

MULTIPLICATION EXACTE MATRICE CREUSE – VECTEUR SUR ARCHITECTURES GPU ET CPU MULTICŒUR

Brice BOYER¹ Jean-Guillaume DUMAS¹
Pascal GIORGI²

¹LJK, Université de Grenoble, France

²LIRMM, Université Montpellier 2, France

JNCF '10

3 mai 2010

Motivations

- Multiplication matrice creuse – vecteur (SPMV) : la base des algorithmes “boîte noire”.

Motivations

- Multiplication matrice creuse – vecteur (SPMV) : la base des algorithmes “boîte noire”.
- GPU et CPU multicœurs se généralisent sur les ordinateurs personnels.

Motivations

- Multiplication matrice creuse – vecteur (SPMV) : la base des algorithmes “boîte noire”.
- GPU et CPU multicœurs se généralisent sur les ordinateurs personnels.
- Les GPU offrent une puissance de calcul à bas coût et de nouvelles manières de traiter les données.

Idées principales

- Utiliser des bibliothèques numériques (cf FFLAS/FFPACK) ;

Idées principales

- Utiliser des bibliothèques numériques (cf FFLAS/FFPACK) ;
- Découper A en sous-matrices, écrire $A = A_1 + \dots + A_k$, tirant avantage d'être sur des corps finis ;

Idées principales

- Utiliser des bibliothèques numériques (cf FFLAS/FFPACK) ;
- Découper A en sous-matrices, écrire $A = A_1 + \dots + A_k$, tirant avantage d'être sur des corps finis ;
- Créer automatiquement des formats hybrides efficaces selon l'architecture utilisée.

Mult. exacte
mat. creuse –
vect. dense –
sur GPU et
CPU
multicœur

**B. Boyer,
J-G Dumas &
P. Giorgi**

1 Introduction

Plan

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

Application

Conclusion

Plan

- 1 Introduction
- 2 Formats de stockage, bibliothèques et techniques existantes

Stockages classiques

Représentation des corps finis

Utilisation des bibliothèques numériques

Plan

- 1 Introduction
- 2 Formats de stockage, bibliothèques et techniques existantes
 - Stockages classiques
 - Représentation des corps finis
 - Utilisation des bibliothèques numériques
- 3 Parallélisation
 - CPU
 - GPU

Plan

- 1 Introduction
- 2 Formats de stockage, bibliothèques et techniques existantes
 - Stockages classiques
 - Représentation des corps finis
 - Utilisation des bibliothèques numériques
- 3 Parallélisation
 - CPU
 - GPU
- 4 Améliorer les formats de stockage
 - Juste à temps
 - Prendre en compte les ± 1
 - Formats Hybrides

Plan

- 1 Introduction
- 2 Formats de stockage, bibliothèques et techniques existantes
 - Stockages classiques
 - Représentation des corps finis
 - Utilisation des bibliothèques numériques
- 3 Parallélisation
 - CPU
 - GPU
- 4 Améliorer les formats de stockage
 - Juste à temps
 - Prendre en compte les ± 1
 - Formats Hybrides
- 5 Applications
 - Bien utiliser le GPU
 - Multivecteurs
 - Calcul du rang/algo par blocs de Wiedemann

Plan

- 1 Introduction
- 2 Formats de stockage, bibliothèques et techniques existantes
 - Stockages classiques
 - Représentation des corps finis
 - Utilisation des bibliothèques numériques
- 3 Parallélisation
 - CPU
 - GPU
- 4 Améliorer les formats de stockage
 - Juste à temps
 - Prendre en compte les ± 1
 - Formats Hybrides
- 5 Applications
 - Bien utiliser le GPU
 - Multivecteurs
 - Calcul du rang/algo par blocs de Wiedemann
- 6 Conclusion

- 1 Introduction
- 2 Formats de stockage, bibliothèques et techniques existantes
 - Stockages classiques
 - Représentation des corps finis
 - Utilisation des bibliothèques numériques
- 3 Parallélisation
 - CPU
 - GPU
- 4 Améliorer les formats de stockage
 - Juste à temps
 - Prendre en compte les ± 1
 - Formats Hybrides
- 5 Applications
 - Bien utiliser le GPU
 - Multivecteurs
 - Calcul du rang/algo par blocs de Wiedemann
- 6 Conclusion

Introduction

- But : une opération SPMV efficace sur de grandes matrices réutilisées abondamment.

Introduction

- **But** : une opération SPMV efficace sur de grandes matrices réutilisées abondamment.
- **Moyens** : utiliser au mieux les architectures CPU multicœur et GPU :

Introduction

- **But** : une opération SPMV efficace sur de grandes matrices réutilisées abondamment.
- **Moyens** : utiliser au mieux les architectures CPU multicœur et GPU :
 - sur CPU, le standard OpenMP ;
 - sur GPU, ATI Stream, NVIDIA Cuda, OpenCL.

