

Projet TOMOX

Accélération de la reconstruction tomographique
sur processeurs graphiques (GPUs)

Nicolas GAC, Ali Mohammad-Djafari, Alexandre Vabre,
Fanny Buyens, Eric Tordjeman, Karim Tadrir

Colloque GPU du GDR MI2B & CERIMED - Obernai
"Calcul intensif sur carte graphique pour l'imagerie moléculaire"

28/29 Mai 2009



UNIVERSITÉ
PARIS-SUD 11



digiteo
Research in information sciences and technologies



Projet TOMOX : Opération de Maturation Technico-Economique(OMTE)

Equipe opérationnelle

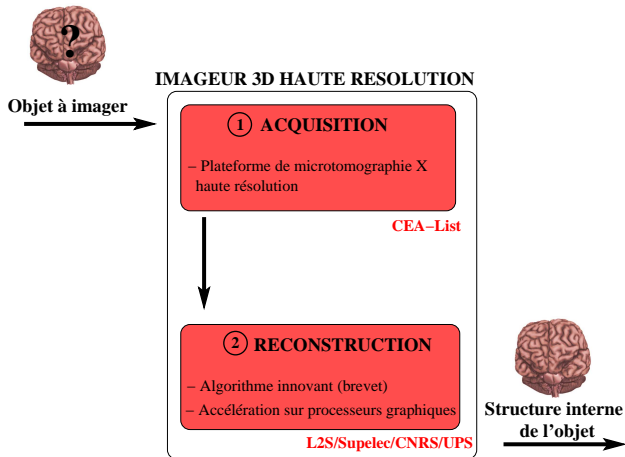
L2S	Ali Mohammad Djafari
CEA LIST	Alexandre Vabre, Fanny Buyens
Digiteo	Eric Tordjeman (Marketing) Nicolas Gac (Développement)
Supélec	Karim Tadrict (Juridique)



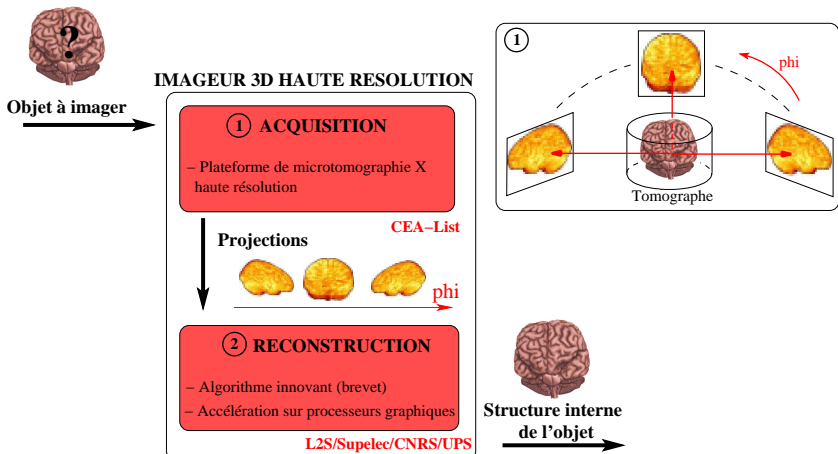
- 1 Objectif : accélérer la reconstruction tomographique
 - Projet TOMOX
 - Algorithme bayésien itératif du L2S/CEA
 - Accélération de la reconstruction
- 2 Tomographie sur GPU : Avant et Après CUDA
 - Avant CUDA ☹
 - Après CUDA ☺
- 3 Parallélisation des opérateurs de projection/rétroprojection
 - Localité des accès mémoire
 - Découpage en threads
 - Projecteur “voxel-driven” ou “ray-driven”
 - Temps GPU
- 4 Conclusion et perspectives

- 1 **Objectif : accélérer la reconstruction tomographique**
 - Projet TOMOX
 - Algorithme bayésien itératif du L2S/CEA
 - Accélération de la reconstruction
- 2 Tomographie sur GPU : Avant et Après CUDA
- 3 Parallélisation des opérateurs de projection/rétroprojection
- 4 Conclusion et perspectives

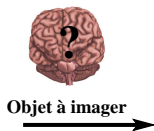
Projet TOMOX : Algorithmes innovants pour la reconstruction en nano-tomographie X



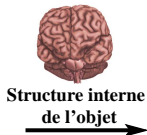
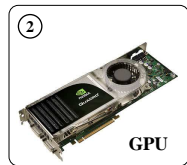
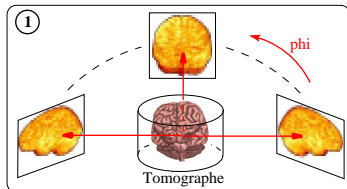
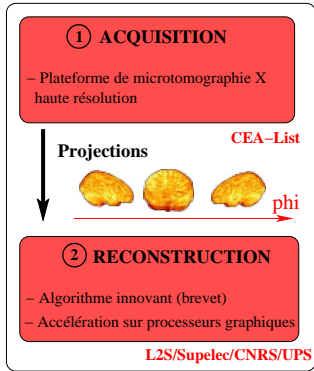
Projet TOMOX : Algorithmes innovants pour la reconstruction en nano-tomographie X



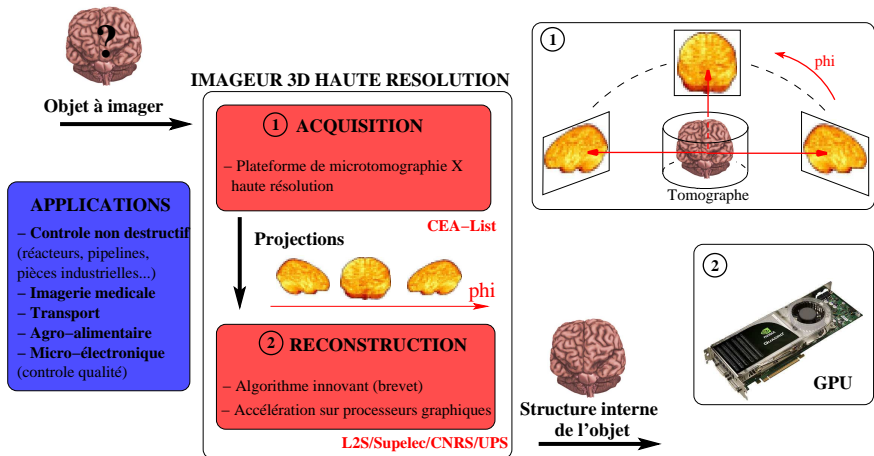
Projet TOMOX : Algorithmes innovants pour la reconstruction en nano-tomographie X



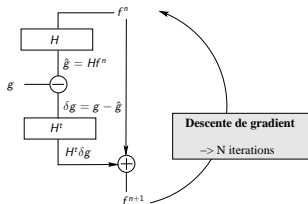
IMAGEUR 3D HAUTE RESOLUTION



Projet TOMOX : Algorithmes innovants pour la reconstruction en nano-tomographie X



