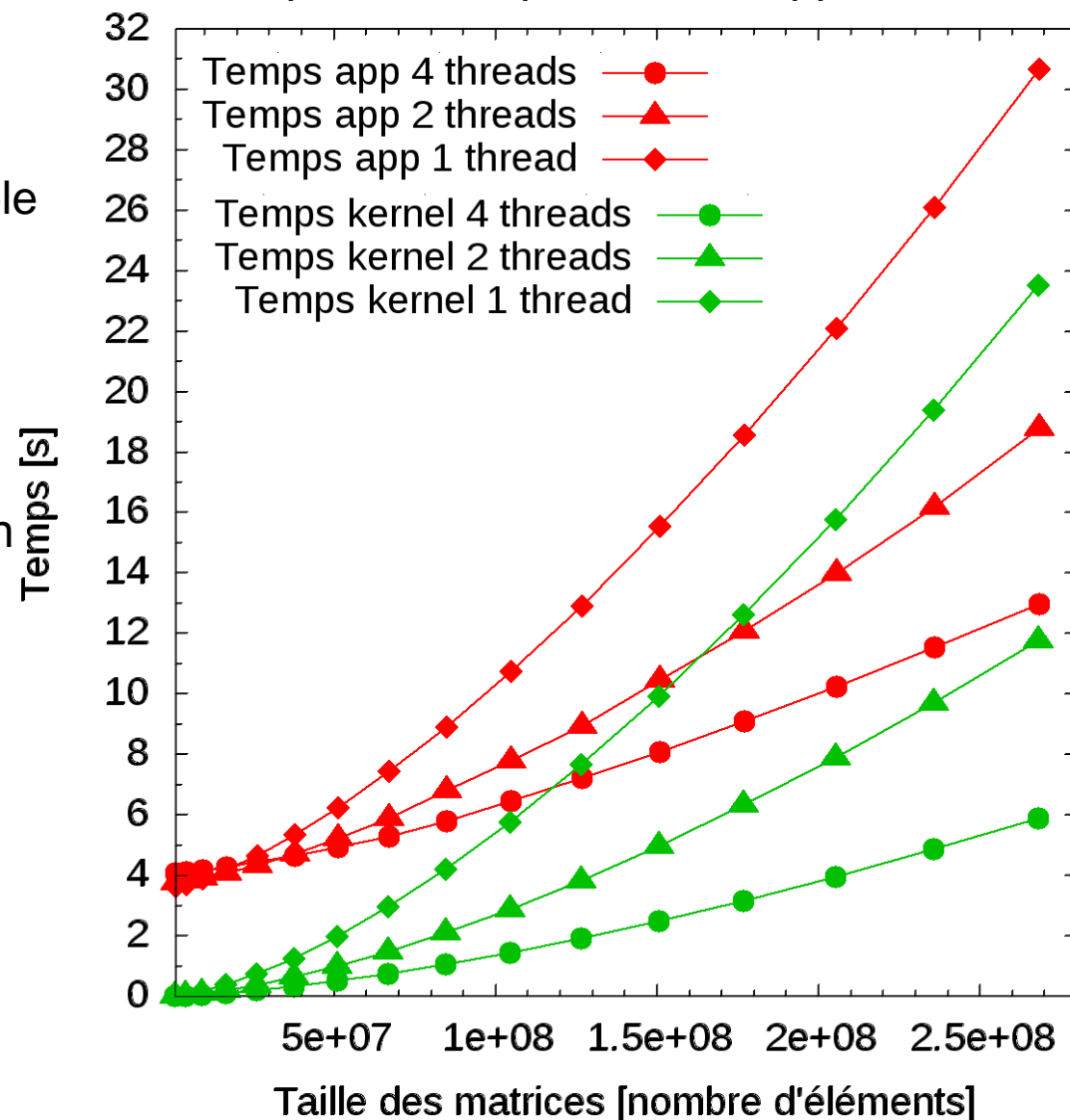
Temps transfert CPU-GPU :

