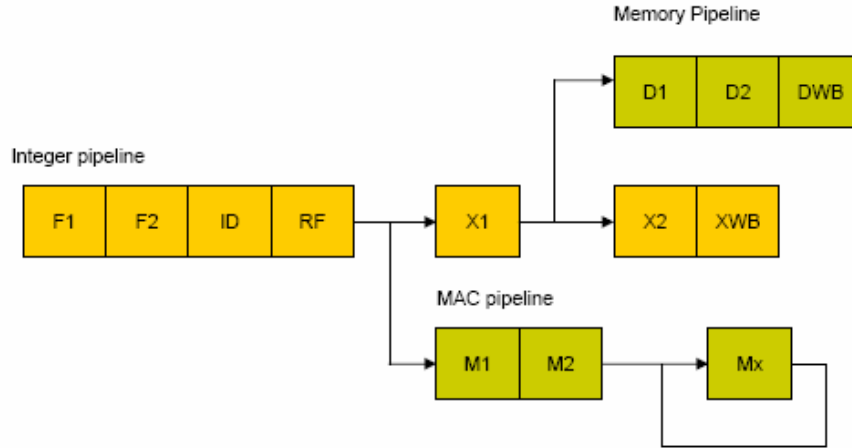


## Examen Décembre 06 - Architectures Avancées

3H – Tous documents autorisés

### PIPELINES



**Figure 1 : pipelines entiers du processeur Xscale**

La figure 1 donne les pipelines entiers du processeur Xscale pour les instructions entières, les instructions mémoires et l'instruction MAC (multiplication accumulation)

La signification des différents étages est

- F1 et F2 : phases d'acquisition (lecture) de l'instruction
- ID : décodage de l'instruction
- RF : lecture des registres
- X1 : calcul adresse mémoire, adresse de branchement, début exécution (EX1).
- X2 : fin exécution (EX2)
- XWB : écriture du résultat dans un registre
- D1 et D2 : phases accès cache données
- DWB : écrire de la donnée dans un registre
- M1 et M2 : phase de multiplication
- Mx : accumulateur  $\leftarrow$  accumulateur + résultat multiplication

#### Question 1 :

a) Donner les latences entre instruction producteur et instruction consommateur en nombre de cycles en supposant que tous les bypass nécessaires soient implémentés. (NB : la valeur n signifie que deux instructions dépendantes peuvent se suivre aux cycles i et i+n)

	Instruction producteur	Instruction consommateur	Latence scalaire
1	UAL Ri, -, -	UAL -, Ri, -	
2	UAL Ri, -, -	ST Ri, ()	
3	LW Ri, ()	UAL -, Ri, -	

b) Quelle est la pénalité pour un branchement mal prédit ?

### OPTIMISATION DE BOUCLES

On utilise le processeur superscalaire défini dans l'annexe 1  
Soit la boucle suivante :

```
float X[128], Y[128], S, tmp ;  
for (i = 0 ; i<128 ; i++){  
    tmp = X[i]*Y[i]  
    if (tmp >0.0) S+=tmp ; }
```

On supposera que l'adresse de X[0] est initialement dans R1, que l'adresse de Y[0] est initialement dans R2. R3 contient initialement le nombre d'itérations de la boucle. On utilisera F0 pour contenir la valeur 0.0.

### Question 2

Quel est en nombre de cycles, le temps d'exécution par itération de la boucle originale sans et avec utilisation des instructions SIMD. ?

### Question 3

Donner par itération de la boucle initiale

- le nombre de cycles sans instructions SIMD avec un déroulage d'ordre 4
- le nombre de cycles avec instructions SIMD avec un déroulage d'ordre 4

### Question 4

- a) Pour la version sans SIMD avec déroulage de boucle d'ordre 4, donner le temps d'exécution total du programme pour les 128 itérations (prendre en compte les instructions avant l'entrée dans la boucle déroulée, et les instructions en sortie de la boucle déroulée)
- b) Pour la version SIMD sans déroulage de boucle, on suppose qu'il n'existe pas d'instruction permettant de faire la somme des quatre floats contenus dans un registre SIMD. Ecrire un programme assembleur qui permette de faire cette somme et de ranger le résultat en mémoire (variable S). Donner le nombre de cycles pour exécuter ce programme.