SANS régularisation bayésienne



$$g = Hf + \epsilon$$

f : volume

g : données du tomographe

H : modèle d'acquisition

ϵ : bruit

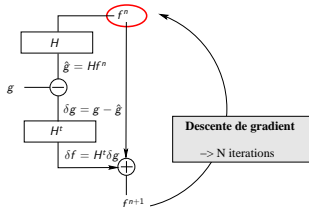
Descente de gradient

$$J(f) = \|g - Hf\|^2$$

$$f^{n+1} = f^n - \alpha \cdot \nabla J(f^n)$$

$$\nabla J(f) = -2 \cdot H^t (g - Hf)$$

SANS régularisation bayésienne



f^n : Estimée du volume

$$g = Hf + \epsilon$$

f : volume

g : données du tomographe

H : modèle d'acquisition

ϵ : bruit

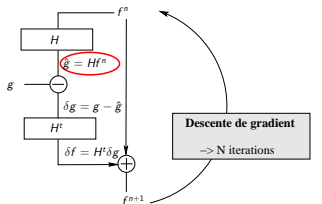
Descente de gradient

$$J(f) = \|g - Hf\|^2$$

$$f^{n+1} = f^n - \alpha \cdot \nabla J(f^n)$$

$$\nabla J(f) = -2 \cdot H^t(g - Hf)$$

SANS régularisation bayésienne



\hat{g} : Estimée des données

$$g = Hf + \epsilon$$

f : volume

g : données du tomographe

H : modèle d'acquisition

ϵ : bruit

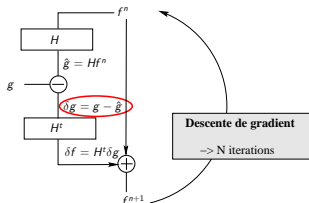
Descente de gradient

$$J(f) = \|g - Hf\|^2$$

$$f^{n+1} = f^n - \alpha \cdot \nabla J(f^n)$$

$$\nabla J(f) = -2 \cdot H^t(g - Hf)$$

SANS régularisation bayésienne



δg : Correction des données

$$g = Hf + \epsilon$$

f : volume

g : données du tomographe

H : modèle d'acquisition

ϵ : bruit

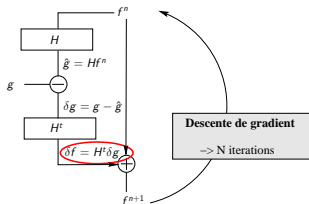
Descente de gradient

$$J(f) = \|g - Hf\|^2$$

$$f^{n+1} = f^n - \alpha \cdot \nabla J(f^n)$$

$$\nabla J(f) = -2 \cdot H^t(g - Hf)$$

SANS régularisation bayésienne



δf : Correction du volume

$$g = Hf + \epsilon$$

f : volume

g : données du tomographe

H : modèle d'acquisition

ϵ : bruit

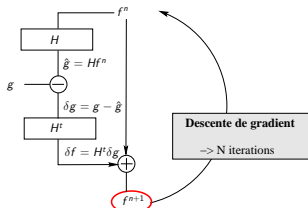
Descente de gradient

$$J(f) = \|g - Hf\|^2$$

$$f^{n+1} = f^n - \alpha \cdot \nabla J(f^n)$$

$$\nabla J(f) = -2 \cdot H^t(g - Hf)$$

SANS régularisation bayésienne



f^{n+1} : Nouvelle estimée du volume

$$g = Hf + \epsilon$$

f : volume

g : données du tomographe

H : modèle d'acquisition

ϵ : bruit

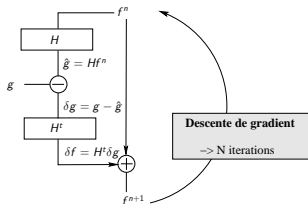
Descente de gradient

$$J(f) = \|g - Hf\|^2$$

$$f^{n+1} = f^n - \alpha \cdot \nabla J(f^n)$$

$$\nabla J(f) = -2 \cdot H^t(g - Hf)$$

AVEC régularisation bayésienne



$$g = Hf + \epsilon$$

f : volume

g : données du tomographe

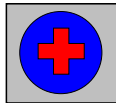
H : modèle d'acquisition

ϵ : bruit

$$\text{Modèle a priori : } f = \cup_k f_k$$

$$f_k = \{r : z(r) = k\}$$

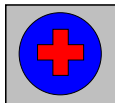
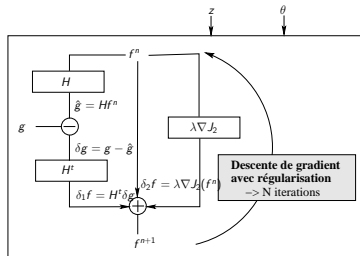
$$f(r)|z(r) = k \sim \mathcal{N}(m_k, \sigma_k^2)$$



Segmentation z

- $z = 0$: air
- $z = 1$: eau
- $z = 2$: metal

AVEC régularisation bayésienne



Segmentation z

- $z = 0$: air
- $z = 1$: eau
- $z = 2$: metal

$$g = Hf + \epsilon$$

f : volume

g : données du tomographe

H : modèle d'acquisition

ϵ : bruit

Modèle a priori : $f = \cup_k f_k$

$$f_k = \{r : z(r) = k\}$$

$$f(r) | z(r) = k \sim \mathcal{N}(m_k, \sigma_k^2)$$

Descente de gradient régularisée

$$J(f) = J_1(f) + J_2(f)$$

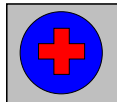
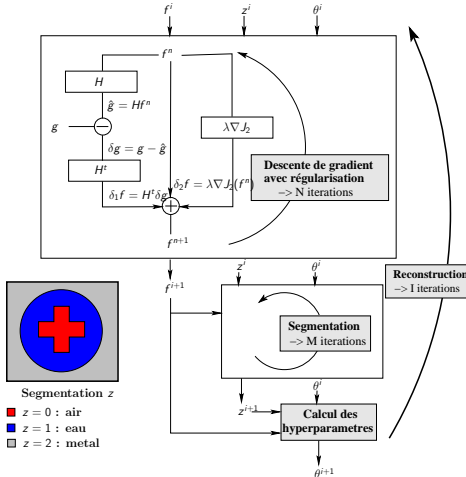
$$J_1(f) = \|g - Hf\|^2$$

$$J_2(f|z, \theta) = \lambda \sum_k \sum_{r \in \mathcal{R}_k} \sum_{r' \in \mathcal{V}(r)} \phi(\bar{f}(r) - \bar{f}(r'))$$

$$\bar{f}(r) = \frac{f(r) - m_k}{\sigma_k}$$

$$f^{n+1} = f^n - \alpha_1 \cdot \nabla J_1(f^n) - \alpha_2 \cdot \nabla J_2(f^n | z, \theta)$$

AVEC régularisation bayésienne



Segmentation z

- $z = 0$: air
- $z = 1$: eau
- $z = 2$: metal

$$g = Hf + \epsilon$$

f : volume

g : données du tomographe

H : modèle d'acquisition

ϵ : bruit

Modèle a priori : $f = \cup_k f_k$

$$f_k = \{r : z(r) = k\}$$

$$f(r) | z(r) = k \sim \mathcal{N}(m_k, \sigma_k^2)$$

Descente de gradient régularisée

$$J(f) = J_1(f) + J_2(f)$$

$$J_1(f) = \|g - Hf\|^2$$

$$J_2(f|z, \theta) = \lambda \sum_k \sum_{r \in \mathcal{R}_k} \sum_{r' \in \mathcal{V}(r)} \phi(\bar{f}(r) - \bar{f}(r'))$$

$$\bar{f}(r) = \frac{f(r) - m_k}{\sigma_k}$$

$$f^{n+1} = f^n - \alpha_1 \cdot \nabla J_1(f^n) - \alpha_2 \cdot \nabla J_2(f^n | z, \theta)$$

Nécessité d'accélérer la reconstruction

Un problème de plus en plus complexe

- Amélioration de la résolution spatiale des scanners
↳ Objectif : *Volume de $\simeq 2048^3$ voxels*
- Utilisation de méthodes itératives
↳ *10 à 100 itérations nécessaires*
- Reconstruction dynamique
↳ *plusieurs frames à reconstruire*

Temps de reconstruction insuffisant sur PCs

- Plusieurs heures voire jours de calcul

Calcul de Hf et $H^t\delta g$: choix de la méthode

① Calcul matriciel

⇒ lecture des coefficients h_{ij} dans la mémoire SDRAM

⚠ volume 2048^3 → matrice $H = 1$ To !

Calcul de Hf et $H^t\delta g$: choix de la méthode

① Calcul matriciel

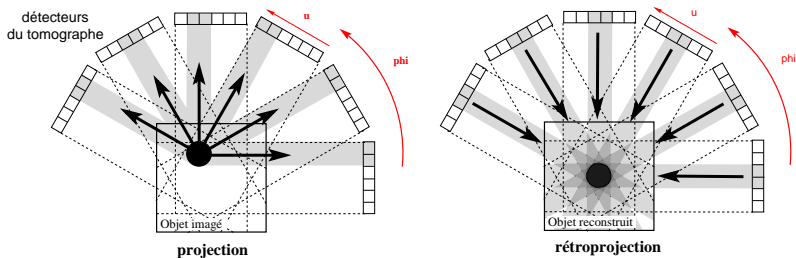
⇒ lecture des coefficients h_{ij} dans la mémoire SDRAM

⚠ volume 2048^3 → matrice $H = 1$ To !