Introduction

- **But** : une opération SPMV efficace sur de grandes matrices réutilisées abondamment.
- **Moyens** : utiliser au mieux les architectures CPU multicœur et GPU :
 - sur CPU, le standard **OpenMP** ;
 - sur GPU, ATI Stream, **NVIDIA Cuda**, OpenCL.

Introduction

- **But** : une opération SPMV efficace sur de grandes matrices réutilisées abondamment.
- **Moyens** : utiliser au mieux les architectures CPU multicœur et GPU :
 - sur CPU, le standard OpenMP ;
 - sur GPU, ATI Stream, NVIDIA Cuda, OpenCL.
- **Techniques** :
 - utiliser des bibliothèques numériques existantes,
 - adapter/optimiser le stockage des matrices selon les algorithmes de multiplication.

Introduction

- **But** : une opération SPMV efficace sur de grandes matrices réutilisées abondamment.
- **Moyens** : utiliser au mieux les architectures CPU multicœur et GPU :
 - sur CPU, le standard OpenMP ;
 - sur GPU, ATI Stream, NVIDIA Cuda, OpenCL.
- **Techniques** :
 - utiliser des bibliothèques numériques existantes,
 - adapter/optimiser le stockage des matrices selon les algorithmes de multiplication.

Opération à implémenter

- $\mathbf{y} = \alpha \mathbf{Ax} + \beta \mathbf{y}$.
- $\mathbf{y} = \alpha^{\top} \mathbf{Ax} + \beta \mathbf{y}$.

Introduction

- **But** : une opération SPMV efficace sur de grandes matrices réutilisées abondamment.
- **Moyens** : utiliser au mieux les architectures CPU multicœur et GPU :
 - sur CPU, le standard OpenMP ;
 - sur GPU, ATI Stream, NVIDIA Cuda, OpenCL.
- **Techniques** :
 - utiliser des bibliothèques numériques existantes,
 - adapter/optimiser le stockage des matrices selon les algorithmes de multiplication.

Opération à implémenter

- $\mathbf{y} = \alpha \mathbf{Ax} + \beta \mathbf{y}$.
- $\mathbf{y} = \alpha^{\top} \mathbf{Ax} + \beta \mathbf{y}$.

On se ramène à l'opération $\mathbf{y} = \mathbf{Ax} + \mathbf{y}$.

Plan

- 1 Introduction
- 2 **Formats de stockage, bibliothèques et techniques existantes**
 - Stockages classiques
 - Représentation des corps finis
 - Utilisation des bibliothèques numériques
- 3 Parallélisation
 - CPU
 - GPU
- 4 Améliorer les formats de stockage
 - Juste à temps
 - Prendre en compte les ± 1
 - Formats Hybrides
- 5 Applications
 - Bien utiliser le GPU
 - Multivecteurs
 - Calcul du rang/algo par blocs de Wiedemann
- 6 Conclusion

Matrices utilisées

nom	mat1916	bibd_81_3	EX5	GL7d15	mpolyout2
row	1916	3240	6545	460261	2410560
col	1916	85320	6545	171375	2086560
nbnz	195985	255960	295680	6080381	15707520
rang	1916	3240	4740	132043	1352011

TABLE: Matrices utilisées dans les tests

Stockages classiques

Exemples

- COO représente la matrice sous forme de triplets.

```
data[k] = A[rowid[k], colid[k]] ;
```

Stockages classiques

Exemples

- COO représente la matrice sous forme de **triplets**.

```
data[k] = A[rowid[k], colid[k]] ;
```

- CSR (version **compressée** de COO) :

```
start[i] ≤ k < start[i + 1] ⇒ data[k] = A[i, colid[k]] ;
```

Stockages classiques

Exemples

- COO représente la matrice sous forme de **triplets**.

$\text{data}[k] = A[\text{rowid}[k], \text{colid}[k]] ;$

- CSR (version **compressée** de COO) :

$\text{start}[i] \leq k < \text{start}[i+1] \Rightarrow \text{data}[k] = A[i, \text{colid}[k]] ;$

- ELL est plus **dense** :

$\text{data}[i, j_0] = A[i, \text{colid}[i, j_0]].$

Stockages classiques

Exemples

- COO représente la matrice sous forme de **triplets**.
 $\text{data}[k] = A[\text{rowid}[k], \text{colid}[k]] ;$
- CSR (version **compressée** de COO) :
 $\text{start}[i] \leq k < \text{start}[i+1] \Rightarrow \text{data}[k] = A[i, \text{colid}[k]] ;$
- ELL est plus **dense** :
 $\text{data}[i, j_0] = A[i, \text{colid}[i, j_0]] .$

Autres exemples

- représentation dense ;
- représentation diagonale (DIA) ;
- représentation en colonnes (CSC) ;
- ...

Représentation de $\mathbb{Z}/m\mathbb{Z}$

Retarder la réduction.

Effectuer le moins possible d'opérations f_{mod} .

Représentation de $\mathbb{Z}/m\mathbb{Z}$

Retarder la réduction.

