

Le calcul flottant avec le processeur intel 64 bits

©Sovanna Tan

Novembre 2013

L'unité de calcul flottant (FPU : Floating Point Unit)

- Au début, elle était un circuit séparé, le coprocesseur mathématique.
- Puis, elle a été intégrée au microprocesseur, mais se program-mait toujours de façon distincte.
- Avec les architectures 64 bits, apparaît un nouveau jeu d'ins-tructions SSE (Streaming SIMD (Single Instruction Multiple Data) Extensions) pour des registres de taille 128 ou 256 bits selon les processeurs.

Les types de nombres flottants

- Simple précision IEEE 754-2008 : 1 bit signe, 8 bits exposant, 23 bits mantisse
- Double précision IEEE 754-2008 : 1 bit signe, 11 bits exposant, 52 bits mantisse
- Precision étendue intel 80 bits : 1 bit signe, 15 bits exposant, 64 bits mantisse
 - Ce format est différent des formats IEEE 754-2008.
 - Dans le processeur Intel 32 bits, tous les calculs étaient effectués en précision étendue puis convertis en simple ou double précision IEEE 754-2008.
- Quadruple précision IEEE 754-2008 : 1 bit signe, 15 bits exposant, 112 bits mantisse
- 256 bits (extrapolation de IEEE 754-2008) : 1 bit signe, 19 bits exposant, 236 bits mantisse

Les registres

- 16 registres généraux **xmm0**, **xmm1**,... et **xmm15**.
- Ces registres peuvent être utilisés avec des instructions
 - pour une seule valeur
 - ou pour un vecteur de valeurs.

Ce vecteur comporte alors soit 4 **floats** soit 2 **doubles**.

- Les processeurs de la série i-core ont 16 registres 256 bits **ymm0**, **ymm1**,... et **ymm15**. Les registres **xmm0**, **xmm1**,... et **xmm15** correspondent aux premiers 128 bits des registres **ymm0**, **ymm1**,... et **ymm15**. Ces derniers peuvent être vu comme des vecteurs de 8 **floats** ou 4 **doubles**.

Les instructions de déplacement de données

- Déplacement d'une seule valeur

movsd <i>dest,src</i>	déplacement d'un flottant double
movss <i>dest,src</i>	déplacement d'un flottant simple (float)

- Déplacement de paquets (plusieurs valeurs simultanément)

movapd <i>dest,src</i>	déplacement de 2 doubles alignés ¹
movaps <i>dest,src</i>	déplacement de 4 floats alignés
movdqa <i>dest,src</i>	déplacement d'un double qword aligné
movupd <i>dest,src</i>	déplacement de 2 doubles non alignés
movups <i>dest,src</i>	déplacement de 4 floats non alignés
movdqu <i>dest,src</i>	déplacement d'un double qword non aligné

- Les noms des instructions pour les registres 256 bits commencent par **v**, **vmovapd**, **vmovaps**,...

1. alignés sur 1 bloc de 16 octets

Calculs flottants

addsd <i>dest,src</i>	addition de 2 doubles
addss <i>dest,src</i>	addition de 2 floats
addpd <i>dest,src</i>	2 doubles + 2 doubles
addps <i>dest,src</i>	4 floats + 4 floats
subsd <i>dest,src</i>	soustraction de 2 doubles
subss <i>dest,src</i>	soustraction de 2 floats
subpd <i>dest,src</i>	2 doubles - 2 doubles
subps <i>dest,src</i>	4 floats - 4 floats
mulsd <i>dest,src</i>	multiplication de 2 doubles
mulss <i>dest,src</i>	multiplication de 2 floats
mulpd <i>dest,src</i>	2 doubles * 2 doubles
mulps <i>dest,src</i>	4 floats * 4 floats
divsd <i>dest,src</i>	division de 2 doubles
divss <i>dest,src</i>	division de 2 floats
divpd <i>dest,src</i>	2 doubles / 2 doubles
divps <i>dest,src</i>	4 floats / 4 floats

- *dest* doit être un registre flottant.