② Opérateurs géométriques

⇒ calcul en ligne des coefficients h_{ij}

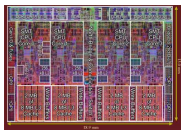
Paire de projection/rétroprojection en tomographie à émission (géométrie parallèle)



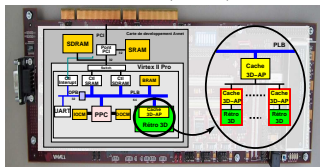
Calcul de Hf et $H^t \delta g$: choix du matériel

High Performance Computing (HPC)

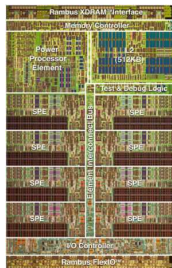
- Parallélisation sur machines multi-processeurs
↳ Efficace sur machine à mémoire distribuée
- Noeuds de calculs performants
↳ processeurs multi-core, many-core ou FPGA/ASIC



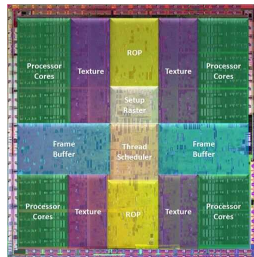
Intel Nehalem (4 coeurs)



SoPc (prototype)

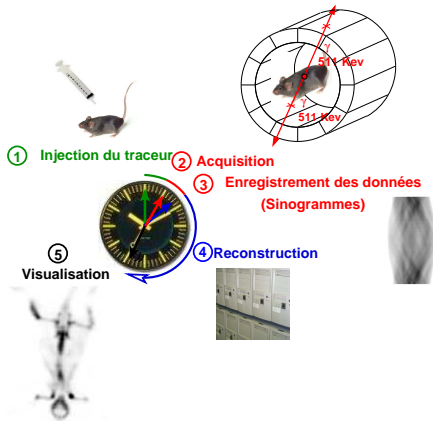


IBM Cell (8+1 coeurs)

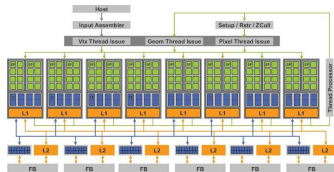


Nvidia GTX 200 (240 coeurs)

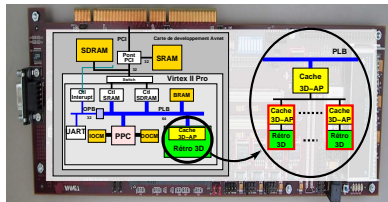
Thèse : "Adéquation Algorithme Architecture pour la reconstruction 3D en imagerie médicale TEP" (Gipsa-lab, Grenoble-INP)



Tomographie TEP



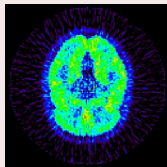
GPU (carte graphique)



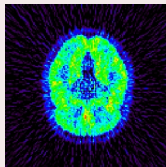
SoPC (prototype)

Exploration architecturale

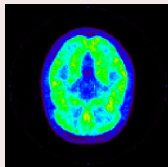
Conclusion de la thèse



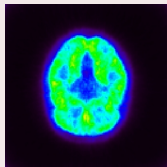
CPU
(non itératif)



FPGA
(non itératif)



CPU
(itératif)



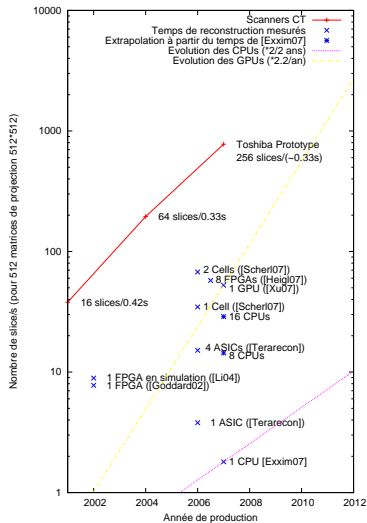
FPGA
(itératif)

Comparaison CPU/GPU/FPGA

	CPU	GPU	FPGA
Temps	3 ^{eme} (*4 P4)	1 ^{er} (*50 P4)	2 ^{eme} (*5 P4)
Efficacité	2 ^{eme} (7 C/op)	3 ^{ieme} (14 C/Op)	1 ^{er} (2 C/Op)

- GPU est l'accélérateur matériel le plus performant
- Validation de notre stratégie d'accès mémoire (cache 3D)

Vitesse d'acquisition // Vitesse de reconstruction



- 1 Objectif : accélérer la reconstruction tomographique
- 2 Tomographie sur GPU : Avant et Après CUDA
 - Avant CUDA ☹️
 - Après CUDA 😊
- 3 Parallélisation des opérateurs de projection/rétroprojection
- 4 Conclusion et perspectives

Avant CUDA : pipeline graphique

Vertex
Shader

Transformation
géométrique

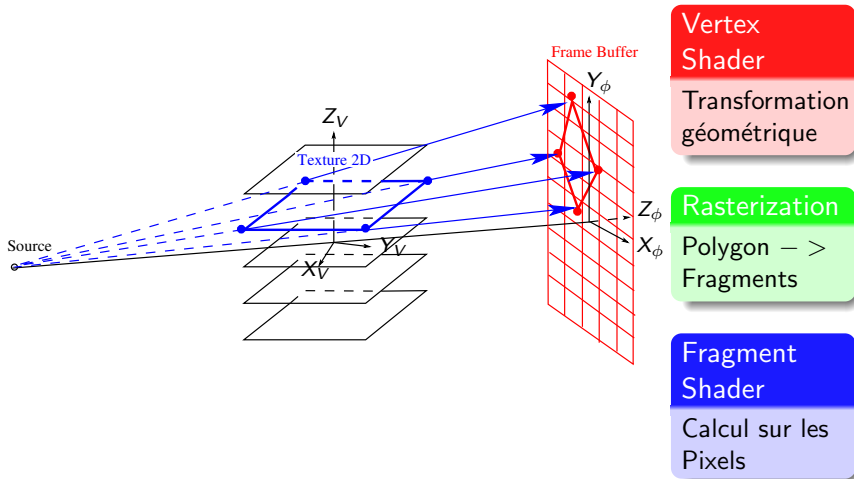
Rasterization

Polygon – >
Fragments

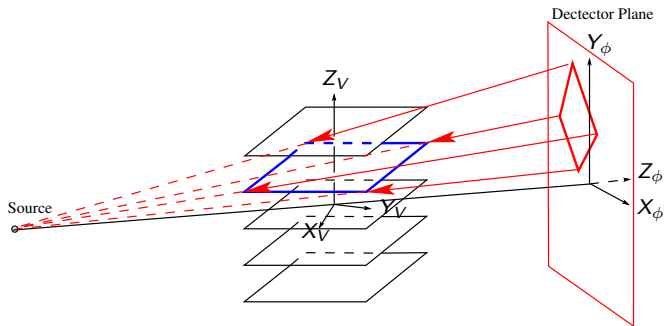
Fragment
Shader

Calcul sur les
Pixels

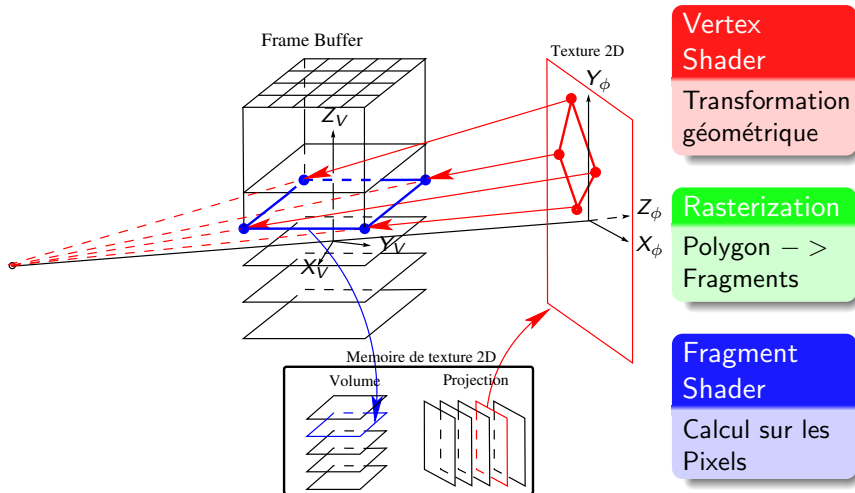
Projection Cone Beam 3D avec le Pipeline Graphique