### **CACHES**

La figure 2 donne le temps, exprimé en cycles par pixel, de la copie d'une image de N x N pixels avec un octet ou un float par pixel sur un processeur Pentium 4. Dans les deux cas, la copie est « vectorisée »

**Question 5 : Expliquer les courbes obtenues.**

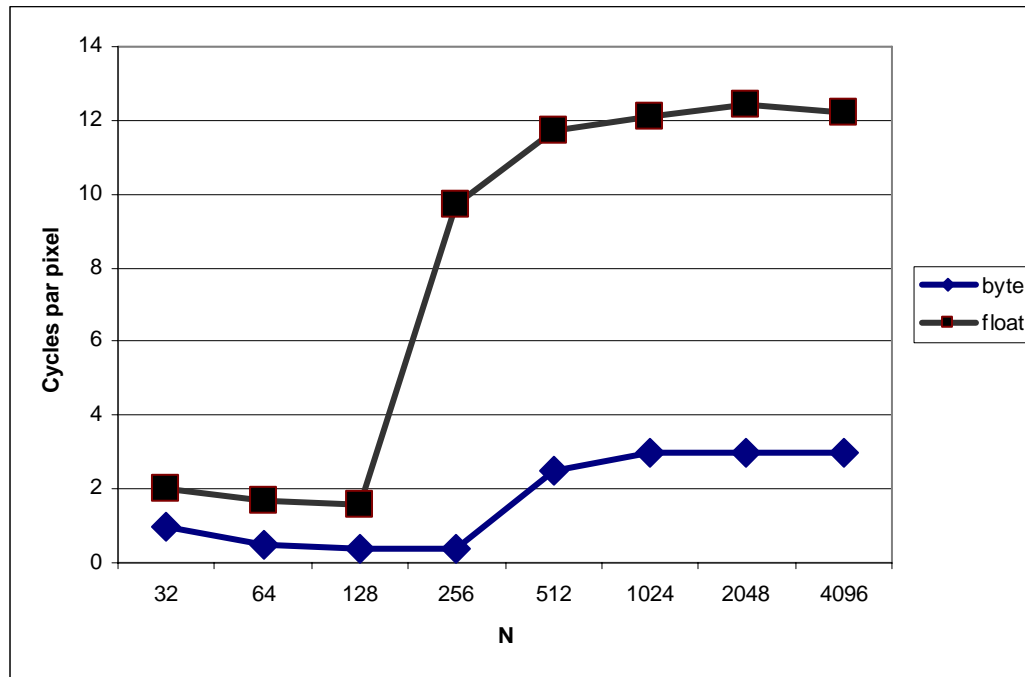


Figure 2 : Temps (en cycles par pixel) de la copie d'une image  $N \times N$  pixels en fonction de  $N$ .

### PIPELINE LOGICIEL AVEC TMS 320C62

Le code assembleur TMS320C62 ci-dessous donne l'itération du pipeline logiciel pour un programme C.

```

LOOP:
    LDW    .D1    *A4++,A2           ; load ai and ai+1
    LDW    .D2    *B4++,B2           ; load bi and bi+1
    [A0]   ADD    .L1    A6,A7,A7     ; sum0+=(ai*bi)
    [B0]   ADD    .L2    B6,B7,B7     ; sum1+=(ai+1*bi+1)
    [A1]   B      .S2    LOOP         ; branch to loop

    MPY    .M1X   A2,B2,A6           ; ai*bi
    MPYH   .M2X   A2,B2,B6           ; ai+1*bi+1
    CMPLT  .L1    0,A6,A0            ; A0 = 1 si A6 > 0
    CMPLT  .L2    0,B6,B0            ; B0 = 1 si B6 > 0
    [A1]   SUB    .S1    A1,1,A1      ; decrement loop counter

    ADD    .L1X   A7,B7,A4           ; sum=sum0+sum1
    
```

Question 6 : Donner le code C correspondant au code assembleur.