Effectuer le moins possible d'opérations f_{mod} .

Représentations de $\mathbb{Z}/m\mathbb{Z}$

– classique : $[[0, m - 1]]$;

Représentation de $\mathbf{Z}/m\mathbf{Z}$

Retarder la réduction.

Effectuer le moins possible d'opérations f_{mod} .

Représentations de $\mathbf{Z}/m\mathbf{Z}$

- classique : $\llbracket 0, m - 1 \rrbracket$;
- centrée : $\llbracket - \lfloor \frac{m-1}{2} \rfloor, \lceil \frac{m-1}{2} \rceil \rrbracket$.

Représentation de $\mathbb{Z}/m\mathbb{Z}$

Retarder la réduction.

Effectuer le moins possible d'opérations `fmod`.

Représentations de $\mathbb{Z}/m\mathbb{Z}$

- classique : $[[0, m - 1]]$;
- centrée : $[[- \lfloor \frac{m-1}{2} \rfloor, \lceil \frac{m-1}{2} \rceil]]$.

Compromis

- coûts différents des `float` ou des `double` ;
- représentation centrée : plus d'accumulations mais réductions plus coûteuse ;
- certains GPU n'ont pas de `double` .

Exemples (CPU)

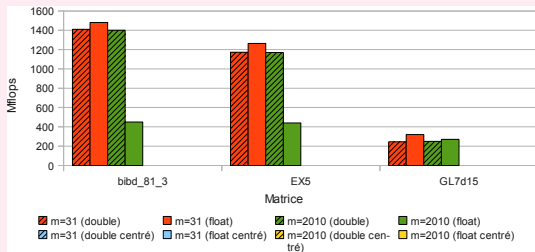


FIGURE: Matrice `bibd_81_3` sur Xeon 3.2GHz au format `ELL_R`

Exemples (CPU)

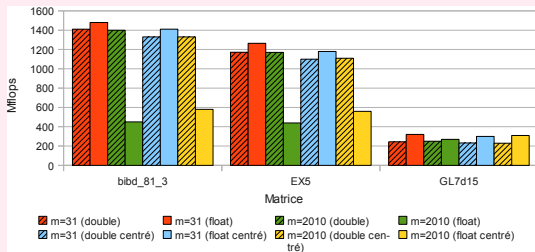


FIGURE: Matrice bibd_81_3 sur Xeon 3.2GHz au format ELL_R

Exemples (GPU)

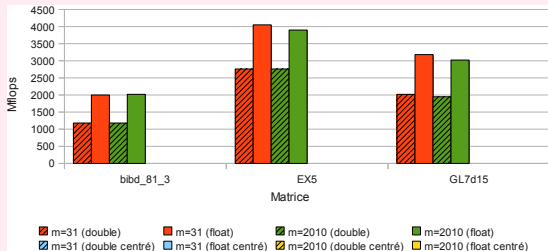


FIGURE: Matrice `bibd_81_3` sur GTX280 au format `ELL_R`

Exemples (GPU)

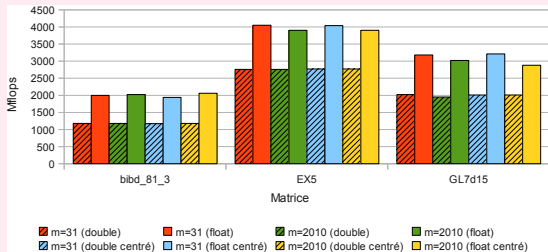


FIGURE: Matrice `bibd_81_3` sur GTX280 au format `ELL_R`

Quelques bibliothèques numériques

Peu de bibliothèques performantes

- Sur le CPU
 - OSKI
 - Sparse BLAS

Quelques bibliothèques numériques

Peu de bibliothèques performantes

- Sur le CPU
 - OSKI
 - Sparse BLAS
- Sur le GPU
 - Cusp

Mult. exacte
mat. creuse –
vect. dense –
sur GPU et
CPU
multicœur

B. Boyer,
J-G Dumas &
P. Giorgi

Utiliser les bibliothèques numériques

Difficultés par rapport au cas dense

extraction de sous-matrices de matrices creuses non triviale.

Intro

Pour
commencer

Stockages
classiques

Les corps

Librairies

Paralléliser

Améliorer le
stockage

Application

Conclusion

Utiliser les bibliothèques numériques

Difficultés par rapport au cas dense

extraction de sous-matrices de matrices creuses non triviale.

Solution

Création de sous-matrices A_i de A avec au plus r éléments non nuls par ligne.

Si $m = 2010$ et le plus grand flottant est $M = 16777215$, alors on peut effectuer $M/m^2 = 4$ accumulations avant réduction. On choisit $r = 4$.

Utiliser les bibliothèques numériques

Difficultés par rapport au cas dense

extraction de sous-matrices de matrices creuses non triviale.

Solution

Création de sous-matrices A_i de A avec au plus r éléments non nuls par ligne.