Les instructions de conversion

cvtss2sd <i>dest,src</i>	convertit un float en double
cvtps2pd <i>dest,src</i>	convertit un paquet de 2 floats en un paquet de 2 doubles
cvtisd2ss <i>dest,src</i>	convertit un double en float
cvtisd2ps <i>dest,src</i>	convertit un paquet de 2 doubles en un paquet de 2 floats
cvtss2si <i>dest,src</i>	convertit un float en dword ou qword
cvtisd2si <i>dest,src</i>	convertit un double en dword ou qword
cvtisi2ss <i>dest,src</i>	convertit un entier en float
cvtisi2sd <i>dest,src</i>	convertit un entier en double

- Lorsque l'opérande source est dans la mémoire, il faut préciser la taille de la donnée à lire.

Les instructions de comparaison

comisd <i>op1,op2</i>	comparaison ordonnée de deux doubles
comiss <i>op1,op2</i>	comparaison ordonnée de deux floats
ucomisd <i>op1,op2</i>	comparaison non ordonnée de deux doubles
ucomiss <i>op1,op2</i>	comparaison non ordonnée de deux floats

- *op1* doit être un registre, *op2* peut être un registre ou une valeur en mémoire.
- Les instructions mettent à jour les drapeaux ZF, PF et CF.
- Les comparaisons non ordonnées lèvent une exception si l'un des opérandes vaut SNaN (Signaling Not a Number i.e. que des 1 dans l'exposant et mantisse de la forme 0x) et les comparaisons ordonnées lèvent une exception si l'un des deux opérandes vaut SNaN ou QNaN (Quiet Not a Number i.e. que des 1 dans l'exposant et mantisse de la forme 1x).

Les branchements conditionnels pour les flottants

- Après une instruction de comparaison de flottants (**comisd**, **comiss**, **ucomisd** ou **ucomiss**) entre *op1* et *op2* :

jb <i>op</i>	branchement à <i>op</i> si $op1 < op2$	CF=1
jbe <i>op</i>	branchement à <i>op</i> si $op1 \leq op2$	CF=1 ou ZF=1
ja <i>op</i>	branchement à <i>op</i> si $op1 > op2$	CF=1
jae <i>op</i>	branchement à <i>op</i> si $op1 \geq op2$	CF=1 ou ZF=1

Arrondis

roundsd <i>dest,src,mode</i>	arrondi <i>src</i> de type double dans <i>dest</i>
roundss <i>dest,src,mode</i>	arrondi <i>src</i> de type float dans <i>dest</i>
roundpd <i>dest,src,mode</i>	arrondi des 2 doubles de <i>src</i> dans <i>dest</i>
roundps <i>dest,src,mode</i>	arrondis des 4 floats de <i>src</i> dans <i>dest</i>

- *dest* doit être un registre flottant.
- *mode* est un opérande immédiat :

0	arrondi à la valeur entière paire la plus proche
1	arrondi à la valeur entière inférieure
2	arrondi à la valeur supérieure
3	arrondi à la valeur absolue inférieure (à la valeur entière la plus proche de 0)

Instructions diverses

minsd <i>dest,src</i>	minimum de 2 doubles (<i>src</i> et <i>dest</i>) dans <i>dest</i>
minss <i>dest,src</i>	minimum de 2 floats (<i>src</i> et <i>dest</i>) dans <i>dest</i>
minpd <i>dest,src</i>	2 calculs de minimum de 2 doubles dans <i>dest</i>
minps <i>dest,src</i>	4 calculs de minimum de 2 floats dans <i>dest</i>
maxsd <i>dest,src</i>	maximum de 2 doubles (<i>src</i> et <i>dest</i>) dans <i>dest</i>
maxss <i>dest,src</i>	maximum de 2 floats (<i>src</i> et <i>dest</i>) dans <i>dest</i>
maxpd <i>dest,src</i>	2 calculs de maximum de 2 doubles dans <i>dest</i>
maxps <i>dest,src</i>	4 calculs de maximum de 2 floats dans <i>dest</i>
sqrtsd <i>dest,src</i>	racine carrée de <i>src</i> de type double dans <i>dest</i>
sqrtps <i>dest,src</i>	racine carrée de <i>src</i> de type float dans <i>dest</i>
sqrtpd <i>dest,src</i>	2 racines carrées de <i>src</i> (2 doubles) dans <i>dest</i>
sqrtps <i>dest,src</i>	4 racines carrées de <i>src</i> (4 floats) dans <i>dest</i>

- *dest* doit être un registre flottant.