- faible
- croissance constante

MPI

- même tendance
- surcoût plus important

OpenMP temps kernel & application



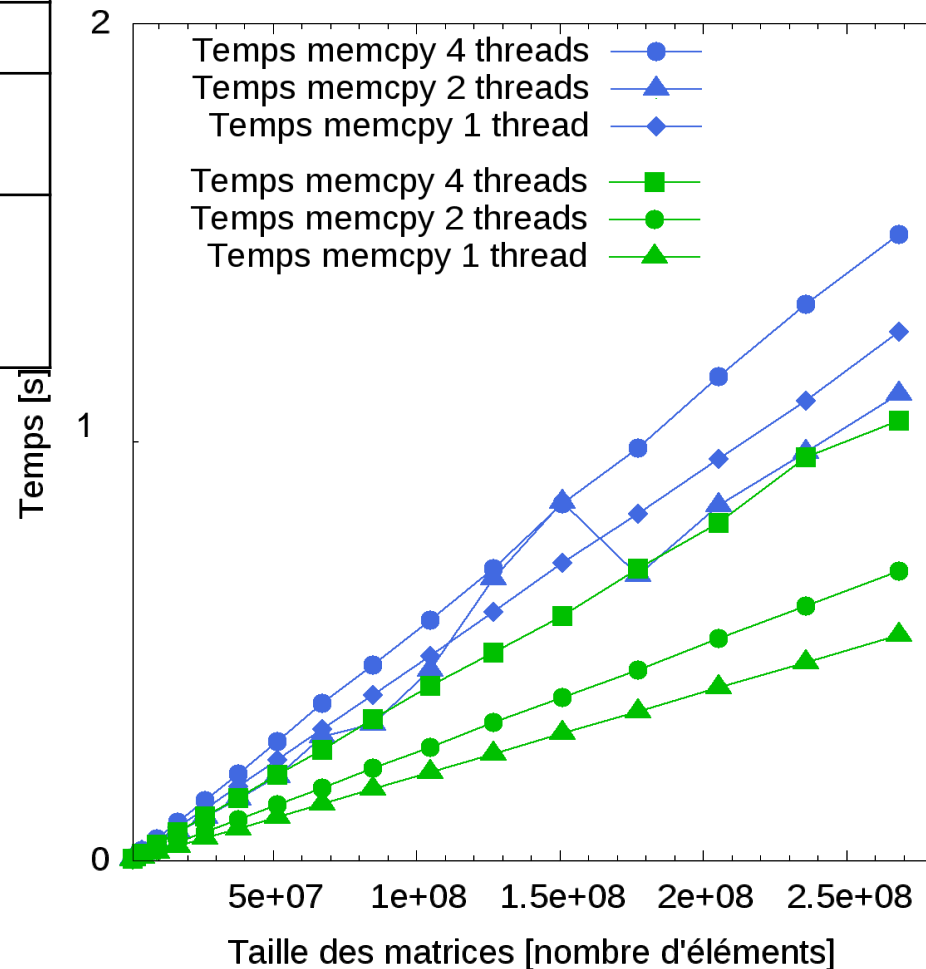
Types d'allocation mémoire

	malloc classique	pinned memory
<i>environ- nement</i>	POSIX	CUDA
<i>fonction</i>	malloc	cudaHostAlloc
<i>localisation</i>	adresse variable (pagination)	même adresse
<i>swap</i>	oui	non + permet des transferts zero-copy