Si $m = 2010$ et le plus grand flottant est $M = 16777215$, alors on peut effectuer $M/m^2 = 4$ accumulations avant réduction. On choisit $r = 4$.

L'algorithme est alors :

```
spmv (y, A, x) {  
    //y += A x  
    foreach submatrix Ai in A do{  
        spmv_num(y, Ai, x);    //y += Ai x  
        reduce (y, m);  
    }  
}
```

Plan

- 1 Introduction
- 2 Formats de stockage, bibliothèques et techniques existantes
 - Stockages classiques
 - Représentation des corps finis
 - Utilisation des bibliothèques numériques
- 3 **Parallélisation**
 - CPU
 - GPU
- 4 Améliorer les formats de stockage
 - Juste à temps
 - Prendre en compte les ± 1
 - Formats Hybrides
- 5 Applications
 - Bien utiliser le GPU
 - Multivecteurs
 - Calcul du rang/algo par blocs de Wiedemann
- 6 Conclusion

Exemple de parallélisation avec OpenMP

```
void spmv_csr(float *data, int *colid, int *start,  
             int lig, float *x, float *y){  
  
    for(int i = 0; i < lig; ++i) {  
        for(int j = start[i]; j < start[i+1]; ++j) {  
            y[i] += data[j] * x[colid[j]];  
        }  
    }  
}
```

Exemple de parallélisation avec OpenMP

```
void spmv_csr(float *data, int *colid, int *start,  
             int lig, float *x, float *y){  
#pragma omp parallel for  
    for(int i = 0; i < lig; ++i) {  
        for(int j = start[i]; j < start[i+1]; ++j) {  
            y[i] += data[j] * x[colid[j]];  
        }  
    }  
}
```

OpenMP : performances

Les processeurs utilisés sont un Intel Xeon 8 cœurs à 3.2GHz et un Core2Duo à 3GHz.

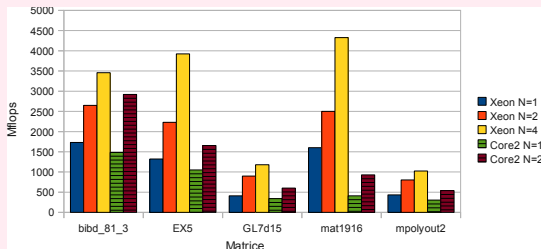


FIGURE: Performances sur un format CSR avec N cœurs utilisés

Parallélisation avec Cuda

```
#define BLOCK 256
__global__ void
spmv_csr_ker(float *data, int *colid, int *start,
             int lig, float *x, float *y){
    int i = blockDim.x * blockIdx.x + threadIdx.x;
    if(i < lig) {
        for(int j = start[i]; j < start[i+1]; ++j) {
            y[i] += data[j] * x[colid[j]];
        }
    }
}

void
spmv_csr_dev(float *data, int *colid, int *start,
             int lig, float *x, float *y){
    int grid = ceil(lig, BLOCK);
    spmv_csr_ker<<<grid, BLOCK>>>(data, colid, start, lig, x, y);
}
```

Parallélisation avec Cuda

```
#define BLOCK 256
__global__ void
spmv_csr_ker(float *data, int *colid, int *start,
             int lig, float *x, float *y){
    int i = blockDim.x * blockIdx.x + threadIdx.x;
    if(i < lig) {
        for(int j = start[i]; j < start[i+1]; ++j) {
            y[i] += data[j] * x[colid[j]];
        }
    }
}

void
spmv_csr_dev(float *data, int *colid, int *start,
             int lig, float *x, float *y){
    int grid = ceil(lig, BLOCK);
    spmv_csr_ker<<<grid, BLOCK>>>(data, colid, start, lig, x, y);
}
```

Exemple de parallélisation avec OpenMP

```
void spmv_csr(float *data, int *colid, int *start,  
             int lig, float *x, float *y){  
    #pragma omp parallel for  
    for(int i = 0; i < lig; ++i) {  
        for(int j = start[i]; j < start[i+1]; ++j) {  
            y[i] += data[j] * x[colid[j]];  
        }  
    }  
}
```

Plan

- 1 Introduction
- 2 Formats de stockage, bibliothèques et techniques existantes
 - Stockages classiques
 - Représentation des corps finis
 - Utilisation des bibliothèques numériques
- 3 Parallélisation
 - CPU
 - GPU
- 4 Améliorer les formats de stockage
 - Juste à temps
 - Prendre en compte les ± 1
 - Formats Hybrides
- 5 Applications
 - Bien utiliser le GPU
 - Multivecteurs
 - Calcul du rang/algo par blocs de Wiedemann
- 6 Conclusion

Précompiler la matrice

Idée

Consiste à compiler un fichier qui contient l'opération $y += Ax$.

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

JIT

±1

Formats Hybrides

Application

Conclusion

Précompiler la matrice

Idée

Consiste à compiler un fichier qui contient l'opération $y += Ax$.

