

Corrigé Examen Novembre 2011 - Architectures Avancées

3H – Tous documents autorisés

Parties indépendantes

OPTIMISATION DE BOUCLES

On utilise le processeur superscalaire défini dans l'annexe 1, ou sa version scalaire

Soit la boucle suivante écrite en assembleur, qui travaille sur des tableaux de « floats » X[256] et Y[256]. Au démarrage R1 contient X[0] et R2 contient l'adresse de Y[0]. R4 contient la valeur 256. F0 contient la valeur 0.5.

```
Boucle :   ADDI R4,R0, -2
           LF F1, 0(R1)
           LF F2, 8(R1)
           FADD F1,F1,F2
           FMUL F1,F1,F0
           SF F1, 4(R2)
           ADDI R1,R1,4
           ADDI R2,R2,4
           ADDI R4,R4,-1
           BNE R4, Boucle
```

Question 1) Donner le programme C équivalent.

```
float x[256], y[256] a=0.5;

main ()
{
  int i ;
  for (i=1;i<255;i++)
    y[i] = (x[i-1]+ x[i+1])*a;
```

Question 2) Donner l'exécution cycle par cycle de la boucle optimisée avec un processeur scalaire. Quel est, en nombre de cycles, le temps d'exécution par itération de la boucle originale ?

Boucle	LF F1, 0(R1)
2	LF F2, 8(R1)
3	ADDI R1,R1,4
4	FADD F1,F1,F2
5	ADDI R2,R2,4
6	ADDI R4,R4,-1
7	FMUL F1,F1,F0
8	
9	
10	SF F1, 0(R2)
11	BNE R4, Boucle

11 cycles

Question 3) Donner l'exécution cycle par cycle de la boucle optimisée en plaçant les instructions dans les différents pipelines E0, E1, FA et FM. Quel est, en nombre de cycles, le temps d'exécution par itération de la boucle originale ?

	E0	E1	FA	FM
Boucle	LD F1, 0(R1)	LD F2, (R2)		
2	ADDI R1,R1,4	ADDI R2,R2,4		
3		ADDI R4,R4,-1	FADD F1,F1,F2	
4				
5				
6				FMUL F1,F1,F0
7				
8				
9	SD F1, 0(R2)	BNE R4,Boucle		
10				

9 cycles

Caches (1^{ère} partie)

Soit le programme C suivant

float X[512], Y[512]

for (i=0 ; i<510 ; i++)

Y[i] = X[i] + X[i+2];

On utilise un cache donnée de 8Ko, avec des lignes de 32 octets. On suppose que les deux tableaux sont rangés l'un derrière l'autre (&Y[0] = &X[511] +4) et que &X[0] = 0. Le cache est à allocation d'écriture (il y a des défauts de cache en écriture) et à réécriture (write back).

Question 4) Quel est le nombre total de défauts de cache dans les deux cas suivants

- cache à correspondance directe ?
- cache associatif deux voies ?

8 floats par ligne

8Ko et lignes de 32 octets : = 2⁸ lignes soit 256 lignes

&X[0] = 0

&Y[0] = 0+512*4 = 2048 = 800H

En CD, X[0] va dans la ligne 0 et Y[0] dans la ligne 128. Pas de conflit

En associatif 2 voies, pas de conflit

i=0 charge X[0] à X[7] dans ligne 0. 1 défaut

i=6 charge de X[8] à X[15] dans ligne 1 : 1 défaut

ensuite 1 défaut toutes les 8 itérations

Il y a donc 1 défaut toutes les 8 itérations en lecture et 1 défaut toutes les 8 itérations en écriture

Soit 1 défaut toutes les 4 itérations et un total de 128 défauts, tant en CD qu'en associativité 2 voies

Caches (2^{ème} partie)

On suppose des caches à écriture allouée. On ne considère que les défauts de démarrage (il n'y a pas de réutilisation des données). Les lignes de cache comprennent p floats

On considère le produit de matrices $Z[N][N] = X[N][N] * Y[N][N]$ avec des données de type « float ».