Pinned memory

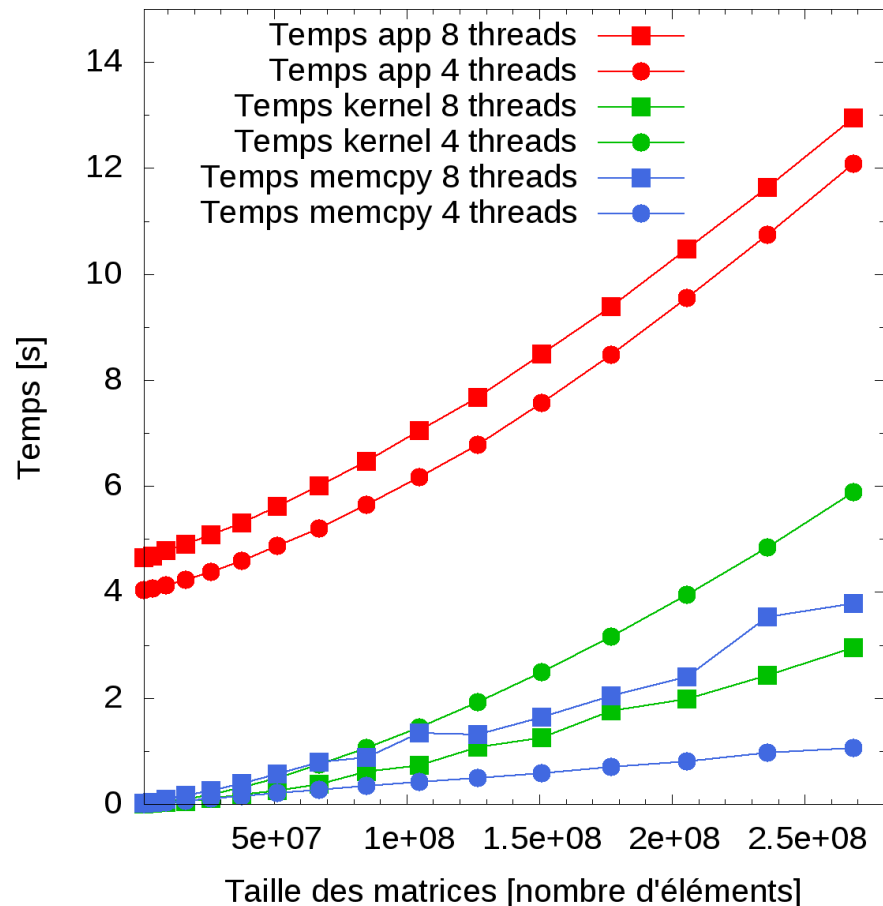
avantages	inconvénients
plus rapide	taille limitée
temps d'accès plus réguliers (pas de swap)	blocage du système en cas de sur- utilisation

OpenMP temps memcopy (*malloc* / *cudaHostAlloc*)

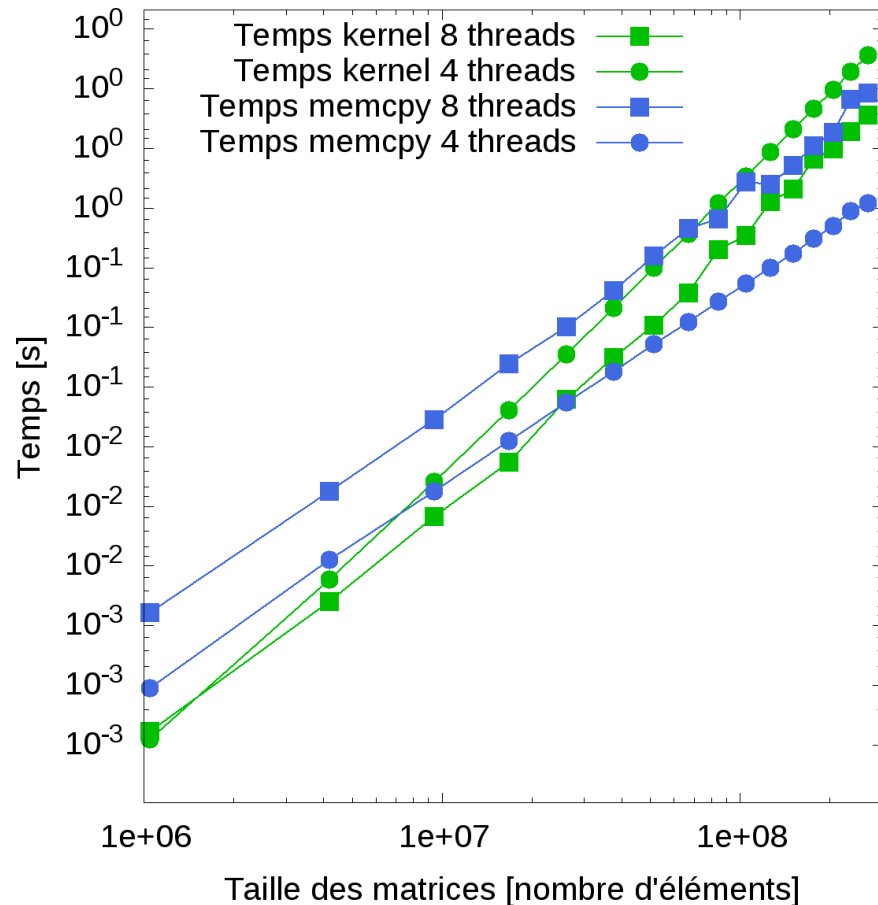


8 threads : concurrence pour ressources

Temps OpenMP – pinned memory



Temps OpenMP – pinned memory
(échelle logarithmique)