Exemple

Si $A = \begin{pmatrix} 2 & 1 \\ 0 & 3 \end{pmatrix}$, on créerait ($m = 27$) :

```
void spmv_jit (float * y, const float * x) {  
    y[0] += 2*x[0] ;  
    y[0] += x[1] ;  
    y[0] = fmod(y[0], 27) ;  
    y[1] += 3*x[1] ;  
    y[1] = fmod(y[1], 27) ;  
}
```

Précompiler la matrice

Idée

Consiste à compiler un fichier qui contient l'opération $y += Ax$.

Exemple

Si $A = \begin{pmatrix} 2 & 1 \\ 0 & 3 \end{pmatrix}$, on créerait ($m = 27$) :

```
void spmv_jit (float * y, const float * x) {  
    y[0] += 2*x[0] ;  
    y[0] += x[1] ;  
    y[0] = fmod(y[0], 27) ;  
    y[1] += 3*x[1] ;  
    y[1] = fmod(y[1], 27) ;  
}
```

Optimisations possibles

Problèmes

Les grosses matrices ne compilent pas.

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

JIT

±1

Formats Hybrides

Application

Conclusion

Problèmes

Les grosses matrices ne compilent pas.

Solution

Découper en sous-matrices que `gcc` peut compiler.

Problèmes

Les grosses matrices ne compilent pas.

Solution

Découper en sous-matrices que `gcc` peut compiler.

Avec `nvcc`, pas de bonnes solutions pour les grosses matrices.

Problèmes

Les grosses matrices ne compilent pas.

Solution

Découper en sous-matrices que `gcc` peut compiler.

Avec `nvcc`, pas de bonnes solutions pour les grosses matrices.

Performances

bibd_81_3	
compiler	61s
SPMV	620Mflops

Prendre en compte les ± 1

Remarque.

Nombreuses applications donnant des matrices contenant beaucoup de 1 et/ou -1 .

Pour les petits m , probablement beaucoup de ± 1 .

Prendre en compte les ± 1

Remarque.

Nombreuses applications donnant des matrices contenant beaucoup de 1 et/ou -1 .

Pour les petits m , probablement beaucoup de ± 1 .

Une matrice creuse est codée

```
template<class Field>
struct SparseMatrix{
    SparseMulMatrix<Field> Mult;
    SparseAddMatrix<Field> Plus;
    SparseAddMatrix<Field> Moin;
};
```

Prendre en compte les ± 1

Remarque.

Nombreuses applications donnant des matrices contenant beaucoup de 1 et/ou -1 .

Pour les petits m , probablement beaucoup de ± 1 .

Une matrice creuse est codée

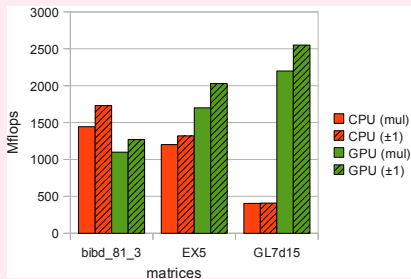
```
template<class Field>
struct SparseMatrix{
    SparseMulMatrix<Field> Mult ;
    SparseAddMatrix<Field> Plus ;
    SparseAddMatrix<Field> Moins ;
};
```

Bénéfices

- pas besoin de stocker les 1 ou les -1 ;
- sauve des multiplications ;
- retarde la réduction.

Prendre en compte les ± 1

Comparaison des performances en fonction de la prise en compte des ± 1 .



Différences entre les formats de stockage

Mult. exacte
mat. creuse –
vect. dense
sur GPU et
CPU
multicœur

**B. Boyer,
J-G Dumas &
P. Giorgi**

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

JIT

±1

Formats Hybrides

Application

Conclusion

Différences entre les formats de stockage

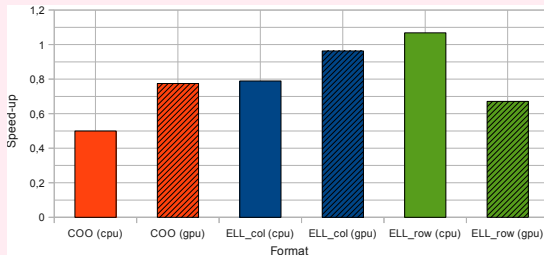


FIGURE: Rapport des performances sur un Xeon 3.2GHz et une GTX280 de divers formats par rapport au format CSR (matrice bibd_81_3).

Format Hybrides

Idée

Découper la matrice A selon des matrices $A_0 + A_1$ encodées différemment.

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

JIT

± 1

Formats Hybrides

Application

Conclusion

Format Hybrides

Idée

Découper la matrice A selon des matrices $A_0 + A_1$ encodées différemment. Par exemple ELL + CSR.

Format Hybrides

Idée

Découper la matrice A selon des matrices $A_0 + A_1$ encodées différemment. Par exemple ELL + CSR.