Question 5) Quel est le nombre total de défauts de cache en fonction de N dans les deux cas suivants (les algorithmes sont rappelés en annexe)

- produit ikj
- produit ijk après transposition

Produit ikj

```
for (i=0;i<N;i++)
    for (k=0;k<N;k++)
    {
        S=Y[i][k]; // défaut en lecture
        Z[i][0]=0.0; //défaut en écriture
        for (j=0; j<N; j++)
            Z[i][j]+=S*X[k][j]; // 2 défauts en lecture/p
                                     itérations
                                     // pas de défaut en écriture
    }
```

Total : $2 N^2 + 2 N^3/p = 2N^2 (N/p+1)$ soit $2N^3/p$ si $N/p \gg 1$

```
for (i=0;i<N;i++)
    for (j=0;j<N;j++)
        XT[j][i]=X[i][j]; //1/p défaut lecture et 1 défaut (écriture)
for (i=0;i<N;i++)
    for (j=0;j<N;j++)
    {
        S=0.0;
        for (k=0; k<N; k++)
            S+=Y[i][k]*XT[j][k]; // 2/p défauts en lecture
        Z[i][j]=S; // 1 défaut écriture
    }
```

Transposition : $(1+1/p)N^2$

PM : $2/p * N^3 + N^2$

Total $2 N^3/p + N^2(2+1/p) = 2N^2(N/p + 1+1/2p)$

Soit $2N^3/p$ si $N/p \gg 1+1/2p$

Instructions SIMD

Soit l'instruction SIMD SSE2 SHUFPS qui travaille sur des doubles (64 bits) dont l'intrinsic est

```
__m128d _mm_shuffle_pd(__m128d a, __m128d b, unsigned int imm8)
```

Soit le define

```
# define SHUFPS (a,b) __m128d _mm_shuffle_pd(__m128d a, __m128d b, unsigned int imm8)
```

Ah correspond aux 64 bits de poids fort de A, Ab correspond aux 64 bits de poids faible de A, Bh correspond aux 64 bits de poids fort de B, Bb correspond aux 64 bits de poids faible de B.

Alors

SHUFPS(A,B,0) = (Bb|Ab)

SHUFPS(A,B,1) = (Bb|Ah)

SHUFPS(A,B,2) = (Bh|Ab)

SHUFPS(A,B,3) = (Bh|Ah)

Soient les déclarations suivantes :

```
_m128d **XS ; //mots de 128 bits contenant des doubles
```

```
double **X ;
```

```
XS=X ;
```

Question 6) Donner le résultat des intrinsics suivants

- SHUFPS (XS [i][j], XS [i][j-1],1)

- SHUFPD (XS [i][j+1], XS [i][j],1)

Un mot de 128 bits $Xs [i][j]$ correspond à deux mots de 64 bits $X[i][2j]$ et $X[i][2j+1]$ en supposant des adresses croissantes de la gauche vers la droite.

$XS[i][j] = X[i][2j] \mid X[i][2j+1]$

$XS[i][j-1] = X[i][2j-2] \mid X[i][2j-1]$

$XS[i][j+1] = X[i][2j+2] \mid X[i][2j+3]$

$SHUFPD (XS[i][j], XS[i][j-1], 1) = X[i][2j-1] \mid X[i][2j]$ // Décalage à gauche d'une position des doubles

$SHUFPD (XS[i][j+1], XS[i][j], 1) = X[i][2j+1] \mid X[i][2j+2]$ // Décalage à droite d'une position des doubles

-

Question 7) Donner le programme C scalaire correspondant au programme C SIMD suivant :

```
_m128d **XS ; **YS, D, tmp, tmp2,tmp3, tmp4 //mots de 128 bits contenant
des doubles
double **X, **Y ;
XS=X ; YS=Y ;
```

```
D= _mm_set1_pd (0.25) ; // crée un mot de deux doubles égaux à 0.25
for (i=1 ; i<255 ; i++)
for (j=1 ; j<127 ; j++)
{
tmp = addpd (ld2d (&XS[i-1][j], ld2d(&XS[i+1][j])) ;
tmp2 = SHUFPD (XS [i][j], XS [i][j-1],1);
tmp3 = SHUFPD (XS [i][j+1],XS [i][j],1));
tmp2 = addpd (tp2,tmp3;
tmp = addpd (tmp;tmp2);
tmp=mulpd (tmp, D);
st2d (&YS[i][j], tmp)
}
```

```
for (i=1 ; i<255 ; i++)
for (j=1 ; j<255 ; j++)
Y[i][j]=0.25 * (X[i-1][j]+X[i+1][j]+X[i][j-1] +X[i][j+1]);
```