Exemple : calcul de la valeur d'un polynôme (1)

```

extern printf, scanf
segment .data
coeff dq 1.0, 2.0, 3.0, 4.0
x dq 2.0
promptDouble db "%lf", 0
prompt1DoubleIn db "4x3+3x2+2x+1_(2) _=%lf", 10, 0
prompt2DoubleIn db "4x3+3x2+2x+1_(%lf) _=%lf", 10, 0
promptEntree db "Entrez x: ", 0
segment .bss
y resq 1
segment .text
global asm_main, horner

asm_main:
    push    rbp
    mov     rbp, rsp
    lea     rdi, [coeff]
    movsdb xmm0, [x]
    mov     esi, 3
    call    horner
    mov     rdi, prompt1DoubleIn
    mov     rax, 1
    call    printf
    mov     rdi, promptEntree
    mov     rax, 0
    call    printf
    mov     rdi, promptDouble
    mov     rsi, y
    mov     rax, 0
    call    scanf

    lea     rdi, [coeff]
    movsdb xmm0, [y]
    mov     esi, 3
    call    horner
    mov     rdi, prompt2DoubleIn
    movsdb xmm1, xmm0
    movsdb xmm0, [y]
    mov     rax, 2
    call    printf
    xor     eax, eax
    leave
    ret

; extrait du livre de Ray Seyfarth
horner: movsdb xmm1, xmm0
; use xmm1 as x
    movsdb xmm0, [rdi+rsi*8]
; accumulator for b_k
    cmp     esi, 0
; is the degree 0?
    jz      done
    more:   sub     esi, 1
    mulsdb xmm0, xmm1
; b_k * x
    addsdb xmm0, [rdi+rsi*8]
; add p_k
    jnz     more
done:     ret

```

Exemple : calcul de la valeur d'un polynôme (2)

```
./horner
```

$$4x^3+3x^2+2x+1 \quad (2) = 49.000000$$

Entrez x : 3

$$4x^3+3x^2+2x+1 \quad (3.000000) = 142.000000$$

```
./horner
```

$$4x^3+3x^2+2x+1 \quad (2) = 49.000000$$

Entrez x : 1.5

$$4x^3+3x^2+2x+1 \quad (1.500000) = 24.250000$$

Algorithme :

- $P(x) = (...(((a_n x + a_{n-1})x + a_{n-2})x + a_{n-3})...)x + a_0$

Exemple : calcul de la somme de deux vecteurs (1)

```

extern printf
segment .data
a      dq  1.0, 2.0, 3.0, 4.0, 5.0, 6.0, 7.0, 8.0
b      dq  2.0, 0.5, 4.0, 0.2, 3.0, 1.5, 0.5, 0.4
floatFormat db "%lf_",0
promptRes db "a*b=",0
promptln db 10,0
segment .bss
c      resq 8
segment .text
global asm_main, horner
asm_main:
    push    rbp
    mov     rbp, rsp
    mov     rcx,0
    mov     rdx,8
    jmp     test1
next1:  movapd xmm0,[a+8*rcx]
        movapd xmm1,[b+8*rcx]
        addpd  xmm1,xmm0
        movapd [c+8*rcx],xmm1
        add     rcx,2
test1:  cmp     rcx,rdx
        jl     next1

    mov     rdi,a
    mov     rsi,8
    call    affichage

    mov     rdi,b
    mov     rsi,8
    call    affichage
    mov     rax,0
    leave
    ret

affichage:
    enter 0,0
    push    r12
    push    r13
    push    r14
    push    rbx
    mov     rbx,rdi
    mov     r13,rsi
    mov     r12,0
    jmp     test3
next3:  mov     rdi,floatFormat
        movsd  xmm0,[rbx+8*r12]
        mov     rax,1
        call   printf
        inc     r12
test3:  cmp     r12,r13
        jl     next3

```

Exemple : calcul de la somme de deux vecteurs (2)

```

mov rdi, promptln
mov rax, 0
call    printf

```

```

pop rbx
pop r14
pop r13
pop r12
mov rax, 0
leave
ret

```

./addvect

1.000000	2.000000	3.000000	4.000000	5.000000	6.000000	7.000000	8.000000
2.000000	0.500000	4.000000	0.200000	3.000000	1.500000	0.500000	0.400000
3.000000	2.500000	7.000000	4.200000	8.000000	7.500000	7.500000	8.400000