Exemple



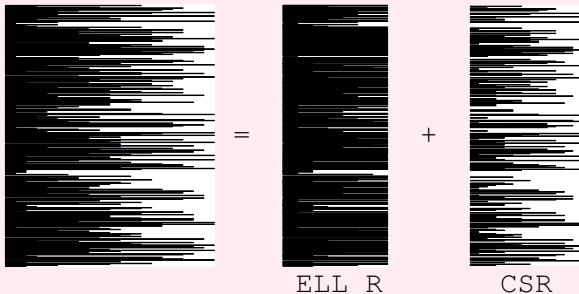
(diagramme de remplissage des lignes)

Format Hybrides

Idée

Découper la matrice A selon des matrices $A_0 + A_1$ encodées différemment. Par exemple ELL + CSR.

Exemple



(diagramme de remplissage des lignes)

Format Hybrides

Algorithme

```
choose_format (SparseMatrix A) {  
    SparseAddMatrix A_p = extrait_pun(A);  
    bool opt_p = essaie_optim(A_p);  
    SparseAddMatrix A_m = extrait_mun(A);  
    bool opt_m = essaie_optim(A_m);  
    SparseMulMatrix A_g = extrait(A, opt_p, opt_m);  
    optim(A_g);  
    retourne (A_p, A_m, A_g);  
}
```

Format Hybrides

Algorithme

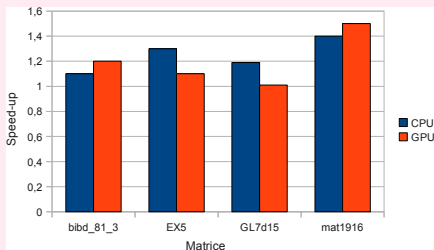
```
choose_format (SparseMatrix A) {  
    SparseAddMatrix A_p = extrait_pun (A);  
    bool opt_p = essaie_optim (A_p);  
    SparseAddMatrix A_m = extrait_mun (A);  
    bool opt_m = essaie_optim (A_m);  
    SparseMulMatrix A_g = extrait (A, opt_p, opt_m);  
    optim (A_g);  
    retourne (A_p, A_m, A_g);  
}
```

optim :

- choisit un format principal (ELL, ELL_R, CSR).
- choisit un format secondaire (si besoin) (CSR, COO_S, COO).

Format Hybrides

Accélération sur diverses matrices



(par rapport au format CSR)

Plan

- 1 Introduction
- 2 Formats de stockage, bibliothèques et techniques existantes
 - Stockages classiques
 - Représentation des corps finis
 - Utilisation des bibliothèques numériques
- 3 Parallélisation
 - CPU
 - GPU
- 4 Améliorer les formats de stockage
 - Juste à temps
 - Prendre en compte les ± 1
 - Formats Hybrides
- 5 Applications
 - Bien utiliser le GPU
 - Multivecteurs
 - Calcul du rang/algo par blocs de Wiedemann
- 6 Conclusion

Problèmes liés à l'implémentation

Comparer les deux approches :

```
spmv_n(y, A, x, n) {  
    A_d = copy_on_gpu(A);  
    y_d = copy_on_gpu(y);  
    x_d = copy_on_gpu(x);  
    for (i=0; i<n; ++i) {  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        x_d = y_d;  
    }  
    y = copy_on_cpu(y_d);  
}
```

Problèmes liés à l'implémentation

Comparer les deux approches :

```
spmv_n(y, A, x, n) {  
    A_d = copy_on_gpu(A);  
    y_d = copy_on_gpu(y);  
    x_d = copy_on_gpu(x);  
    for (i=0; i<n; ++i) {  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        x_d = y_d;  
    }  
    y = copy_on_cpu(y_d);  
}
```

Problèmes liés à l'implémentation

Comparer les deux approches :

```
spmv_n(y, A, x, n) {  
    A_d = copy_on_gpu(A);  
    y_d = copy_on_gpu(y);  
    x_d = copy_on_gpu(x);  
    for (i=0; i<n; ++i) {  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        x_d = y_d;  
    }  
    y = copy_on_cpu(y_d);  
}
```

Problèmes liés à l'implémentation

Comparer les deux approches :

```
spmv_n(y, A, x, n) {  
    A_d = copy_on_gpu(A);  
    y_d = copy_on_gpu(y);  
    x_d = copy_on_gpu(x);  
    for (i=0; i<n; ++i) {  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        x_d = y_d;  
    }  
    y = copy_on_cpu(y_d);  
}
```

Problèmes liés à l'implémentation

Comparer les deux approches :

```
spmv_n(y, A, x, n) {  
    A_d = copy_on_gpu(A);  
    y_d = copy_on_gpu(y);  
    x_d = copy_on_gpu(x);  
    for (i=0; i<n; ++i) {  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        x_d = y_d;  
    }  
    y = copy_on_cpu(y_d);  
}
```