PROGRAMMATION OPENMP

Soit le programme C ci-dessous (Jacobi).

```
#include <stdio.h>
double X[128][128], Xold[128][128];

/* version séquentielle */

int main(void)
{
    int i,j,k, maxit, size;
    double tol, error, resid;
    tol = 0.0001;
    error = 10*tol;
    maxit=50000;
    size=128*128;

    for (i=0; i<128;i++)
        for (j=0; j<128; j++)
            X[i][j] =(double) (i*i+j*j+i+j);
    k=0;
```

```

while (k<maxit & error>tol)
{
    error = 0.0;
    {
        for (i = 0; i<128;i++)
            for (j=0; j<128 ; j++)
                Xold[i][j]=X[i][j];
    }
    for (i=1; i<128-1; i++)
        for (j=1 ; j<128-1; j++) {
            X[i][j] = 0.25*(Xold[i-1][j]+Xold[i+1][j]+Xold[i][j-1]
                +Xold[i][j+1]);
            resid=X[i][j] - Xold[i][j];
            error+=resid*resid;}
        k++;
        error=sqrt(error)/((double)size;}
}
}

```

Question 8) Donner une version OpenMP du programme.

```

#include <stdio.h>
#include <omp.h>

double X[128][128], Xold[128][128];
/* version OpenMP */
int main(void)
{
    int i,j,k, maxit, size;
    double tol, error, resid;
    tol = 0.0001;
    error = 10*tol;
    maxit=50000;
    size=128*128;
    for (i=0; i<128;i++)
        for (j=0; j<128; j++)
            X[i][j] =(double) (i*i+j*j+i+j);
    k=0;
    omp_set_num_threads(2);
    while (k<maxit & error>tol)
    {
        error = 0.0;
        #pragma omp parallel
        {
            for (i = 0; i<128;i++)
                for (j=0; j<128 ; j++)
                    Xold[i][j]=X[i][j];
        }
        #pragma omp parallel private (resid) reduction (+:error)
        {
            for (i=1; i<128-1; i++)
                for (j=1 ; j<128-1; j++)
                    {
                        X[i][j] = 0.25 * (Xold[i-1][j]+Xold[i+1][j]+Xold[i][j-1]
                            +Xold[i][j+1]);
                        resid=X[i][j] - Xold[i][j];
                        error+=resid*resid;
                    }
            k++;
            error=sqrt(error)/((double)size;
        }
    }
}
}

```

LOI D'AMDAHL

Question 9) Quel est la fraction maximale de la partie séquentielle d'un programme pour obtenir une efficacité parallèle de 50% avec 16 processeurs ?

$$EP = Acc/n = \frac{1}{16 * (S + \frac{1-S}{16})} = 0,5$$

$$16S+1-S=15S+1 = 2$$

$$15S = 1$$

$$S = 1/15 = 6,67\%$$

PREDICTION DE BRANCHEMENT

Un programme a cinq branchements conditionnels notés Branch 1 à Branch 5. Le comportement des branchements est le suivant pour une exécution du programme (P pour branchement pris et N pour branchement non pris).

Branch 1 : P-P-P-P-P

Branch 2 : N-N-N

Branch 3 : P-N-P-N-P-N-P-N

Branch 4 : P-P-P-N-N-N

Branch 5 : P-P-P-N-P-P-P-N-P

On suppose que le comportement de chaque branchement reste le même pour chaque exécution du programme. Les prédicteurs 1 bit et 2 bits sont toujours initialisés dans le même état avant chaque séquence d'exécution d'un des branchements.