Question 7 : Quel est l'intervalle inter-itération (II) ? Justifier sa valeur.

### SIMD IA-32

Le jeu d'instruction IA-32 contient des instructions SIMD de conversion d'entiers 32 bits (int) en flottants 32 bits simple précision (float) : CVTDQ2PS .

On considère des variables 128 bits (`_mm128i`) `v0, v1, v2, v3, v4, v5, v6, v7` pour les entiers et des variables 128 bits (`_mm128`) `f0, f1, f2` et `f3` pour les flottants simple précision.

**Question 8 :** En supposant que la variable `v0` contient 16 octets (`unsigned char`) correspondant à des pixels (niveaux de gris), donner la suite des instructions SIMD IA-32 (Annexe 2) nécessaires pour convertir les 16 données 8 bits en 16 données de type `float` dans les variables `v4` à `v7`. Quelle opération faut-il alors effectuer pour avoir des données flottantes comprises entre 0.0 et 1.0 ?

### **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 (R0=0) de 32 bits et 32 registres flottants (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. Les chargements ne peuvent pas passer devant les rangements en attente.

#### JEU D'INSTRUCTIONS (extrait)

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

**Table 1 : instructions disponibles**

La Table 2 donne la latence entre une instruction source et une instruction destination, dans le cas de dépendances de données. La valeur 1 est le cas où les deux instructions peuvent se succéder normalement, d'un cycle  $i$  au cycle  $i+1$ .

Latences	<i>Source</i>	UAL	LF/SF (données)	FADD/ FMAX/ FMIN	FMUL	FDIV	FSQRT
		<i>Destination</i>					
UAL		1	2				
LF/ST (adresses)		1	3				
SF (données)		1	2	3	4	20 (NP)	20 (NP)
Opération flottante			2	3	4	20 (NP)	20 (NP)

**Table 2 : latences des instructions**

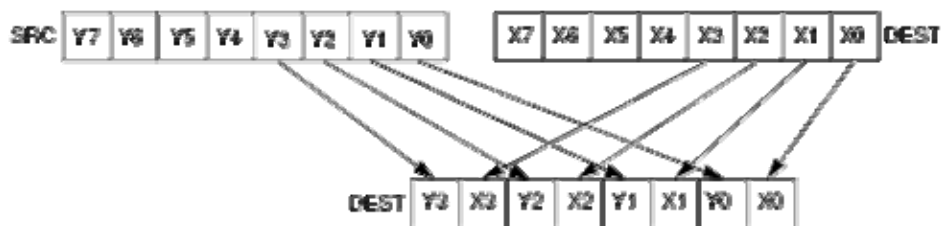
Le processeur a également 32 registres de 128 bits  $S_0$  à  $S_7$  pouvant contenir chacun 4 flottants simple précision et les instructions SIMD données dans la Table 3. Cette table donne également les latences des instructions SIMD et le pipeline utilisé dans le cas superscalaire.