Problèmes liés à l'implémentation

Comparer les deux approches :

```
spmv_n(y, A, x, n) {  
    A_d = copy_on_gpu(A);  
    y_d = copy_on_gpu(y);  
    x_d = copy_on_gpu(x);  
    for (i=0; i<n; ++i) {  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        x_d = y_d;  
    }  
    y = copy_on_cpu(y_d);  
}  
  
n_spmv(A, n) {  
    randomize(y);  
    randomize(x);  
    A_d = copy_on_gpu(A);  
    for (i=0; i<n; ++i) {  
        y_d = copy_on_gpu(y);  
        x_d = copy_on_gpu(x);  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        y = copy_on_cpu(y_d);  
    }  
}
```

Problèmes liés à l'implémentation

Comparer les deux approches :

```
spmv_n(y, A, x, n) {  
    A_d = copy_on_gpu(A);  
    y_d = copy_on_gpu(y);  
    x_d = copy_on_gpu(x);  
    for (i=0; i<n; ++i) {  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        x_d = y_d;  
    }  
    y = copy_on_cpu(y_d);  
}  
  
n_spmv(A, n) {  
    randomize(y);  
    randomize(x);  
    A_d = copy_on_gpu(A);  
    for (i=0; i<n; ++i) {  
        y_d = copy_on_gpu(y);  
        x_d = copy_on_gpu(x);  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        y = copy_on_cpu(y_d);  
    }  
}
```

Problèmes liés à l'implémentation

Comparer les deux approches :

```
spmv_n(y, A, x, n) {  
    A_d = copy_on_gpu(A);  
    y_d = copy_on_gpu(y);  
    x_d = copy_on_gpu(x);  
    for (i=0; i<n; ++i) {  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        x_d = y_d;  
    }  
    y = copy_on_cpu(y_d);  
}  
  
n_spmv(A, n) {  
    randomize(y);  
    randomize(x);  
    A_d = copy_on_gpu(A);  
    for (i=0; i<n; ++i) {  
        y_d = copy_on_gpu(y);  
        x_d = copy_on_gpu(x);  
        spmv_dev(y_d, A_d, x_d); // spmv on GPU  
        y = copy_on_cpu(y_d);  
    }  
}
```

Comparaison

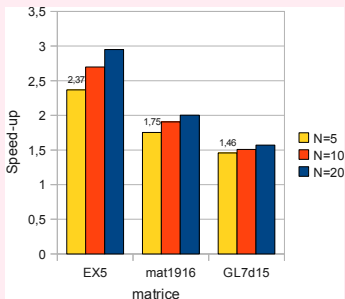


FIGURE: Accélération de $y \leftarrow A^N x$ par rapport à $N = 1$ pour $N = 5, 10, 20$ sur GTX280.

Mult. exacte
mat. creuse –
vect. dense
sur GPU et
CPU
multicœur

B. Boyer,
J-G Dumas &
P. Giorgi

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

Application

Bien utiliser le GPU

Multivecteurs

Rang/Wiedemann

Conclusion

Algorithmes par bloc

Multivecteurs

Mult. exacte
mat. creuse –
vect. dense
sur GPU et
CPU
multicœur

B. Boyer,
J-G Dumas &
P. Giorgi

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

Application

Bien utiliser le GPU

Multivecteurs

Rang/Wiedemann

Conclusion

Algorithmes par bloc

Multivecteurs

Pour des algorithmes par blocs ;

Algorithmes par bloc

Multivecteurs

Pour des algorithmes par blocs ;
implémentation particulière.

Mult. exacte
mat. creuse –
vect. dense
sur GPU et
CPU
multicœur

B. Boyer,
J-G Dumas &
P. Giorgi

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

Application

Bien utiliser le GPU

Multivecteurs

Rang/Wiedemann

Conclusion

Algorithmes par bloc

Multivecteurs

Pour des algorithmes par blocs ;
implémentation particulière.

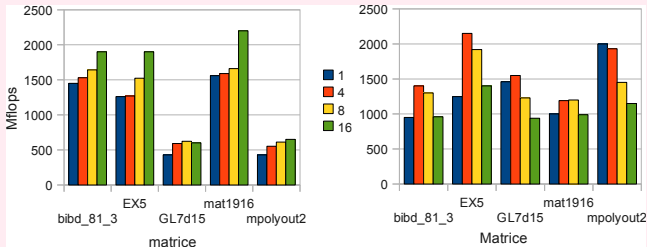


FIGURE: Performances de SPMV sur des multivecteurs de taille 1,4,8 et 16 (CPU à gauche/GPU à droite).