Rétroprojection Cone Beam 3D

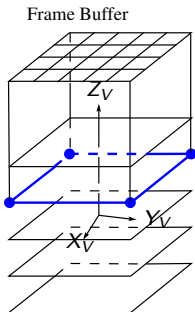


Rétroprojection Cone Beam 3D avec le Pipeline Graphique

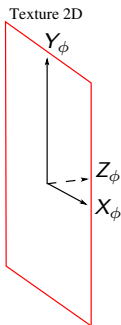
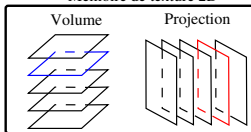


Rétroprojection Cone Beam 3D avec le Pipeline Graphique

Source
o



Memoire de texture 2D



Projection des
Vertex

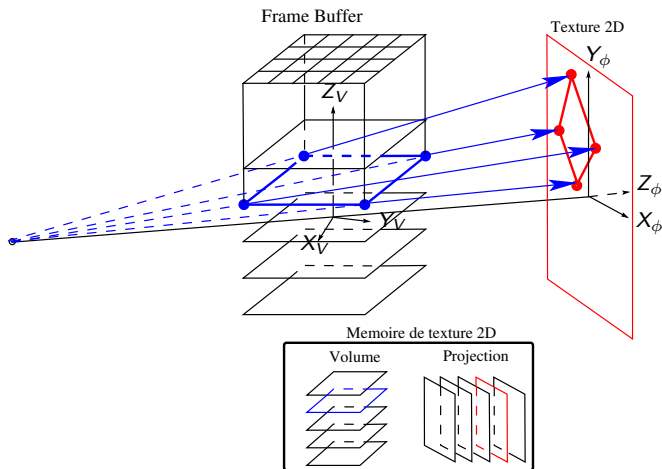
Projection des
Fragments

Indice de
Texture

Echantillonnage
de texture

Accumulation

Rétroprojection Cone Beam 3D avec le Pipeline Graphique



Projection des
Vertex

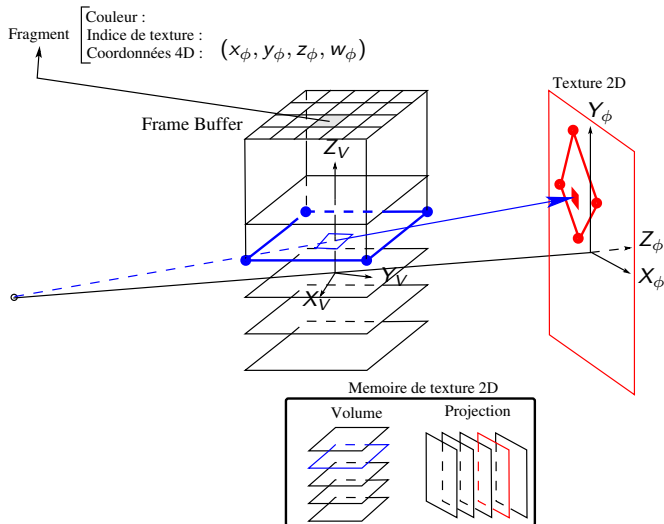
Projection des
Fragments

Indice de
Texture

Echantillonnage
de texture

Accumulation

Rétroprojection Cone Beam 3D avec le Pipeline Graphique



Projection des
Vertex

Projection des
Fragments

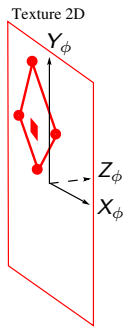
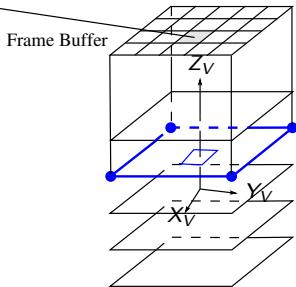
Indice de
Texture

Echantillonnage
de texture

Accumulation

Rétroprojection Cone Beam 3D avec le Pipeline Graphique

Fragment $\left[\begin{array}{l} \text{Couleur :} \\ \text{Indice de texture : } \left(x_{\phi} / w_{\phi}, y_{\phi} / w_{\phi} \right) \\ \text{Coordonnées 4D : } \left(X_{\phi}, y_{\phi}, z_{\phi}, w_{\phi} \right) \end{array} \right.$



Projection des Vertex

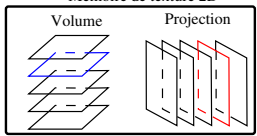
Projection des Fragments

Indice de Texture

Echantillonnage de texture

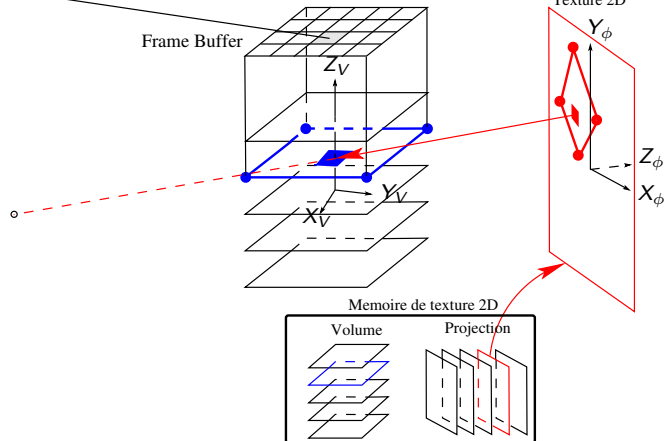
Accumulation

Memoire de texture 2D



Rétroprojection Cone Beam 3D avec le Pipeline Graphique

Fragment $\left[\begin{array}{l} \text{Couleur : } bin_{interp} \\ \text{Indice de texture : } (x_\phi / w_\phi, y_\phi / w_\phi) \\ \text{Coordonnées 4D : } (x_\phi, y_\phi, z_\phi, w_\phi) \end{array} \right.$



Projection des
Vertex

Projection des
Fragments

Indice de
Texture

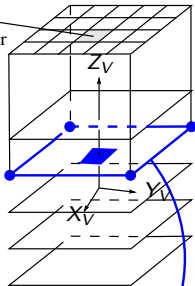
Echantillonnage
de texture

Accumulation

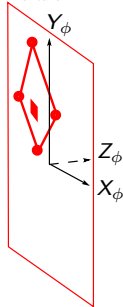
Rétroprojection Cone Beam 3D avec le Pipeline Graphique

Fragment $\left[\begin{array}{l} \text{Couleur : } bin_{interp} \\ \text{Indice de texture : } (x_\phi / w_\phi, y_\phi / w_\phi) \\ \text{Coordonnées 4D : } (x_\phi, y_\phi, z_\phi, w_\phi) \end{array} \right.$

Frame Buffer



Texture 2D



Projection des Vertex

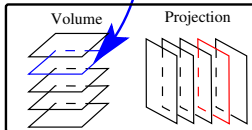
Projection des Fragments

Indice de Texture

Echantillonnage de texture

Accumulation

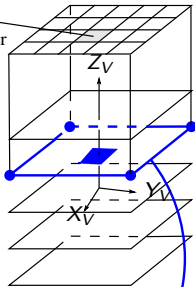
Memoire de texture 2D



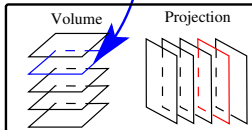
Rétroprojection Cone Beam 3D avec CUDA

Fragment $\left[\begin{array}{l} \text{Couleur : } bin_{interp} \\ \text{Indice de texture : } (x_\phi / w_\phi, y_\phi / w_\phi) \\ \text{Coordonnées 4D : } (x_\phi, y_\phi, z_\phi, w_\phi) \end{array} \right.$