PLF	PLF $S_i$ , dép.(Ra)	2	E0	Charge quatre flottants simple précision à partir de l'adresse (Ra + dépl.16 bits avec ES).
PSF	PSF $S_i$ , dép.(Ra)	2	E0	Range en mémoire à partir de l'adresse (Ra + dépl.16 bits avec ES) quatre flottants simple précision
PFADD	PFADD $S_d, S_a, S_b$	3	FA	$S_d \leftarrow S_a + S_b$ sur quatre « floats »
PFSUB	PFSUB $S_d, S_a, S_b$	3	FA	$S_d \leftarrow S_a - S_b$ sur quatre « floats »
PFMUL	PFMUL $S_d, S_a, S_b$	4	FM	$S_d \leftarrow S_a \times S_b$ sur quatre « floats »
PFMULS	PFMULS $S_d, S_a, S_b$	4	FM	$S_f \leftarrow S_a \times S_b$ (chaque élément de $S_a$ est multiplié par $S_b$ et rangé dans l'élément correspondant de $S_f$ ) sur quatre « floats »
PLB	PLB $S_i$ , dép.(Ra)	2	E0	Charge seize octets à partir de l'adresse (Ra + dépl.16 bits avec ES).
PSB	PSB $S_i$ , dép.(Ra)	2	E0	Range en mémoire à partir de l'adresse (Ra + dépl.16 bits avec ES) seize octets
PFMAX	PFMAX $S_d, S_a, S_b$	1	FA	$S_d \leftarrow S_a \max S_b$ : maximum sur 4 floats
PFMIN	PFMIN $S_d, S_a, S_b$	1	FA	$S_d \leftarrow S_a \min S_b$ : minimum sur 4 floats

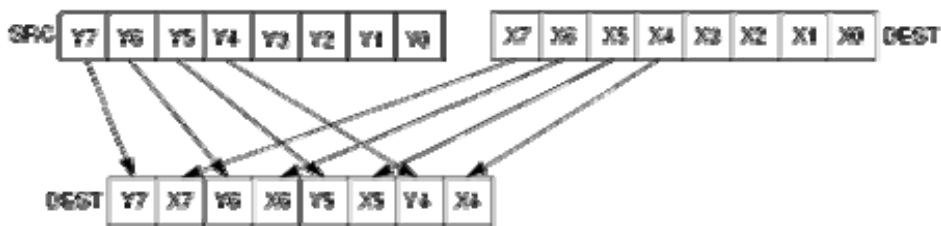
**Table 3 : Instructions SIMD**

**ANNEXE 2 : Instructions SIMD IA-32 utilisables**

CVTDQ2PS	_mm_cvtepi32_ps (a)	Convertit 4 entiers 32 bits signés en 32 bits flottants
MULPS	_mm_mul_ps(a,b)	Quatre multiplications flottantes 32 bits
DIVPS	_mm_div_ps (a,b)	Quatre divisions flottantes 32 bits
MOVDQA	_mm_load_si128 (*p)	Chargement aligné de 128 bits
MOVDQA	_mm_store_si128 (*p,a)	Rangement aligné de 128 bits
MULPS	_mm_mul_ps (a, b)	Quatre multiplications flottantes 32 bits
PACKUSWB	_mm_packs_epu16 (m1, m2)	Compacte avec saturation shorts en octets non signés
PMAXUB	_mm_max_epu8 ( a, b)	Max des octets non signés source et destination
PMINUB	_mm_min_epu8 ( a, b)	Min des octets non signés source et destination
POR	_mm_or_si128 (a, b)	Ou logique parallèle
PSRLDQ	_mm_srli_si128 (a, imm)	Décalage logique gauche de « imm » octets
PSSLDQ	_mm_slli_si128 (a, imm)	Décalage logique droite de « imm » octets
PUNPCKHBW	_mm_unpacklo_epi8 ( m1, m2)	Entrelace les octets (haut) de la destination et la source
PUNPCKHWD	_mm_unpacklo_epi16( m1, m2)	Entrelace les shorts (haut) de la destination et la source
PUNPCKLBW	_mm_unpacklo_epi8 ( m1, m2)	Entrelace les octets (bas) de la destination et la source
PUNPCKLWD	_mm_unpacklo_epi16 ( m1, m2)	Entrelace les shorts (bas) de la destination et la source
PXOR	_mm_xor_si128 (a, b)	Ou exclusif parallèle
	_mm_set_ps1 (float w)	Quatre fois le flottant 32 bits w dans un mot de 128 bits



PUNPCKLBW



PUNPCKHBW