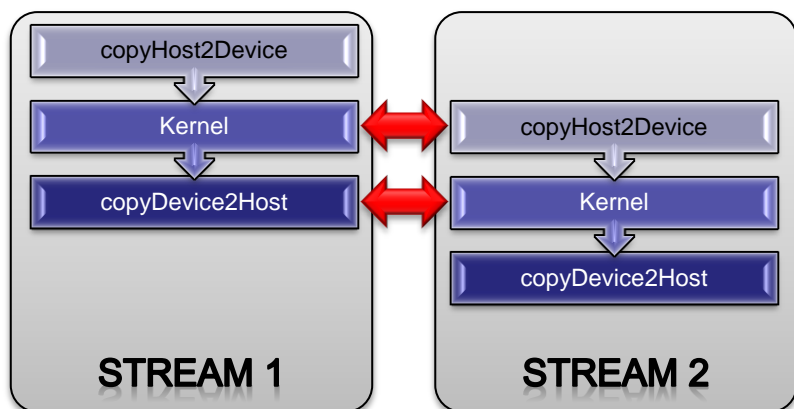


- temps kernel divisés par 2 (comme prévu)
- temps mémoire doublés à cause de la concurrence
- plus irrégulier (voir l'échelle logarithmique)

Recouvrement des transferts mémoire GPU

Principe :

- superposer les transferts CPU-GPU & l'exécution kernel
- *streams* (concept proche des threads)
- restriction : *pinned memory obligatoire*