Question 10) Donner le nombre de bonnes prédictions pour chaque branchement avec les prédicteurs suivants :

a) Toujours pris

B1 : 5/5 – B2 : 0/3 – B3 : 4/8 - B4 : 3/6 – B5 : 7/9

b) Toujours non pris

B1 : 0/5 – B2 : 3/3 – B3 : 4/8 - B4 : 3/6 – B5 : 2/9

c) Prédicteur 1 bit, initialisé à prédit pris.

B1 : 5/5 – B2 : 2/3 – B3 : 1/8 - B4 : 5/6 – B5 : 5/9

d) Prédicteur 2 bits, initialisé à prédit faiblement pris.

B1 : 5/5 – B2 : 2/3 – B3 : 1/8 - B4 : 5/6 – B5 : 7/9

Annexe 1

Soit un processeur superscalaire à ordonnancement statique qui a les caractéristiques suivantes :

- les instructions sont de longueur fixe (32 bits)

- Il a 32 registres entiers (dont R0=0) de 32 bits et 32 registres flottants (de F0 à F31) de 32 bits.

- Il peut lire et exécuter 4 instructions par cycle.

- L'unité entière contient deux pipelines d'exécution entière sur 32 bits, soit deux additionneurs, deux décaleurs. Tous les bypass possibles sont implantés.

- L'unité flottante contient un pipeline flottant pour l'addition et un pipeline flottant pour la multiplication. - L'unité Load/Store peut exécuter jusqu'à deux chargements par cycle, mais ne peut effectuer qu'un load et un store simultanément. Elle ne peut effectuer qu'un seul store par cycle.

- Il dispose d'un mécanisme de prédiction de branchement qui permet de "brancher" en 1 cycle si la prédiction est correcte. Les sauts et branchements ne sont pas retardés.

La Table 1 donne les instructions disponibles et le pipeline qu'elles utilisent : E0 et E1 sont les deux pipelines entiers, FA est le pipeline flottant de l'addition et FM le pipeline flottant de la multiplication. Les instructions peuvent être exécutées simultanément si elles utilisent chacune un pipeline séparé.

L'addition et la multiplication flottante sont pipelinées. La division flottante n'est pas pipelinée (une division ne peut commencer que lorsque la division précédente est terminée).

L'ordonnancement est statique.

JEU D'INSTRUCTIONS (extrait)

LF	LF Fi, dép.(Ra)	2	E0 ou E1	$F_i \leftarrow M(Ra + \text{dépl.16 bits avec ES})$
SF	SF Fi, dép.(Ra)	0	E0	$F_i \rightarrow M(Ra + \text{dépl.16 bits avec ES})$
ADD	ADD Rd,Ra, Rb	1	E0 ou E1	$Rd \leftarrow Ra + Rb$
ADDI	ADDI Rd, Ra, IMM	1	E0 ou E1	$Rd \leftarrow Ra + \text{IMM-16 bits avec ES}$
SUB	SUB Rd,Ra, Rb	1	E0 ou E1	$Rd \leftarrow Ra - Rb$
FADD	FADD Fd, Fa, Fb	3	FA	$Fd \leftarrow Fa + Fb$
FMUL	FMUL Fd, Fa, Fb	3	FM	$Fd \leftarrow Fa \times Fb$
BEQ	BEQ Ri, dépl	1	E1	si $R_i=0$ alors $CP \leftarrow NCP + \text{depl}$
BNE	BNE Ri, dépl	1	E1	si $R_i \neq 0$ alors $CP \leftarrow NCP + \text{depl}$

Table 1 : instructions disponibles

ANNEXE 2 : Instructions SIMD IA-32 utilisables

#define ld2d(a)	_mm_load_pd(&a)	chargement aligné de deux doubles
#define st2d(a, b)	_mm_store_pd(&a, b)	rangement aligné de deux doubles
#define addpd(a,b)	_mm_add_pd(a,b)	Addition SIMD de 2 doubles
#define mulpd(a,b)	_mm_mul_pd(a,b)	Multiplication SIMD de 2 doubles
#define setd (a)	_mm_set1_pd (a)	Mot constitué de deux doubles a