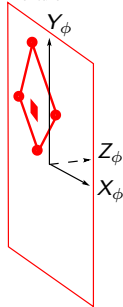
Frame Buffer



Memoire de texture 2D



Texture 2D



Projection des
Fragments

Indice de
Texture

Echantillonnage
de texture

Accumulation

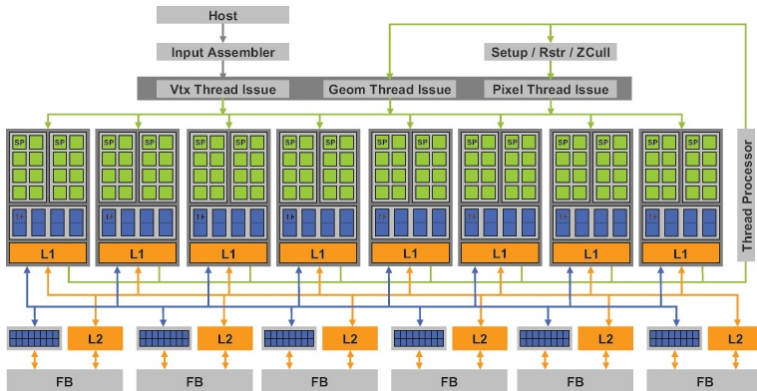
OpenGL (Pipeline Graphique) ou CUDA ?

Intérêt relatif de la programmation OpenGL vs CUDA


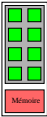
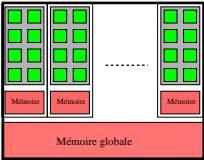
- Plus grande complexité de programmation en OpenGL
- Dans le meilleur des cas, accélération d'un facteur 3

Ref	Feldkamp en OpenGL/CUDA (sur 8800 GTX)
[Xu 07]	accélération d'un facteur 3
[Knaup 08]	accélération d'un facteur 1.2

Après CUDA : plein de threads !



Découpage en threads

	Matériel	Logiciel	Exécution	
	un Processeur Élémentaire (PE)	un thread	séquentielle	(a)
	un processeur multi-threads	un bloc de threads	parallèle (SIMT)	(b)
	une carte GPU (device)	une grille de threads; (kernel)	parallèle (MIMD) mémoire centralisée	(c)

Programmation GPU

1 Parallélisation de l'algorithme

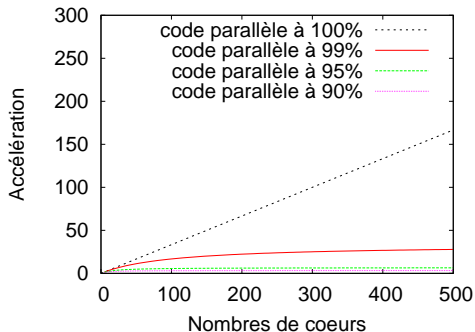
➤ nourrir en threads (plus ou moins indépendants) le GPU

n coeurs (1 Ghz)

vs

1 coeur (3 Ghz)

taux de parallélisation	accélération GTX 200 (240 coeurs)
100 %	80
99 %	24
95 %	6
90 %	3



Programmation GPU

① Parallélisation de l'algorithme

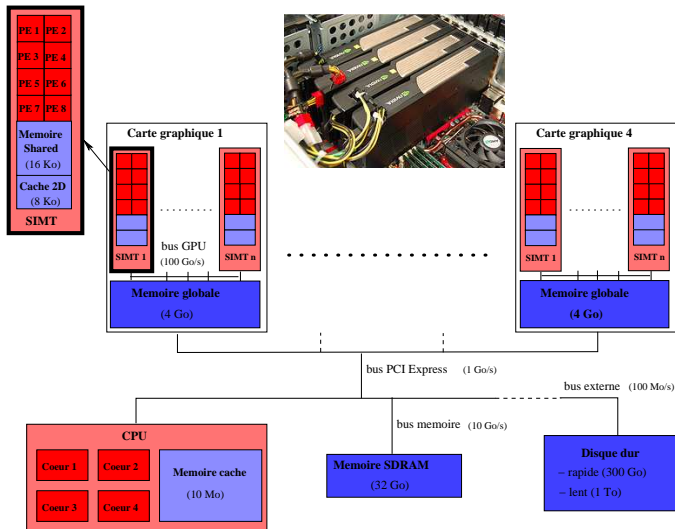
⇒ nourrir en threads (plus ou moins indépendants) le GPU

② Implémentation GPU


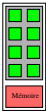
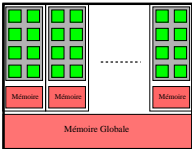
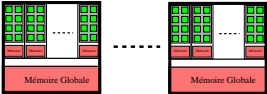
Selon l'intensité arithmétique du code (puissance de calcul exploitée / débit des données), l'exécution sera soit *memory bound* soit *computation bound* (ex : calcul X^k [Kirschenmann 08])

⇒ optimisation du code portera alors soit sur les **accès mémoire** ou soit sur la **complexité arithmétique**

Supercalculateur personnel



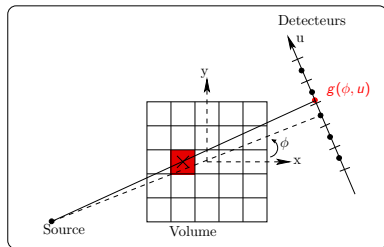
Découpage en threads et en grilles (kernels)

Matériel	Logiciel	Exécution	
 <p>un Processeur Élémentaire (PE)</p>	un thread	séquentielle	(a)
 <p>un processeur multi-threads</p>	un bloc de threads	parallèle (SIMT)	(b)
 <p>une carte GPU (device)</p>	une grille de threads; (kernel)	parallèle (MIMD) mémoire centralisée	(c)
 <p>PC multi-carte</p>	threads du PC hôte via librairie pthread (un thread CPU = un kernel GPU)	parallèle (MIMD) mémoire distribuée	(d)

- 1 Objectif : accélérer la reconstruction tomographique
- 2 Tomographie sur GPU : Avant et Après CUDA
- 3 **Parallélisation des opérateurs de projection/rétroprojection**
 - Localité des accès mémoire
 - Découpage en threads
 - Projecteur "voxel-driven" ou "ray-driven"
 - Temps GPU
- 4 Conclusion et perspectives

Rétroprojection 2D : algorithme

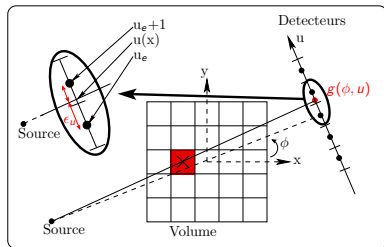
CALCUL DES COORDONNEES



```
for (xn, yn) in Volume do  
  for phi = 0 to phimax - 1 do  
    // Calcul des coordonnées  
    u(phi, xn, yn) = ...  
    // Accumulation  
    f*(xn, yn)+ = g(u, φ)  
  end for  
end for
```

Rétroprojection 2D : interpolation linéaire

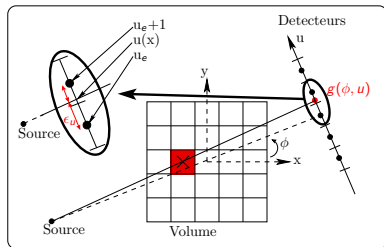
CALCUL DES COORDONNEES



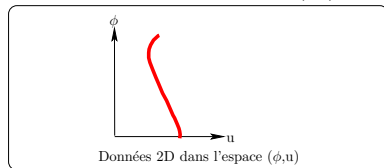
```
for (xn, yn) in Volume do
  for phi = 0 to phi_max - 1 do
    // Calcul des coordonnées
    u(phi, xn, yn) = ...
    // Interpolation linéaire
    g_interp = (1 - epsilon_u) * g(phi, u_e) +
              epsilon_u * g(phi, u_e + 1)
    // Accumulation
    f*(xn, yn) += g_interp
  end for
end for
```