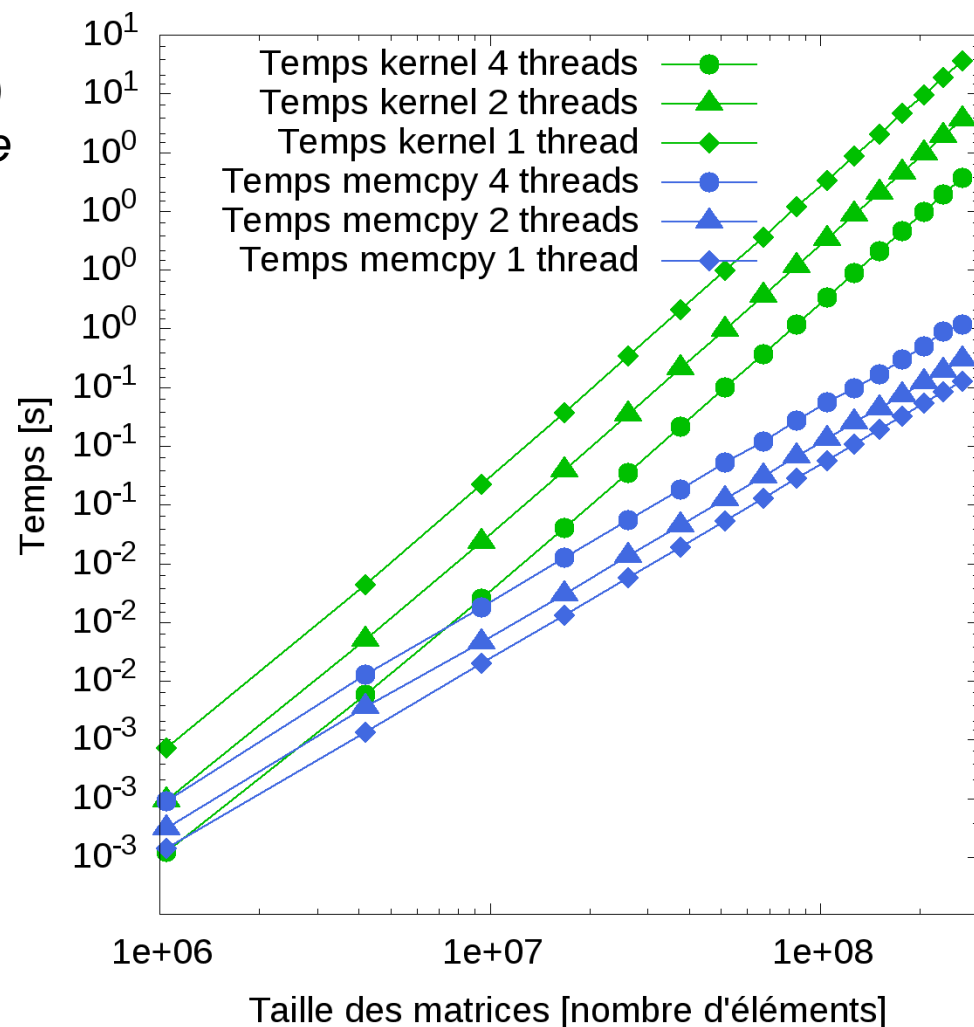


Dans notre cas :

- Temps kernel \gg temps transferts
 \Rightarrow aucune amélioration

- **Toujours envisager cette possibilité !**

Temps OpenMP – pinned memory (échelle logarithmique)



Accélérations & limites des applications

Limites des applications

- limites mémoire vive hôte
vs. mémoire device
- différentes stratégies de partitionnement
⇒ différentes limites mémoire
(implémentation MPI ⇒
2 matrices pour le buffering)
⇒ OpenMP est meilleur

Dim. matrice	MPI				OpenMP			
	1	2	4	8	1	2	4	8
18432	Green	Green	Green	Green	Green	Green	Green	Green
20480	Red	Green	Green	Green	Red	Green	Green	Green
22528	Red	Red	Red	Red	Red	Green	Green	Red
24576	Red	Red	Red	Red	Red	Red	Green	Red

Accélérations

- 4 threads/processus = utilisation optimale du Tesla S1070 (4 cartes)
- dim. matrices $\leq 18432 \times 18432$
(plus large ⇒ aucun résultat pour un seul thread/processus)
- efficacité : OpenMP = 64% vs. MPI = 50.5%
- OpenMP meilleur que MPI (augmentation en taille matrice ; + surcoût MPI)

Accélérations (4 threads/processus vs. 1 thread/processus)

Matrix dim	1024	2048	4096	6144	8192	10240	12288	14336	16384	18432
OpenMP	0.89	0.90	0.96	1.11	1.34	1.63	1.92	2.18	2.40	2.57
MPI	0.96	0.97	1.03	1.16	1.37	1.60	1.83	2.05	2.08	2.02

Efficacité

- limite, mais juste par la taille mémoire
- autrement, croissance linéaire


*Performances obtenues suivant
la taille des données*

Matrix dim. n (x 1024)	Operations n^3 (x 1024 ³)	Time t (s)	Measured GFlops
1	1	4.03	0.27
2	8	4.06	2.12
10	1000	6.16	174.31
12	1728	7.57	245.10
22	10648	23.85	479.38
24	13824	29.40	504.88

Sommaire

- I. Introduction (GPU & multiGPU)
- II. Gestion des architectures multi-cartes
- III. Résultats expérimentaux
- IV. Conclusions
 - *Discussion finale*
 - *Perspectives*

Conclusions & perspectives

- aujourd'hui la norme : 1 CPU - 1 GPU
- nouveau challenge : **multiGPU**
 - pas de communication directe entre les cartes graphiques 
 - management extérieur \Rightarrow OpenMP / MPI
 - OpenMP est le plus approprié
 - ❖ effort de programmation raisonnable
 - ❖ meilleures performances en termes de temps d'exécution et de scalabilité
 - stratégies d'optimisation additionnelles :
pinned memory, recouvrement

Perspectives

= grappes de nœuds multiGPU, avec trois niveaux :

- communication inter-nœuds \Rightarrow processus MPI (mémoire distribuée)
- gestion d'un nœud multiGPU \Rightarrow threads OpenMP (mémoire partagée)
- calcul, au sein des GPU

Programmation multiGPU OpenMP *versus* MPI

Questions ?

Contact: gabriel.noaje@univ-reims.fr