Algorithmme

3 phases :

1. Calcul de la suite $S_i = {}^T Y A^i Y$;

Algorithmme

Mult. exacte
mat. creuse –
vect. dense
sur GPU et
CPU
multicœur

B. Boyer,
J-G Dumas &
P. Giorgi

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

Application

Bien utiliser le GPU

Multivecteurs

Rang/Wiedemann

Conclusion

3 phases :

1. Calcul de la suite $S_i = {}^T Y A^i Y$;
2. Calcul (σ -base) du polynôme minimal P_A^x de $\sum_i S_i x^i$;

3 phases :

1. Calcul de la suite $S_i = {}^T Y A^i Y$;
2. Calcul (σ -base) du polynôme minimal P_A^x de $\sum_i S_i x^i$;
3. Interpolation du déterminant dans
 $\text{rg}(A) = \deg(\det P_A^x) - \text{codeg}(\det P_A^x)$.

Parallélisation

1. $S_i = {}^T Y A^i Y$: on utilise SPMV.

Parallélisation

1. $S_i = {}^T Y A^i Y$: on utilise SPMV.
2. σ -base : réduction à la multiplication de matrices polynomiales
calcul en parallèle des DFT et des multiplications point à point.

Parallélisation

1. $S_i = {}^T Y A^i Y$: on utilise SPMV.
2. σ -base : réduction à la multiplication de matrices polynomiales
calcul en parallèle des DFT et des multiplications point à point.
3. interpolation : évaluations de $\det P_A^x$ en parallèle.

Résultats

Mult. exacte
mat. creuse –
vect. dense
sur GPU et
CPU
multicœur

B. Boyer,
J-G Dumas &
P. Giorgi

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

Application

Bien utiliser le GPU
Multivecteurs
Rang/Wiedemann

Conclusion

Matrice	mat1916		bibd_81_3		EX5	
nb cœurs	1	8	1	8	1	8
Seq-LB	15.09	3.08	47.73	12.41	84.21	20.22
Seq-SPMV	5.02	0.91	41.28	7.56	49.66	7.36
σ -base	9.02	1.64	18.45	3.63	37.45	8.39
Interpolation	0.37	0.29	1.07	0.82	2.29	1.75
Total-LB	24.48	5.01	67.25	16.86	123.95	30.36
Total-SPMV	14.41	2.84	60.80	12.01	89.40	17.50

TABLE: Rang mod 65521 avec OpenMP sur un Xeon E5345, 8 × 2.33GHz

Plan

- 1 Introduction
- 2 Formats de stockage, bibliothèques et techniques existantes
 - Stockages classiques
 - Représentation des corps finis
 - Utilisation des bibliothèques numériques
- 3 Parallélisation
 - CPU
 - GPU
- 4 Améliorer les formats de stockage
 - Juste à temps
 - Prendre en compte les ± 1
 - Formats Hybrides
- 5 Applications
 - Bien utiliser le GPU
 - Multivecteurs
 - Calcul du rang/algo par blocs de Wiedemann
- 6 Conclusion

Mult. exacte
mat. creuse –
vect. dense –
sur GPU et
CPU
multicœur

B. Boyer,
J-G Dumas &
P. Giorgi

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

Application

Conclusion

Conclusion

- Nous obtenons de très bonnes performances avec OpenMP et Cuda ;
- la création de formats hybrides dépendant du type de parallélisme et de la répartition des éléments non nuls par ligne permet de meilleures performances ;
- un algorithme heuristique permet de choisir efficacement des formats rapides ;
- la partie CPU de cette librairie s'est intégrée facilement dans LINBOX.

Conclusion/Perspectives

Pour continuer...

– SPMV

- implémentation sur des corps finis généraux, sur $\mathbf{Z}/2\mathbf{Z}$;
- parallélisation sur multi-GPU ;
- méthodes hybrides CPUs/GPUs ;
- utilisation de METIS pour partitionner/réordonner la matrice ;
- utilisation de Sparse BLAS.
- intégrer la librairie dans LINBOX ;

Conclusion/Perspectives

Pour continuer...

– SPMV

- implémentation sur des corps finis généraux, sur $\mathbf{Z}/2\mathbf{Z}$;
- parallélisation sur multi-GPU ;
- méthodes hybrides CPUs/GPUs ;
- utilisation de METIS pour partitionner/réordonner la matrice ;
- utilisation de Sparse BLAS.
- intégrer la librairie dans LINBOX ;

– calcul du rang

- Utilisation des GPU pour SPMV et les σ -bases ;
- Se rapprocher des performances optimales sur multi-cœurs ;
- Redessiner l'algorithme des σ -bases en vue d'un meilleur parallélisme.
- Étendre l'algorithme aux corps finis généraux.

Mult. exacte
mat. creuse –
vect. dense –
sur GPU et
CPU
multicœur

B. Boyer,
J-G Dumas &
P. Giorgi

Intro

Pour
commencer

Paralléliser

Améliorer le
stockage

Application

Conclusion

Merci !

MULTIPLICATION EXACTE MATRICE CREUSE – VECTEUR SUR ARCHITECTURES GPU ET CPU MULTICŒUR

Brice BOYER¹ Jean-Guillaume DUMAS¹
Pascal GIORGI²

¹LJK, Université de Grenoble, France

²LIRMM, Université Montpellier 2, France

JNCF '10

3 mai 2010