Rétroprojection 2D : accès aux données dispersés

CALCUL DES COORDONNEES



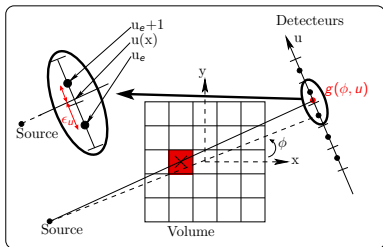
ACCES AUX DONNEES $g(\phi, u)$



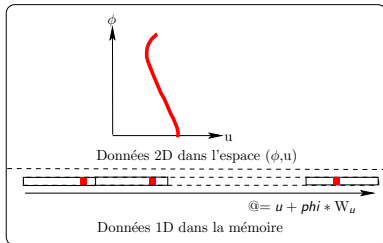
```
for (xn, yn) in Volume do
  for phi = 0 to phi_max - 1 do
    // Calcul des coordonnées
    u(phi, xn, yn) = ...
    // Interpolation linéaire
    g_interp = (1 - epsilon_u) * g(phi, u_e) +
              epsilon_u * g(phi, u_e + 1)
    // Accumulation
    f*(xn, yn) += g_interp
  end for
end for
```

Rétroprojection 2D : accès aux données dispersés

CALCUL DES COORDONNEES



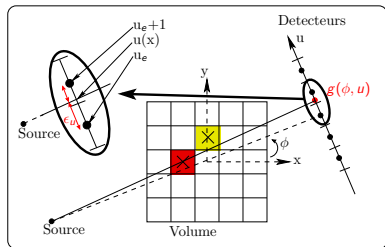
ACCES AUX DONNEES $g(\phi, u)$



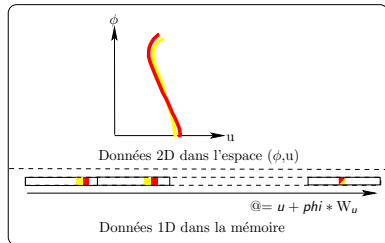
```
for (xn, yn) in Volume do
  for phi = 0 to phi_max - 1 do
    // Calcul des coordonnées
    u(phi, xn, yn) = ...
    // Interpolation linéaire
    g_interp = (1 - epsilon_u) * g(phi, u_e) +
              epsilon_u * g(phi, u_e + 1)
    // Accumulation
    f*(xn, yn) += g_interp
  end for
end for
```


Rétroprojection 2D : accès aux données dispersés

CALCUL DES COORDONNEES



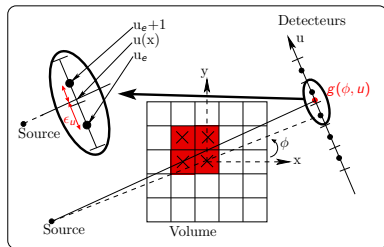
ACCES AUX DONNEES $g(\phi, u)$



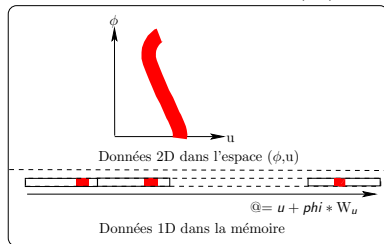
```
for (xn, yn) in Volume do
  for phi = 0 to phi_max - 1 do
    // Calcul des coordonnées
    u(phi, xn, yn) = ...
    // Interpolation linéaire
    g_interp = (1 - epsilon_u) * g(phi, u_e) +
              epsilon_u * g(phi, u_e + 1)
    // Accumulation
    f*(xn, yn) += g_interp
  end for
end for
```

Rétroprojection 2D par bloc : accès aux données localisés

CALCUL DES COORDONNEES

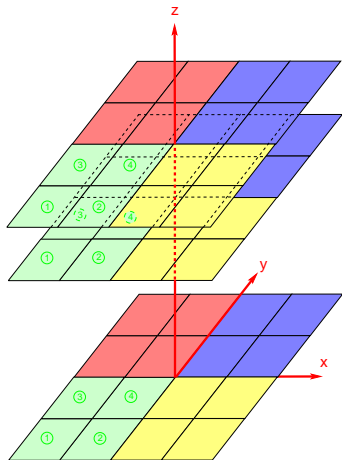


ACCES AUX DONNEES $g(\phi, u)$



```
for (Bx, By) in Volume do
  for phi = 0 to phi_max - 1 do
    for (xn, yn) in Bloc do
      // Calcul des coordonnées
      u(phi, xn, yn) = ...
      // Interpolation linéaire
      g_interp = (1 - epsilon_u) *
        g(phi, u_e) + epsilon_u * g(phi, u_e + 1)
      // Accumulation
      f*(xn, yn) += g_interp
    end for
  end for
end for
```

Découpage en threads de la rétroprojection 3D



(a) Calcul séquentiel sur un PE

- Boucle sur z
- Boucle sur ϕ

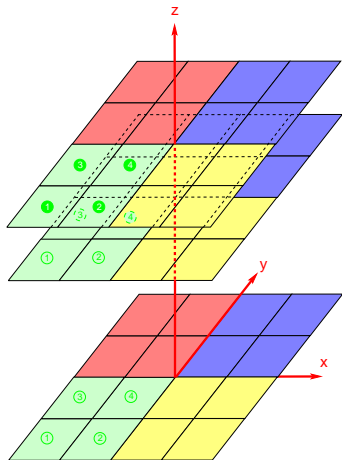
(b) Calcul parallèle sur un proc. SIMT

- Boucle sur (x,y)

(c) Calcul parallèle sur une carte

- Boucle sur les blocs (B_x, B_y, B_z)

Découpage en threads de la rétroprojection 3D



(a) Calcul séquentiel sur un PE

- Boucle sur z
- Boucle sur ϕ

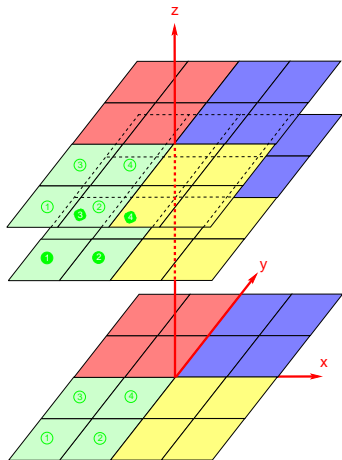
(b) Calcul parallèle sur un proc. SIMT

- Boucle sur (x,y)

(c) Calcul parallèle sur une carte

- Boucle sur les blocs (B_x, B_y, B_z)

Découpage en threads de la rétroprojection 3D



(a) Calcul séquentiel sur un PE

- Boucle sur z
- Boucle sur ϕ

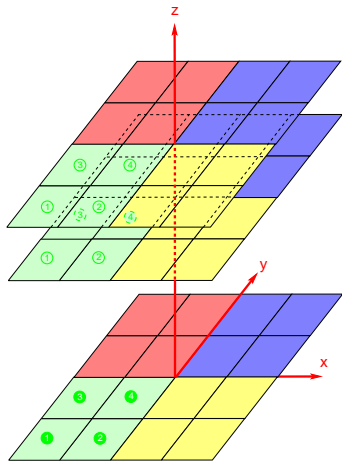
(b) Calcul parallèle sur un proc. SIMT

- Boucle sur (x,y)

(c) Calcul parallèle sur une carte

- Boucle sur les blocs (B_x, B_y, B_z)

Découpage en threads de la rétroprojection 3D



(a) Calcul séquentiel sur un PE

- Boucle sur z
- Boucle sur ϕ

(b) Calcul parallèle sur un proc. SIMT

- Boucle sur (x,y)

(c) Calcul parallèle sur une carte

- Boucle sur les blocs (B_x, B_y, B_z)

Pseudo-code parallèle de la rétroprojection 3D

```
for (Bx.id,By.id,Bz.id)=(0,0) to (NBx - 1, NBy - 1, NBz - 1) c
  for (Tx.id,Ty.id)=(0,0) to (NTx-1,NTy-1) b
    xn=Bx.id · NTx+Tx.id a
    yn=By.id · NTy+Ty.id

    for zn = znfirst to znstop
      for phi = phistart to phistop

        un(xn, yn, phi) = ...
        Ax,y,φ = ...
        vn(xn, yn, zn, phi) = Ax,y,φ · zn
        @texun,phi = un + phi · Nun
        B(xn, yn, zn)+ = tex2D(@texun,phi, vn)
        zn + +
      end

      // Accumulation du Bloc BBx.id,By.id
      // de la mémoire shared dans la mémoire globale
    end
  end
end
```

Calcul incrémental de vn

```
for (Bx.id,By.id,Bz.id)=(0,0) to (NBx - 1, NBy - 1, NBz - 1) c
  for (Tx.id,Ty.id)=(0,0) to (NTx-1,NTy-1) b
    xn=Bx.id · NTx+Tx.id a
    yn=By.id · NTy+Ty.id

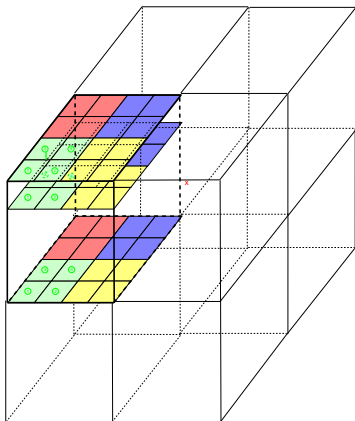
    for phi = phistart to phistop

      un(xn, yn, phi) = ...
      Ax,y,φ = ...
      vn(xn, yn, zn, phi) = Ax,y,φ · zn0
      @texun,phi = un + phi · Nun

      for zn = znfirst to znstop
        B(xn, yn, zn)+ = tex2D(@texun,phi, vn)
        vn+ = Ax,y,φ
        zn + +
      end

      // Accumulation du Bloc BBx.id,By.id
      // de la mémoire shared dans la mémoire globale
    end
  end
end
```


Découpage en threads et en grilles de la rétroprojection 3D



(a) Calcul séquentiel sur un PE

- Boucle sur z
- Boucle sur ϕ

(b) Calcul parallèle sur un proc. SIMT

- Boucle sur (x,y)

(c) Calcul parallèle sur une carte

- Boucle sur les blocs (B_x, B_y, B_z)

(d) Calcul parallèle sur PC multi-cartes

- Boucle sur les grilles (G_x, G_y)
- Boucle sur la grille G_ϕ

Pseudo-code parallèle de la rétroprojection 3D

```

for (Gx.id,Gy.id)=(0,0) to (NGx - 1, NGy - 1)
for Gφ.id=0 to NGφ - 1
//Copie des Plans Pphi_start..Pphi_stop de la mémoire CPU dans la mémoire globale du GPU
for (Bx.id,By.id,Bz.id)=(0,0) to (NBx - 1, NBy - 1, NBz - 1)
for (Tx.id,Ty.id)=(0,0) to (NTx-1,NTy-1)
xn=Bx.id · NTx+Tx.id
yn=By.id · NTy+Ty.id

for phi = phistart to phistop

un(xn, yn, phi) = ...
Ax,y,φ = ...
vn(xn, yn, zn, phi) = Ax,y,φ · zn0
@texun,phi = un + phi · Nun

for zn = znfirst to znstop
B(xn, yn, zn)+ = tex2D(@texun,phi, vn)
vn+ = Ax,y,φ
zn + +
end

// Accumulation du Bloc BBx.id,By.id,Bz.id
// de la mémoire shared dans la mémoire globale
end
end
end
//Accumulation de la grille G(Gx.id,Gy.id) de la mémoire globale dans la mémoire du PC hôte
end
    
```

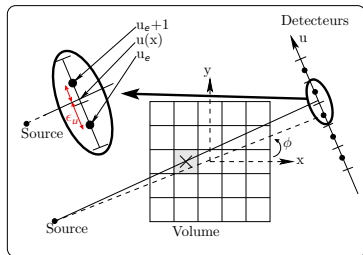
Compromis pour la taille des boucles

<i>Boucles</i>	<i>Commentaires</i>
$(B_x, B_y) / (T_x, T_y)$	$(T_x, T_y) \nearrow$ parallélisation sur un proc. SIMT $(T_x, T_y) \nearrow$ accès local via le cache 2D (8Ko) $(T_x, T_y) \searrow$ taille du cache 2D (8Ko) \rightarrow dimension u de la tuile $(T_x, T_y) \searrow$ stockage du bloc dans la mémoire shared (16Ko)
$(G_x, G_y) / (B_x, B_y)$	$(B_x, B_y) \searrow$ stockage de la grille dans la mémoire globale (4Go) $(B_x, B_y) \nearrow$ minimisation des transferts mémoire CPU \leftarrow \rightarrow GPU
$G\phi / \phi$	$\phi \nearrow$ augmentation de la taille du thread $\phi \nearrow$ accumulation dans la mémoire shared (pas dans la mémoire CPU) $\phi \searrow$ stockage des plans dans la mémoire globale (4Go) $\phi \searrow$ minimisation des transferts mémoire CPU \leftarrow \rightarrow GPU
B_z / z_n	$z_n \searrow$ stockage du bloc dans la mémoire shared (16Ko) $z_n \nearrow$ minimisation transfert mem. shared \leftarrow \rightarrow mem. globale $z_n \nearrow$ bénéfice du calcul incrémental
T_x, T_y B_x, B_y	$T_x \cdot T_y = 256$ threads $B_x \cdot B_y =$ multiple de 32

Ordres de grandeur

Volume	512^3	1024^3	2048^3
zn		16	
Tx, Ty		16·16	
Bx, By	32·32	32·32	32·32
Bz	16	32	64
Gx, Gy	1·1	2·2	4·4
un bloc	16 Ko	16 Ko	16 Ko
une grille	256 Mo	512 Mo	1 Go
ϕ	512	512	128
$G\phi$	1	2	16
Tuile	16·16=1 Ko		
Plan	$512 \cdot 256$ = 512 Ko	$1024 \cdot 512$ = 2 Mo	$2048 \cdot 1024$ = 8 Mo
Memoire globale pour projection	256 Mo	1 Go	1Go

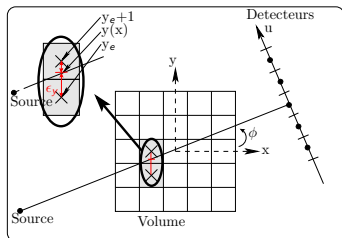
Projecteur 2D "voxel-driven"



```
for (xn, yn) in Volume do  
  for phi = 0 to phi_max - 1 do  
    // Calcul des coordonnées  
    u(phi, xn, yn) = ...  
    // Calcul des coefficients  
    d'interpolation  
    C_0 = 1 - epsilon_u  
    C_1 = epsilon_u  
    // Dépôt de la dose  
    g(u_e, phi)+ = C_0 * V(xn, yn)  
    g(u_e + 1, phi)+ = C_1 * V(xn, yn)  
  end for  
end for
```

Projecteur 2D "ray-driven" (Algorithme 2D de [Joseph 82])

```
for (un, phi) in Projection do  
  for xn = 0 to xnmax - 1 do  
    // Calcul des coordonnées  
    yn(xn, un, phi) = ...  
    // Interpolation linéaire  
    finterp = (1 - εy) · f(xn, yne) +  
              εy · f(xn, yne + 1)  
    // Accumulation  
    g*(un, phi) + = finterp  
  end for  
end for
```



Débit de calcul de la paire projection/rétroprojection

Reconstruction sur une GTX 295 (240 SPs @1.2 Ghz)

- volume de 264^3 voxels (champs de vue cylindrique)
- 64 projections de 256^2 pixels

	Temps	Giga opérations par seconde (Gop/s)
Rétroprojecteur VIB	63 ms	9,9 Gop/s
Projecteur de Joseph	63 ms	15,9 Gop/s

- 1 Objectif : accélérer la reconstruction tomographique
- 2 Tomographie sur GPU : Avant et Après CUDA
- 3 Parallélisation des opérateurs de projection/rétroprojection
- 4 Conclusion et perspectives**

Conclusion : Atteindre l'équilibre entre débit de calcul et débit mémoire

① Exploiter au maximum la puissance de calcul

- Parallélisation de l'algorithme
- Réduction de la complexité arithmétique
- Utilisation de l'interpolateur *hardware* de texture

Conclusion : Atteindre l'équilibre entre débit de calcul et débit mémoire

① Exploiter au maximum la puissance de calcul

- Parallélisation de l'algorithme
 - ↳ **Découpage en threads et en kernels**
- Réduction de la complexité arithmétique
- Utilisation de l'interpolateur *hardware* de texture

Conclusion : Atteindre l'équilibre entre débit de calcul et débit mémoire

① Exploiter au maximum la puissance de calcul

- Parallélisation de l'algorithme
↳ **Découpage en threads et en kernels**
- Réduction de la complexité arithmétique
↳ **Calcul incrémental des coordonnées**
- Utilisation de l'interpolateur *hardware* de texture

Conclusion : Atteindre l'équilibre entre débit de calcul et débit mémoire

① Exploiter au maximum la puissance de calcul

- Parallélisation de l'algorithme
↳ **Découpage en threads et en kernels**
- Réduction de la complexité arithmétique
↳ **Calcul incrémental des coordonnées**
- Utilisation de l'interpolateur *hardware* de texture
↳ **Interpolation bi-linéaire**

② Favoriser la localité spatio-temporelle des données

- Lecture des données via le cache 2D de texture
 - réutilisation au maximum des données
 - accès à des données 2D
- Ecriture dans la mémoire shared

Conclusion : Atteindre l'équilibre entre débit de calcul et débit mémoire

① Exploiter au maximum la puissance de calcul

- Parallélisation de l'algorithme
↳ **Découpage en threads et en kernels**
- Réduction de la complexité arithmétique
↳ **Calcul incrémental des coordonnées**
- Utilisation de l'interpolateur *hardware* de texture
↳ **Interpolation bi-linéaire**

② Favoriser la localité spatio-temporelle des données

- Lecture des données via le cache 2D de texture
 - réutilisation au maximum des données
 - accès à des données 2D
- Ecriture dans la mémoire shared
↳ **Harmoniser l'ordre des boucles avec la structure des données**

Travaux en cours

① Algorithme itératif avec paire projection/rétroprojection GPU

Interface Matlab

② Cohérence paire projection/rétroprojection

Etude de la convergence de l'algorithme itératif pour plusieurs projecteurs couplés avec un rétroprojecteur voxel driven

③ Projecteur ray driven avec échantillonnage régulier

Parallélisation sur GPU à l'aide des textures 3D

Perspectives

Parallélisation sur serveur 4 Tesla (S1070)

- Décomposition de l'Espace Image (DEI) : volume distribué et projections centralisées
- Décomposition de l'Espace de Projection (DEP) : volume centralisé et projections distribuées

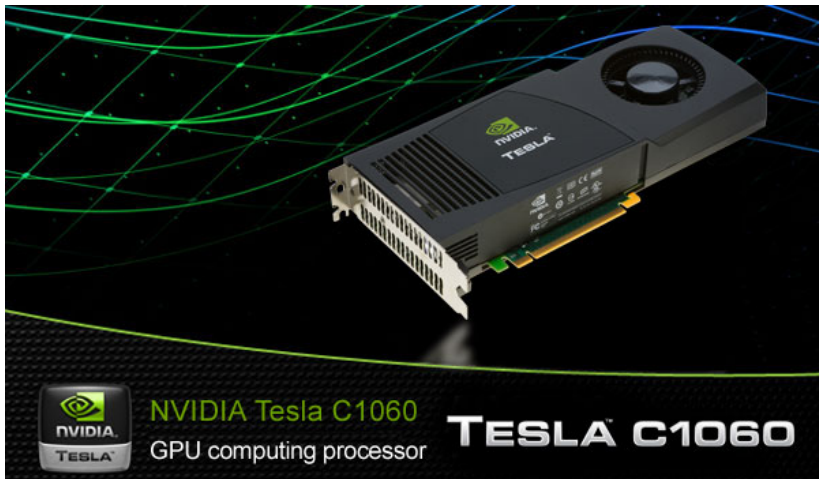
Données réelles du CEA

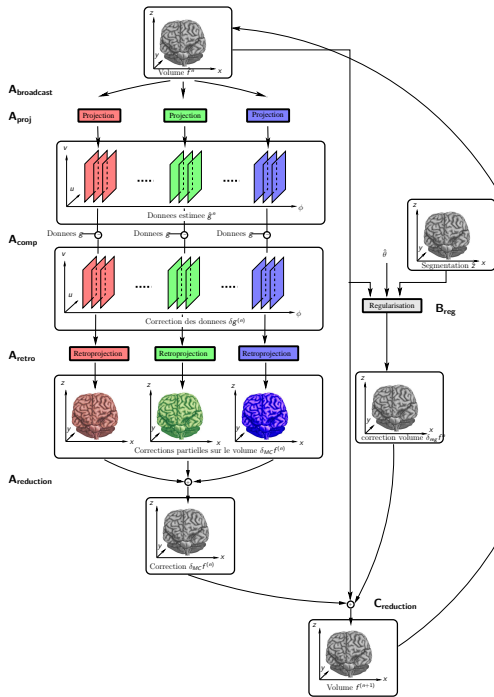
Objectif : reconstruction de volumes 2048^3

Etape de segmentation sur GPU

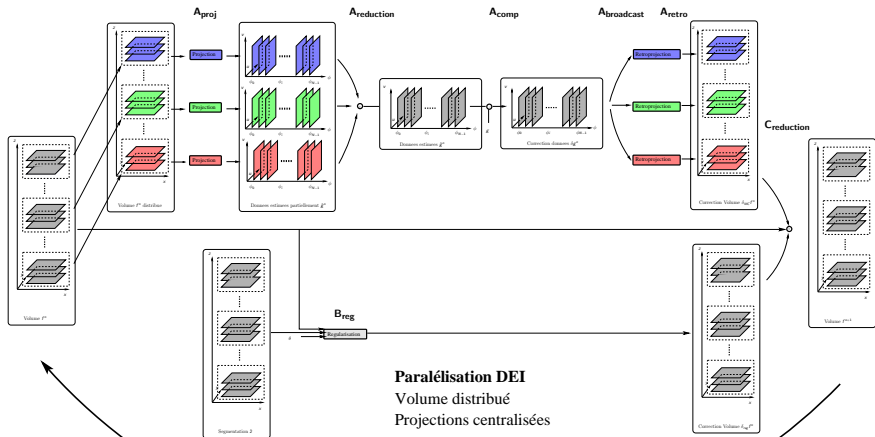
Utile ou pas utile ?

Merci de votre attention !





Parallélisation DEP
 Volume centralisé
 Projections centralisées



Références



N. Gac, S.Mancini, M.Desvignes & D.Houzet.

High Speed 3D Tomography on CPU, GPU and FPGA.

EURASIP Journal on Embedded systems, vol. Special issue : Design and Architectures for Signal Image Processing, 2008.



Peter M. Joseph.

An Improved Algorithm for Reprojecting Rays through Pixel Images.

vol. 1, no. 3, pages 192–196, Nov. 1982.



Wilfried Kirschenmann.

Parallélisation d'un solveur de neutronique sur GPU.

In 2ième Journée Thème Émergeant GPGPU (GDR ASR), December 2008.



M. Knaup & M. Kachelrieß.

GPU-Based Parallel-Beam and Cone-Beam GPU Forward- and Backprojection using CUDA.

In MIC, 2008.



Fang Xu & Klaus Mueller.

Real-time 3D computed tomographic reconstruction using commodity graphics hardware.

Physics in Medicine and Biology, vol. 52, no. 12, pages 3405–3419, 2007.